
Preface

For quantitative researchers working in an investment bank, the process of writing a fixed income model usually has two stages. First, a theoretical framework for yield curve dynamics is specified, using the language of mathematics (especially stochastic calculus) to ensure that the underlying model is well-specified and internally consistent. Second, in order to use the model in practice, the equations arising from the first step need to be turned into a working implementation on a computer. While specification of the theoretical model may be seen as the difficult part, in quantitative finance applications the second step is technically and intellectually often more challenging than the first. In the implementation phase, not only does one need to translate abstract ideas into computer code, one also needs to ensure that the resulting numbers being produced are meaningful to a trading desk, are stable and robust, are in line with market observations, and are produced in a timely manner. Many of these requirements are, as it turns out, extremely challenging, and not only demand a strong knowledge of actual market practices (which tend to deviate in significant ways from “textbook” theory), but also require application of a large arsenal of techniques from applied mathematics, chiefly approximation methods and numerical techniques.

While there are many good introductory books on fixed income derivatives on the market, when we hire people who have read them we find that they still require significant training before they become productive members of our quantitative research teams. For one, while existing literature covers some aspects of the first step above, advanced approaches to specifying yield curve dynamics are typically not covered in sufficient detail. More importantly, there is simply too little said in the literature about the process of getting the theory to work in the real world of trading and risk management. An important goal of our book series is to close these gaps in the literature.

As we write this in early 2010, financial markets are still reeling from a severe crisis that has, at least in part, been blamed on over-the-counter (OTC) options markets, the venue where complex derivative securities are transacted. Stricter regulation of some types of OTC derivatives currently seems all but inevitable, and many common OTC securities may in the future either be outlawed or traded only on public exchanges. In the wake of the crisis, opinion of financial engineers and bankers has hit an all-time low, with many in the public convinced that they are peddlers of toxic waste or “weapons of financial destruction”. All things considered, the present may therefore seem like an inauspicious moment to launch a series of monographs on the pricing and risk management of interest rate derivatives. We disagree, for several reasons. First, in defense of OTC derivatives we note that although they certainly can be used inappropriately to create excessive leverage and risk, many complex (or “exotic”) derivatives serve as innovative and cost-effective vehicles for bank clients to reduce their financial risk. Second, irrespective of what will ultimately transpire on the regulatory front, it has become obvious that going forward both regulators and market participants need a better grasp of the management and characterization of complex financial risk. This is perhaps particularly true for the quantitative research professionals (the “quants”, in common parlance) who recently have been taken to task by the press for the failure of their models and their inability to predict the credit crisis. While this simplistic characterization is actually quite unfair, there is no doubt that many derivatives models that worked well enough before the credit crisis are no longer adequate. Indeed, even the simple task of pricing a basic interest rate swap — possibly the simplest of all interest rate derivatives — has recently required major methodology revisions¹. If nothing else, a severe crisis serves to expose weaknesses in the foundation on which models are built, allowing one to reinforce it for future storms. In this light, we feel that the time is just about right for a comprehensive, practical, and up-to-date exposition of interest rate modeling and risk management².

The three volumes of *Interest Rate Modeling* are aimed primarily at practitioners working in the area of interest rate derivatives, but much of the material is quite general and, we believe, will also hold significant appeal to researchers working in other asset classes. Students and academics interested in financial engineering and applied work will find the material particularly useful for its description of real-life model usage and for its expansive discussion of model calibration, approximation theory, and numerical methods. In preparing the books we have drawn on nearly 30 years of combined industry experience, and much of the material has never been exposed in book form before.

¹We cover this in Chapter 6.

²We ought to note that interest rate derivatives (unlike *credit* derivatives) so far have not been directly implicated in the financial crisis.

Quantitative finance attracts students and practitioners from many different academic fields, and with varying levels of preparation in mathematics and computation. (Case in point: L.B.G.A was originally a robotics engineer and V.V.P a probabilist.) To cater to a broad audience, we have kept the exposition fairly informal; graduate students in applied fields such as engineering and physics should feel at home with the level (or lack) of rigor used in the book. We have relied on a proposition-proof format throughout, largely because this facilitates easier cross-referencing in a long text, but acknowledge that the format is occasionally more formal than the results themselves. For instance, we tend to skip over technical regularity conditions in our proofs and also frequently list approximate results in propositions without explicitly specifying the sense in which they approximate true values. Although the exposition is largely self-contained, some previous knowledge of basic option pricing principles (e.g., at the level of Hull [2006]) may be useful.

Interest Rate Modeling divides into three separate volumes. *Volume I* provides the theoretical and computational foundations for the series, emphasizing the construction of efficient grid- and simulation-based methods for contingent claims pricing. Numerical methods serve an extremely important role in the text, so we develop this topic to an advanced level suitable for professional-quality model implementations. Placing this material early in the text allows us to incorporate it into our discussion of individual models in subsequent chapters. The second part of *Volume I* is dedicated to local-stochastic volatility modeling and to the construction of vanilla models for individual swap and Libor rates. Although the focus is eventually turned toward fixed income securities, much of the material in this volume applies to a broad capital market setting and will be of interest to anybody working in the general area of asset pricing.

Volume II is dedicated to in-depth study of term structure models of interest rates. While providing a thorough analysis of classical short rate models, the primary focus of the volume is on multi-factor stochastic volatility dynamics, in the setups of both the separable HJM and Libor market models. Implementation techniques are covered in detail, as are strategies for model parameterization and calibration to market data.

The first half of *Volume III* contains a detailed study of several classes of fixed income securities, ranging from simple vanilla options to highly exotic cancelable and path-dependent trades. The analysis is done in product-specific fashion, covering, among other subjects, risk characterization, calibration strategies, and valuation methods. In its second half, *Volume III* studies the general topic of derivative portfolio risk management, with a particular emphasis on the challenging problem of computing smooth price sensitivities to market input perturbations.

Although much of the material in *Interest Rate Modeling* is focused on the technical and theoretical issues surrounding model implementation on a computer, it is impractical for us to delve into the exercise of writing actual

computer routines. Fortunately, there are several specialized books on how to write good quant code, see, e.g., Hyer [2010] and Joshi [2004]. Both of these books work with C++ which is still the most common computer language used in professional quant libraries. For those that choose to work with C++, we wholeheartedly endorse books by Scott Meyers (see, e.g., Meyers [2005]) and Andrei Alexandrescu (see, e.g., Sutter and Alexandrescu [2004]) as guides to sound and maintainable code.

During the six year process of writing this book series, we have received encouragement and constructive criticism from many people. We particularly wish to thank Peter Carr, Peter Forsyth, Alexandre Antonov, Peter Jäckel, Dominique Bang, Martin Dahlgren, Neil Oliver, Patrick Roome, Regis van Steenkiste, Natasha Bushueva and many members of the research teams at Barclays Capital and Bank of America Merrill Lynch. Natalia Kryzhanovskaya meticulously proofread our first draft, and contributed greatly to the harmonization of notation across what turned out to be a very long manuscript. All remaining errors are, of course, entirely our own. Speaking of errors: with nearly 20,000 equations, it is probable that a few typos remain, despite our best efforts to weed them out. A list of errata will be maintained on www.andersen-piterbarg-book.com where supplemental material and news will also be posted on a running basis. We greatly appreciate reporting of typos or factual errors to our web address, and will list the names of all those who contribute to error spotting in future editions of *Interest Rate Modeling*.

Lastly, we owe a great debt of gratitude to our families for their support and patience, even when our initial plans for a brief book on tips and tricks for working quants ballooned into something more ambitious that consumed many evenings and weekends over the last six years.

London, New York,
June 2004 — August 2010

*Leif B.G Andersen
Vladimir V. Piterbarg*

Table of Contents for All Volumes

VOLUME I Foundations and Vanilla Models

Part I Foundations

1	Introduction to Arbitrage Pricing Theory	3
1.1	The Setup	3
1.2	Trading Gains and Arbitrage	7
1.3	Equivalent Martingale Measures and Arbitrage	8
1.4	Derivative Security Pricing and Complete Markets	10
1.5	Girsanov's Theorem	12
1.6	Stochastic Differential Equations	14
1.7	Explicit Trading Strategies and PDEs	16
1.8	Kolmogorov's Equations and the Feynman-Kac Theorem	18
1.9	Black-Scholes and Extensions	21
1.9.1	Basics	21
1.9.2	Alternative Derivation	25
1.9.3	Extensions	27
1.9.3.1	Deterministic Parameters and Dividends	27
1.9.3.2	Stochastic Interest Rates	28
1.10	Options with Early Exercise Rights	30
1.10.1	The Markovian Case	32

1.10.2	Some General Bounds	34
1.10.3	Early Exercise Premia	36
2	Finite Difference Methods	43
2.1	1-Dimensional PDEs: Problem Formulation	43
2.2	Finite Difference Discretization	45
2.2.1	Discretization in x -Direction. Dirichlet Boundary Conditions	45
2.2.2	Other Boundary Conditions	47
2.2.3	Time-Discretization	49
2.2.4	Finite Difference Scheme	50
2.3	Stability	52
2.3.1	Matrix Methods	52
2.3.2	Von Neumann Analysis	53
2.4	Non-Equidistant Discretization	56
2.5	Smoothing and Continuity Correction	58
2.5.1	Crank-Nicolson Oscillation Remedies	58
2.5.2	Continuity Correction	58
2.5.3	Grid Shifting	59
2.6	Convection-Dominated PDEs	60
2.6.1	Upwinding	61
2.6.2	Other Techniques	62
2.7	Option Examples	63
2.7.1	Continuous Barrier Options	63
2.7.2	Discrete Barrier Options	65
2.7.3	Coupon-Paying Securities and Dividends	67
2.7.4	Securities with Early Exercise	68
2.7.5	Path-Dependent Options	69
2.7.6	Multiple Exercise Rights	70
2.8	Special Issues	72
2.8.1	Mesh Refinements for Multiple Events	72
2.8.2	Analytics at the Last Time Step	75
2.8.3	Analytics at the First Time Step	76
2.9	Multi-Dimensional PDEs: Problem Formulation	78
2.10	Two-Dimensional PDE with No Mixed Derivatives	79
2.10.1	Theta Method	80
2.10.2	The Alternating Direction Implicit (ADI) Method	81
2.10.3	Boundary Conditions and Other Issues	84
2.11	Two-Dimensional PDE with Mixed Derivatives	85
2.11.1	Orthogonalization of the PDE	85
2.11.2	Predictor-Corrector Scheme	88
2.12	PDEs of Arbitrary Order	91

3 Monte Carlo Methods	93
3.1 Fundamentals	93
3.1.1 Generation of Random Samples	95
3.1.1.1 Inverse Transform Method	96
3.1.1.2 Acceptance-Rejection Method	97
3.1.1.3 Composition	99
3.1.2 Correlated Gaussian Samples	100
3.1.2.1 Cholesky Decomposition	101
3.1.2.2 Eigenvalue Decomposition	102
3.1.3 Principal Components Analysis (PCA)	103
3.2 Generation of Sample Paths	104
3.2.1 Example: Asian Basket Options in Black-Scholes Economy	104
3.2.2 Discretization Schemes, Convergence, and Stability	106
3.2.3 The Euler Scheme	108
3.2.3.1 Linear-Drift SDEs	110
3.2.3.2 Log-Euler Scheme	110
3.2.4 The Implicit Euler Scheme	111
3.2.4.1 Implicit Diffusion Term	112
3.2.5 Predictor-Corrector Schemes	113
3.2.6 Ito-Taylor Expansions and Higher-Order Schemes	114
3.2.6.1 Ordinary Taylor Expansion of ODEs	115
3.2.6.2 Ito-Taylor Expansions	116
3.2.6.3 Milstein Second-Order Discretization Scheme	117
3.2.7 Other Second-Order Schemes	119
3.2.8 Bias vs. Monte Carlo Error	120
3.2.9 Sampling of Continuous Process Extremes	122
3.2.10 PCA and Bridge Construction of Brownian Motion Paths	126
3.2.10.1 Brownian Bridge and Quasi-Random Sequences	126
3.2.10.2 PC Construction	128
3.3 Sensitivity Computations	129
3.3.1 Finite Difference Estimates	129
3.3.1.1 Black-Scholes Delta	129
3.3.1.2 General Case	131
3.3.2 Pathwise Estimate	133
3.3.2.1 Black-Scholes Delta	133
3.3.2.2 General Case	134
3.3.2.3 Sensitivity Path Generation	136
3.3.3 Likelihood Ratio Method	136
3.3.3.1 Black-Scholes Delta	137
3.3.3.2 General Case	138
3.3.3.3 Euler Schemes	138

3.3.3.4	Some Remarks	139
3.4	Variance Reduction Techniques	140
3.4.1	Variance Reduction and Efficiency	141
3.4.2	Antithetic Variates	141
3.4.2.1	The Gaussian Case	141
3.4.2.2	General Case	143
3.4.3	Control Variates	143
3.4.3.1	Basic Idea	143
3.4.3.2	Non-Linear Controls	145
3.4.4	Importance Sampling	146
3.4.4.1	Basic Idea	146
3.4.4.2	Density Formulation	147
3.4.4.3	Importance Sampling and SDEs	149
3.4.4.4	More on SDE Path Simulation	150
3.4.4.5	Rare Event Simulation and Linearization	152
3.5	Some Notes on Bermudan Security Pricing	156
3.5.1	Basic Idea	156
3.5.2	Parametric Lower Bound Methods	157
3.5.3	Parametric Lower Bound: An Example	158
3.5.4	Regression-Based Lower Bound	159
3.5.5	Upper Bound Methods	160
3.5.6	Confidence Intervals	161
3.5.7	Other Methods	162
3.A	Appendix: Constants for Φ^{-1} Algorithm	163
4	Fundamentals of Interest Rate Modeling	165
4.1	Fixed Income Notations	165
4.1.1	Bonds and Forward Rates	165
4.1.2	Futures Rates	167
4.1.3	Annuity Factors and Par Rates	168
4.2	Fixed Income Probability Measures	169
4.2.1	Risk Neutral Measure	170
4.2.2	T -Forward Measure	172
4.2.3	Spot Measure	173
4.2.4	Terminal and Hybrid Measures	174
4.2.5	Swap Measures	175
4.3	Multi-Currency Markets	176
4.3.1	Notations and FX Forwards	176
4.3.2	Risk Neutral Measures	177
4.3.3	Other Measures	178
4.4	The HJM Analysis	179
4.4.1	Bond Price Dynamics	179
4.4.2	Forward Rate Dynamics	180
4.4.3	Short Rate Process	181
4.5	Examples of HJM Models	182

4.5.1	The Gaussian Model	182
4.5.2	Gaussian HJM Models with Markovian Short Rate	185
4.5.3	Log-Normal HJM Models	187
5	Fixed Income Instruments	189
5.1	Fixed Income Markets and Participants	189
5.2	Certificates of Deposit and Libor Rates	192
5.3	Forward Rate Agreements (FRA)	193
5.4	Eurodollar Futures	194
5.5	Fixed-for-Floating Swaps	195
5.6	Libor-in-Arrears Swaps	198
5.7	Averaging Swaps	199
5.8	Caps and Floors	199
5.9	Digital Caps and Floors	201
5.10	European Swaptions	201
5.10.1	Cash-Settled Swaptions	203
5.11	CMS Swaps, Caps and Floors	204
5.12	Bermudan Swaptions	205
5.13	Exotic Swaps and Structured Notes	206
5.13.1	Libor-Based Exotic Swaps	207
5.13.2	CMS-Based Exotic Swaps	208
5.13.3	Multi-Rate Exotic Swaps	208
5.13.4	Range Accruals	209
5.13.5	Path-Dependent Swaps	210
5.14	Callable Libor Exotics	211
5.14.1	Definitions	211
5.14.2	Pricing Callable Libor Exotics	213
5.14.3	Types of Callable Libor Exotics	214
5.14.4	Callable Snowballs	214
5.14.5	CLEs Accreting at Coupon Rate	214
5.14.6	Multi-Tranches	215
5.15	TARNs and Other Trade-Level Features	215
5.15.1	Knock-out Swaps	216
5.15.2	TARNs	216
5.15.3	Global Cap	217
5.15.4	Global Floor	217
5.15.5	Pricing and Trade Representation Challenges	217
5.16	Volatility Derivatives	218
5.16.1	Volatility Swaps	218
5.16.2	Volatility Swaps with a Shout	219
5.16.3	Min-Max Volatility Swaps	220
5.16.4	Forward Starting Options and Other Forward Volatility Contracts	220
5.A	Appendix: Day Counting Rules and Other Trivia	221
5.A.1	Libor Rate Definitions	222
5.A.2	Swap Payments	223

Part II Vanilla Models

6	Yield Curve Construction and Risk Management	227
6.1	Notations and Problem Definition	228
6.1.1	Discount Curves	228
6.1.2	Matrix Formulation	230
6.1.3	Construction Principles and Yield Curves	230
6.2	Yield Curve Fitting with N -Knot Splines	232
6.2.1	C^0 Yield Curves: Bootstrapping	232
6.2.1.1	Piecewise Linear Yields	233
6.2.1.2	Piecewise Flat Forward Rates	234
6.2.2	C^1 Yield Curves: Hermite Splines	236
6.2.3	C^2 Yield Curves: Twice Differentiable Cubic Splines	238
6.2.4	C^2 Yield Curves: Twice Differentiable Tension Splines	241
6.3	Non-Parametric Optimal Yield Curve Fitting	243
6.3.1	Norm Specification and Optimization	243
6.3.2	Choice of λ	246
6.3.3	Example	247
6.4	Managing Yield Curve Risk	248
6.4.1	Par-Point Approach	249
6.4.2	Forward Rate Approach	250
6.4.3	From Risks to Hedging: The Jacobian Approach ..	252
6.4.4	Cumulative Shifts and other Common Tricks ..	254
6.5	Various Topics in Discount Curve Construction	256
6.5.1	Curve Overlays and Turn-of-Year Effects	256
6.5.2	Cross-Currency Curve Construction	257
6.5.2.1	Basic Problem	257
6.5.2.2	Separation of Discount and Forward Rate Curves	258
6.5.2.3	Cross-Currency Basis Swaps	260
6.5.2.4	Modified Curve Construction Algorithm ..	261
6.5.3	Tenor Basis and Multi-Index Curve Group Construction	263
6.A	Appendix: Spline Theory	268
6.A.1	Hermite Spline Theory	268
6.A.2	C^2 Cubic Splines	271
6.A.3	C^2 Exponential Tension Splines	272
7	Vanilla Models with Local Volatility	275
7.1	General Framework	276
7.1.1	Model Dynamics	276
7.1.2	Volatility Smile and Implied Density	276

7.1.3	Choice of φ	277
7.2	CEV Model	278
7.2.1	Basic Properties	278
7.2.2	Call Option Pricing	280
7.2.3	Regularization	282
7.2.4	Displaced Diffusion Models	283
7.3	Quadratic Volatility Model	285
7.3.1	Case 1: Two Real Roots to the Left of $S(0)$	285
7.3.2	Case 2: One Real Root to the Left of $S(0)$	289
7.3.3	Extensions and Other Root Configurations	289
7.4	Finite Difference Solutions for General φ	290
7.4.1	Multiple λ and T	291
7.4.2	Forward Equation for Call Options	291
7.5	Asymptotic Expansions for General φ	293
7.5.1	Expansion around Displaced Log-Normal Process	293
7.5.2	Expansion around Gaussian Process	296
7.6	Extensions to Time-Dependent φ	297
7.6.1	Separable Case	297
7.6.2	Skew Averaging	298
7.6.2.1	Examples	302
7.6.2.2	A Caveat About the Process Domain	304
7.6.3	Skew and Convexity Averaging by Small-Noise Expansion	305
7.6.4	Numerical Example	309
8	Vanilla Models with Stochastic Volatility I	313
8.1	Model Definition	313
8.2	Model Parameters	315
8.3	Basic Properties	316
8.4	Fourier Integration	322
8.4.1	General Theory	322
8.4.2	Applications to SV Model	325
8.4.3	Numerical Implementation	328
8.4.4	Refinements of Numerical Implementation	330
8.4.5	Fourier Integration for Arbitrary European Payoffs	334
8.5	Integration in Variance Domain	337
8.6	CEV-Type Stochastic Volatility Models and SABR	341
8.7	Numerical Examples: Volatility Smile Statics	343
8.8	Numerical Examples: Volatility Smile Dynamics	345
8.9	Hedging in Stochastic Volatility Models	350
8.9.1	Hedge Construction, Delta and Vega	350
8.9.2	Minimum Variance Delta Hedging	353
8.9.3	Minimum Variance Hedging: an Example	354
8.A	Appendix: General Volatility Processes	356

9	Vanilla Models with Stochastic Volatility II	359
9.1	Fourier Integration with Time-Dependent Parameters	359
9.2	Asymptotic Expansion with Time-Dependent Volatility	362
9.3	Averaging Methods	366
9.3.1	Volatility Averaging	367
9.3.2	Skew Averaging	369
9.3.3	Volatility of Variance Averaging	370
9.3.4	Calibration by Parameter Averaging	372
9.4	PDE Method	377
9.4.1	PDE Formulation	377
9.4.2	Range for Stochastic Variance	378
9.4.3	Discretizing Stochastic Variance	379
9.4.4	Boundary Conditions for Stochastic Variance	381
9.4.5	Range for Underlying	382
9.4.6	Discretizing the Underlying	383
9.5	Monte Carlo Method	383
9.5.1	Exact Simulation of Variance Process	384
9.5.2	Biased Taylor-Type Schemes for Variance Process	385
9.5.2.1	Euler Schemes	385
9.5.2.2	Higher-Order Schemes	385
9.5.3	Moment Matching Schemes for Variance Process	386
9.5.3.1	Log-normal Approximation	386
9.5.3.2	Truncated Gaussian	387
9.5.3.3	Quadratic-Exponential	388
9.5.3.4	Summary of QE Algorithm	390
9.5.4	Broadie-Kaya Scheme for the Underlying	390
9.5.5	Other Schemes for the Underlying	392
9.5.5.1	Taylor-Type Schemes	392
9.5.5.2	Simplified Broadie-Kaya	392
9.5.5.3	Martingale Correction	392
9.A	Appendix: Proof of Proposition 9.3.4	393
9.B	Appendix: Coefficients for Asymptotic Expansion	397

VOLUME II Term Structure Models

Part III Term Structure Models

10 One-Factor Short Rate Models I	401
10.1 The One-Factor Gaussian Short Rate Model	402
10.1.1 The Ho-Lee Model	402
10.1.1.1 Notations and First Steps	402
10.1.1.2 Fitting the Term Structure of Discount Bonds	403
10.1.1.3 Analysis and Comparison with HJM Approach	405
10.1.2 The Mean-Reverting GSR Model	407
10.1.2.1 The Vasicek Model	407
10.1.2.2 The General One-Factor GSR Model . . .	409
10.1.2.3 Time-Stationarity and Caplet Hump . . .	412
10.1.3 European Option Pricing	414
10.1.3.1 The Jamshidian Decomposition	414
10.1.3.2 Gaussian Swap Rate Approximation . . .	416
10.1.4 Swaption Calibration	417
10.1.5 Finite Difference Methods	418
10.1.5.1 PDE and Spatial Boundary Conditions .	419
10.1.5.2 Determining Spatial Boundary Conditions from PDE	420
10.1.5.3 Upwinding	421
10.1.6 Monte Carlo Simulation	421
10.1.6.1 Exact Discretization	421
10.1.6.2 Approximate Discretization	423
10.1.6.3 Using other Measures for Simulation . .	424
10.2 The Affine One-Factor Model	425
10.2.1 Basic Definitions	425
10.2.1.1 SDE	425
10.2.1.2 Regularity Issues	426
10.2.1.3 Volatility Skew	426
10.2.1.4 Time-Dependent Parameters	427
10.2.2 Discount Bond Pricing and Extended Transform .	427
10.2.2.1 Constant Parameters	428
10.2.2.2 Piecewise Constant Parameters	430
10.2.3 Discount Bond Calibration	431
10.2.3.1 Change of Variables	431
10.2.3.2 Algorithm for $\omega(t)$	432
10.2.4 European Option Pricing	433

10.2.5	Swaption Calibration	435
10.2.5.1	Basic Problem	435
10.2.5.2	Calibration Algorithm	436
10.2.6	Quadratic One-Factor Model	437
10.2.7	Numerical Methods for the Affine Short Rate Model	437
11	One-Factor Short Rate Models II	439
11.1	Log-Normal Short Rate Models	439
11.1.1	The Black-Derman-Toy Model	439
11.1.2	Black-Karasinski Model	441
11.1.3	Issues in Log-Normal Models	441
11.1.4	Sandmann-Sondermann Transformation	442
11.2	Other Short Rate Models	445
11.2.1	Power-Type Models and Empirical Model Estimation	445
11.2.2	The Black Shadow Rate Model	446
11.2.3	Spanned and Unspanned Stochastic Volatility: the Fong and Vasicek Model	448
11.3	Numerical Methods for General One-Factor Short Rate Models	449
11.3.1	Finite Difference Methods	450
11.3.2	Calibration to Initial Yield Curve	451
11.3.2.1	Forward Induction	452
11.3.2.2	Forward-from-Backward Induction	453
11.3.2.3	Yield Curve and Volatility Calibration	455
11.3.2.4	The Dybvig Parameterization	457
11.3.2.5	Link to HJM Models	458
11.3.2.6	The Hagan and Woodward Parameterization	459
11.3.3	Monte Carlo Simulation	462
11.3.3.1	SDE Discretization	462
11.3.3.2	Practical Issues with Monte Carlo Methods	464
11.A	Appendix: Markov-Functional Models	466
11.A.1	State Process and Numeraire Mapping	466
11.A.2	Libor MF Parameterization	467
11.A.3	Swap MF Parameterization	469
11.A.4	Non-Parametric Calibration	470
11.A.5	Numerical Implementation	471
11.A.6	Comments and Comparisons	472

12 Multi-Factor Short Rate Models	473
12.1 The Gaussian Model	474
12.1.1 Development from Separability Condition	474
12.1.1.1 Mean-Reverting State Variables	475
12.1.1.2 Further Changes of Variables	479
12.1.2 Classical Development	481
12.1.2.1 Diagonalization of Mean Reversion Matrix	482
12.1.3 Correlation Structure	484
12.1.4 The Two-Factor Gaussian Model	485
12.1.4.1 Some Basics	485
12.1.4.2 Variance and Correlation Structure	486
12.1.4.3 Volatility Hump	487
12.1.4.4 Another Formulation of the Two-Factor Model	488
12.1.5 Multi-Factor Statistical Gaussian Model	491
12.1.6 Swaption Pricing	496
12.1.6.1 Jamshidian Decomposition	496
12.1.6.2 Gaussian Swap Rate Approximation	500
12.1.7 Calibration via Benchmark Rates	501
12.1.8 Monte Carlo Simulation	504
12.1.9 Finite Difference Methods	505
12.2 The Affine Model	506
12.2.1 Introduction	506
12.2.2 Basic Model	507
12.2.3 Regularity Issues	508
12.2.4 Discount Bond Prices	509
12.2.5 Some Concrete Models	511
12.2.5.1 Fong-Vasicek Model	511
12.2.5.2 Longstaff-Schwartz Model	512
12.2.5.3 Multi-Factor CIR Models	513
12.2.6 Brief Notes on Option Pricing	514
12.3 The Quadratic Gaussian Model	514
12.3.1 Quadratic Gaussian Models are Affine	515
12.3.2 The Basics	516
12.3.3 Parameterization	518
12.3.3.1 Smile Generation	518
12.3.3.2 Quadratic Term	519
12.3.3.3 Linear Term	521
12.3.4 Swaption Pricing	522
12.3.4.1 State Vector Distribution Under the Annuity Measure	522
12.3.4.2 Exact Pricing of European Swaptions	523
12.3.4.3 Approximations for European Swaptions	524
12.3.5 Calibration	527

12.3.6	Spanned Stochastic Volatility	528
12.3.7	Numerical Methods	528
12.A	Appendix: Quadratic Forms of Gaussian Vectors	528
13	The Quasi-Gaussian Model	533
13.1	One-Factor Quasi-Gaussian Model	533
13.1.1	Definition	533
13.1.2	Local Volatility	535
13.1.3	Swap Rate Dynamics	536
13.1.4	Approximate Local Volatility Dynamics for Swap Rate	537
13.1.4.1	Simple Approximation	538
13.1.4.2	Advanced Approximation	538
13.1.5	Linear Local Volatility	541
13.1.6	Linear Local Volatility for a Swaption Strip	543
13.1.7	Volatility Calibration	544
13.1.8	Mean Reversion Calibration	546
13.1.8.1	Effects of Mean Reversion	546
13.1.8.2	Calibrating Mean Reversion to Volatility Ratios	548
13.1.8.3	Calibrating Mean Reversion to Inter-Temporal Correlations	551
13.1.8.4	Final Comments on Mean Reversion Calibration	553
13.1.9	Numerical Methods	554
13.1.9.1	Direct Integration	554
13.1.9.2	Finite Difference Methods	556
13.1.9.3	Monte Carlo Simulation	559
13.1.9.4	Single-State Approximations	559
13.2	One-Factor Quasi-Gaussian Model with Stochastic Volatility	563
13.2.1	Definition	563
13.2.2	Swap Rate Dynamics	564
13.2.3	Volatility Calibration	566
13.2.4	Mean Reversion Calibration	567
13.2.5	Non-Zero Correlation	567
13.2.6	PDE and Monte Carlo Methods	568
13.3	Multi-Factor Quasi-Gaussian Model	568
13.3.1	General Multi-Factor Model	568
13.3.2	Local and Stochastic Volatility Parameterization ..	570
13.3.3	Swap Rate Dynamics and Approximations	572
13.3.4	Volatility Calibration	577
13.3.5	Mean Reversions, Correlations, and Numerical Schemes	578
13.A	Appendix: Density Approximation	579

13.A.1	Simplified Forward Measure Dynamics	579
13.A.2	Effective Volatility	580
13.A.3	The Forward Equation for Call Options	581
13.A.4	Asymptotic Expansion	582
13.A.5	Proof of Theorem 13.1.14	583
14	The Libor Market Model I	585
14.1	Introduction and Setup	586
14.1.1	Motivation and Historical Notes	586
14.1.2	Tenor Structure	587
14.2	LM Dynamics and Measures	587
14.2.1	Setting	587
14.2.2	Probability Measures	588
14.2.3	Link to HJM Analysis	591
14.2.4	Separable Deterministic Volatility Function	592
14.2.5	Stochastic Volatility	594
14.2.6	Time-Dependence in Model Parameters	597
14.3	Correlation	597
14.3.1	Empirical Principal Components Analysis	598
14.3.1.1	Example: USD Forward Rates	599
14.3.2	Correlation Estimation and Smoothing	600
14.3.2.1	Example: Fit to USD Data	603
14.3.3	Negative Eigenvalues	604
14.3.4	Correlation PCA	605
14.3.4.1	Example: USD Data	607
14.3.4.2	Poor Man's Correlation PCA	608
14.4	Pricing of European Options	608
14.4.1	Caplets	609
14.4.2	Swaptions	610
14.4.3	Spread Options	613
14.4.3.1	Term Correlation	614
14.4.3.2	Spread Option Pricing	615
14.5	Calibration	615
14.5.1	Basic Principles	615
14.5.2	Parameterization of $\ \lambda_k(t)\ $	616
14.5.3	Interpolation on the Whole Grid	617
14.5.4	Construction of $\lambda_k(t)$ from $\ \lambda_k(t)\ $	619
14.5.4.1	Covariance PCA	620
14.5.4.2	Correlation PCA	620
14.5.4.3	Discussion and Recommendation	621
14.5.5	Choice of Calibration Instruments	621
14.5.6	Calibration Objective Function	624
14.5.7	Sample Calibration Algorithm	626
14.5.8	Speed-Up Through Sub-Problem Splitting	627
14.5.9	Correlation Calibration to Spread Options	629

14.5.10	Volatility Skew Calibration	631
14.6	Monte Carlo Simulation	631
14.6.1	Euler-Type Schemes	632
14.6.1.1	Analysis of Computational Effort	633
14.6.1.2	Long Time Steps	634
14.6.1.3	Notes on the Choice of Numeraire	636
14.6.2	Other Simulation Schemes	636
14.6.2.1	Special-Purpose Schemes with Drift Predictor-Corrector	637
14.6.2.2	Euler Scheme with Predictor-Corrector	638
14.6.2.3	Lagging Predictor-Corrector Scheme	638
14.6.2.4	Further Refinements of Drift Estimation	640
14.6.2.5	Brownian-Bridge Schemes and Other Ideas	641
14.6.2.6	High-Order Schemes	643
14.6.3	Martingale Discretization	644
14.6.3.1	Deflated Bond Price Discretization	645
14.6.3.2	Comments and Alternatives	646
14.6.4	Variance Reduction	647
14.6.4.1	Antithetic Sampling	647
14.6.4.2	Control Variates	648
14.6.4.3	Importance Sampling	648
15	The Libor Market Model II	651
15.1	Interpolation	651
15.1.1	Back Stub, Simple Interpolation	652
15.1.2	Back Stub, Arbitrage-Free Interpolation	653
15.1.3	Back Stub, Gaussian Model	655
15.1.4	Front Stub, Zero Volatility	656
15.1.5	Front Stub, Exogenous Volatility	657
15.1.6	Front Stub, Simple Interpolation	660
15.1.7	Front Stub, Gaussian Model	661
15.2	Advanced Swap Pricing via Markovian Projection	662
15.2.1	Advanced Formula for Swap Rate Volatility	664
15.2.2	Advanced Formula for Swap Rate Skew	666
15.2.3	Skew and Smile Calibration in LM Models	668
15.3	Near-Markov LM Models	670
15.4	Swap Market Models	670
15.5	Evolving Separate Discount and Forward Rate Curves	672
15.5.1	Basic Ideas	673
15.5.2	HJM Extension	674
15.5.3	Applications to LM Models	677
15.5.4	Deterministic Spread	681
15.6	SV Models with Non-Zero Correlation	681
15.7	Multi-Stochastic Volatility Extensions	683

15.7.1	Introduction	683
15.7.2	Setup	684
15.7.3	Pricing Caplets and Swaptions	685
15.7.4	Spread Options.	686
15.7.5	Another Use of Multi-Dimensional Stochastic Volatility	687

VOLUME III Products and Risk Management

Part IV Products

16 Single-Rate Vanilla Derivatives	691
16.1 European Swaptions	691
16.1.1 Smile Dynamics	692
16.1.2 Adjustable Backbone.	693
16.1.3 Stochastic Volatility Swaption Grid	696
16.1.4 Calibrating Stochastic Volatility Model to Swaptions	697
16.1.5 Some Other Interpolation Rules	699
16.2 Caps and Floors.	700
16.2.1 Basic Problem	700
16.2.2 Setup and Norms	701
16.2.3 Calibration Procedure.	702
16.3 Terminal Swap Rate Models	703
16.3.1 TSR Basics	703
16.3.2 Linear TSR Model	705
16.3.3 Exponential TSR Model	708
16.3.4 Swap-Yield TSR Model	709
16.4 Libor-in-Arrears	710
16.5 Libor-with-Delay	713
16.5.1 Swap-Yield TSR Model	714
16.5.2 Other Terminal Swap Rate Models	715
16.5.3 Approximations Inspired by Term Structure Models	715
16.5.4 Applications to Averaging Swaps	716
16.6 CMS and CMS-Linked Cash Flows	717
16.6.1 The Replication Method for CMS.	718
16.6.2 Annuity Mapping Function as a Conditional Expected Value	720
16.6.3 Swap-Yield TSR Model	722
16.6.4 Linear and Other TSR Models	722
16.6.5 The Quasi-Gaussian Model	724
16.6.6 The Libor Market Model	725
16.6.7 Correcting Non-Arbitrage-Free Methods	728
16.6.8 Impact of Annuity Mapping Function and Mean Reversion	729
16.6.9 CDF and PDF of CMS Rate in Forward Measure. .	730
16.6.10 SV Model for CMS Rate	734

16.6.11	Dynamics of CMS Rate in Forward Measure	735
16.6.12	Cash-Settled Swaptions	738
16.7	Quanto CMS	740
16.7.1	Overview	740
16.7.2	Modeling the Joint Distribution of Swap Rate and Forward Exchange Rate	742
16.7.3	Normalizing Constant and Final Formula	743
16.8	Eurodollar Futures	744
16.8.1	Fundamental Results on Futures.....	745
16.8.2	Motivations and Plan	747
16.8.3	Preliminaries.....	748
16.8.4	Expansion Around the Futures Value	748
16.8.5	Forward Rate Variances	751
16.8.6	Forward Rate Correlations	753
16.8.7	The Formula	754
16.9	Convexity and Moment Explosions	755
17	Multi-Rate Vanilla Derivatives	759
17.1	Introduction to Multi-Rate Vanilla Derivatives	759
17.2	Marginal Distributions and Reference Measure	761
17.3	Dependence Structure via Copulas	762
17.3.1	Introduction to Gaussian Copula Method	762
17.3.2	General Copulas	764
17.3.3	Archimedean Copulas	766
17.3.4	Making Copulas from Other Copulas	767
17.4	Copula Methods for CMS Spread Options	770
17.4.1	Normal Model for the Spread	770
17.4.2	Gaussian Copula for Spread Options	771
17.4.3	Spread Volatility Smile Modeling with the Power Gaussian Copula	774
17.4.4	Copula Implied From Spread Options	775
17.5	Rates Observed at Different Times.....	778
17.6	Numerical Methods for Copulas	779
17.6.1	Numerical Integration Methods.....	780
17.6.2	Dimensionality Reduction for CMS Spread Options	783
17.6.3	Dimensionality Reduction for Other Multi-Rate Derivatives	785
17.6.4	Dimensionality Reduction by Conditioning.....	787
17.6.5	Dimensionality Reduction by Measure Change ...	791
17.6.6	Monte Carlo Methods	793
17.7	Limitations of the Copula Method	795
17.8	Stochastic Volatility Modeling for Multi-Rate Options....	796
17.8.1	Measure Change by Drift Adjustment	797
17.8.2	Measure Change by CMS Caplet Calibration	798
17.8.3	Impact of Correlations on the Spread Smile	799

XXVIII Contents

17.8.4	Connection to Term Structure Models	800
17.9	CMS Spread Options in Term Structure Models	802
17.9.1	Libor Market Model	802
17.9.2	Quadratic Gaussian Model	804
17.A	Appendix: Implied Correlation in Displaced Log-Normal Models	805
17.A.1	Preliminaries	805
17.A.2	Implied Log-Normal Correlation	806
17.A.3	A Few Numerical Results	807
18	Callable Libor Exotics	809
18.1	Model Calibration for Callable Libor Exotics	809
18.1.1	Risk Factors for CLEs	810
18.1.2	Model Choice and Calibration	813
18.2	Valuation Theory	814
18.2.1	Preliminaries	814
18.2.2	Recursion for Callable Libor Exotics	815
18.2.3	Marginal Exercise Value Decomposition	816
18.3	Monte Carlo Valuation	817
18.3.1	Regression-Based Valuation of CLEs, Basic Scheme	817
18.3.2	Regression for Underlying	819
18.3.3	Valuing CLE as a Cancelable Note	821
18.3.4	Using Regressed Variables for Decision Only	822
18.3.5	Regression Valuation with Boundary Optimization	824
18.3.6	Lower Bound via Regression Scheme	825
18.3.7	Iterative Improvement of Lower Bound	827
18.3.8	Upper Bound	830
18.3.8.1	Basic Ideas	830
18.3.8.2	Nested Simulation (NS) Algorithm	831
18.3.8.3	Bias and Computational Cost of NS Algorithm	834
18.3.8.4	Confidence Intervals and Practical Usage	836
18.3.8.5	Non-Analytic Exercise Values	837
18.3.8.6	Improvements to NS Algorithm	839
18.3.8.7	Other Upper Bound Algorithms	841
18.3.9	Regression Variable Choice	842
18.3.9.1	State Variables Approach	842
18.3.9.2	Explanatory Variables	843
18.3.9.3	Explanatory Variables with Convexity	846
18.3.10	Regression Implementation	848
18.3.10.1	Automated Explanatory Variable Selection	848
18.3.10.2	Suboptimal Point Exclusion	850
18.3.10.3	Two Step Regression	851

18.3.10.4	Robust Implementation of Regression Algorithm	852
18.4	Valuation with Low-Dimensional Models	856
18.4.1	Single-Rate Callable Libor Exotics	856
18.4.2	Calibration Targets for the Local Projection Method	856
18.4.3	Review of Suitable Local Models	857
18.4.4	Defining a Suitable Analog for Core Swap Rates	859
18.4.5	PDE Methods for Path-Dependent CLEs	861
18.4.5.1	CLEs Accreting at Coupon Rate	862
18.4.5.2	Snowballs	864
19	Bermudan Swaptions	867
19.1	Definitions	867
19.2	Local Projection Method	868
19.3	Smile Calibration	870
19.4	Amortizing, Accreting, Other Non-Standard Swaptions	872
19.4.1	Relationship Between Non-Standard and Standard Swap Rates	874
19.4.2	Same-Tenor Approach	875
19.4.3	Representative Swaption Approach	876
19.4.4	Basket Approach	879
19.4.5	Super-Replication for Non-Standard Bermudan Swaptions	882
19.4.6	Zero-Coupon Bermudan Swaptions	886
19.4.7	American Swaptions	887
19.4.7.1	American Swaptions vs. High- Frequency Bermudan Swaptions	888
19.4.7.2	The Proxy Libor Rate Method	889
19.4.7.3	The Libor-as-Extra-State Method	890
19.4.8	Mid-Coupon Exercise	891
19.5	Flexi-Swaps	892
19.5.1	Purely Global Bounds	893
19.5.2	Purely Local Bounds	893
19.5.3	Marginal Exercise Value Decomposition	895
19.5.4	Narrow Band Limit	896
19.6	Monte Carlo Valuation	897
19.6.1	Regression Methods	897
19.6.2	Parametric Boundary Methods	898
19.6.2.1	Sample Exercise Strategies for Bermudan Swaptions	898
19.6.2.2	Some Numerical Tests	901
19.6.2.3	Additional Comments	904
19.7	Other Topics	904

19.7.1	Robust Bermudan Swaption Hedging with European Swaptions	904
19.7.2	Carry and Exercise	907
19.7.3	Fast Pricing via Exercise Premia Representation	908
19.A	Appendix: Forward Volatility and Correlation	912
19.B	Appendix: A Primer on Moment Matching.....	913
19.B.1	Basics	913
19.B.2	Example 1: Asian Option in BSM Model	914
19.B.3	Example 2: Basket Option in BSM Model	916
20	TARNs, Volatility Swaps, and Other Derivatives	919
20.1	TARNs	919
20.1.1	Definitions and Examples	919
20.1.2	Valuation and Risk with Globally Calibrated Models	921
20.1.3	Local Projection Method	922
20.1.4	Volatility Smile Effects	923
20.1.5	PDE for TARNs.....	925
20.2	Volatility Swaps	927
20.2.1	Local Projection Method	928
20.2.2	Shout Options	929
20.2.3	Min-Max Volatility Swaps	932
20.2.4	Impact of Volatility Dynamics on Volatility Swaps	934
20.3	Forward Swaption Straddles	939
21	Out-of-Model Adjustments	945
21.1	Adjusting the Model	946
21.1.1	Calibration to Coupons	946
21.1.2	Adjusters	948
21.1.3	Path Re-Weighting	950
21.1.4	Proxy Model Method	955
21.1.5	Asset-Based Adjustments	957
21.1.6	Mapping Function Adjustments	959
21.2	Adjusting the Market	959
21.3	Adjusting the Trade	960
21.3.1	Fee Adjustments	961
21.3.2	Fee Adjustment Impact on Exotic Derivatives	962
21.3.3	Strike Adjustment	963

Part V Risk Management

22	Introduction to Risk Management	969
22.1	Risk Management and Sensitivity Computations	970
22.1.1	Basic Information Flow	970
22.1.2	Risk: Theory and Practice	972
22.1.3	Example: the Black-Scholes Model	974
22.1.4	Example: Black-Scholes Model with Time-Dependent Parameters	977
22.1.5	Actual Risk Computations	979
22.1.6	What about Θ_{prm} and Θ_{num} ?	980
22.1.7	A Note on Trading P&L and the Computation of Implied Volatility	981
22.2	P&L Analysis	984
22.2.1	P&L Predict	985
22.2.2	P&L Explain	987
22.2.2.1	Waterfall Explain	987
22.2.2.2	Bump-and-Reset Explain	988
22.3	Value-at-Risk	989
22.A	Appendix: Alternative Proof of Lemma 22.1.1	992
23	Payoff Smoothing and Related Methods	995
23.1	Issues with Discretization Schemes	995
23.1.1	Problems with Grid Dimensioning	996
23.1.2	Grid Shifts Relative to Payout	996
23.1.3	Additional Comments	999
23.2	Basic Techniques	1000
23.2.1	Adaptive Integration	1000
23.2.2	Adding Singularities to the Grid	1001
23.2.3	Singularity Removal	1003
23.2.4	Partial Analytical Integration	1004
23.3	Payoff Smoothing For Numerical Integration and PDEs	1006
23.3.1	Introduction to Payoff Smoothing	1006
23.3.2	Payoff Smoothing in One Dimension	1008
23.3.2.1	Box Smoothing	1009
23.3.2.2	Other Smoothing Methods	1012
23.3.3	Payoff Smoothing in Multiple Dimensions	1013
23.4	Payoff Smoothing for Monte Carlo	1016
23.4.1	Tube Monte Carlo for Digital Options	1016
23.4.2	Tube Monte Carlo for Barrier Options	1018
23.4.3	Tube Monte Carlo for Callable Libor Exotics	1023
23.4.4	Tube Monte Carlo for TARNs	1023
23.A	Appendix: Delta Continuity of Singularity-Enlarged Grid Method	1024
23.B	Appendix: Conditional Independence for Tube Monte Carlo	1026

24 Pathwise Differentiation	1029
24.1 Pathwise Differentiation: Foundations	1029
24.1.1 Callable Libor Exotics	1029
24.1.1.1 CLE Greeks	1030
24.1.1.2 Keeping the Exercise Time Constant	1032
24.1.1.3 Noise in CLE Greeks	1033
24.1.2 Barrier Options	1034
24.2 Pathwise Differentiation for PDE Based Models	1038
24.2.1 Model and Setup	1038
24.2.2 Bucketed Deltas	1039
24.2.3 Survival Density	1042
24.3 Pathwise Differentiation for Monte Carlo Based Models	1045
24.3.1 Pathwise Derivatives of Forward Libor Rates	1045
24.3.2 Pathwise Deltas of European Options	1048
24.3.2.1 Pathwise Deltas of the Numeraire	1048
24.3.2.2 Pathwise Deltas of the Payoff	1049
24.3.3 Adjoint Method For Greeks Calculation	1050
24.3.4 Pathwise Delta Approximation for Callable Libor Exotics	1052
24.4 Notes on Likelihood Ratio and Hybrid Methods	1054
25 Importance Sampling and Control Variates	1057
25.1 Importance Sampling In Short Rate Models	1057
25.2 Payoff Smoothing by Importance Sampling	1059
25.2.1 Binary Options	1059
25.2.2 TARNs	1062
25.2.3 Removing the First Digital	1062
25.2.4 Smoothing All Digitals by One-Step Survival Conditioning	1063
25.2.5 Simulating Under the Survival Measure Using Conditional Gaussian Draws	1066
25.2.6 Generalized Trigger Products in Multi-Factor LM Models	1068
25.3 Model-Based Control Variates	1071
25.3.1 Low-Dimensional Markov Approximation for LM models	1072
25.3.2 Two-Dimensional Extension	1075
25.3.3 Approximating Volatility Structure	1076
25.3.4 Markov Approximation as a Control Variate	1078
25.4 Instrument-Based Control Variates	1080
25.5 Dynamic Control Variates	1084
25.6 Control Variates and Risk Stability	1087

26 Vegas in Libor Market Models	1089
26.1 Basic Problem of Vega Computations	1089
26.2 Review of Calibration	1091
26.3 Vega Calculation Methods	1092
26.3.1 Direct Vega Calculations	1092
26.3.1.1 Definition and Analysis	1092
26.3.1.2 Numerical Example	1095
26.3.2 What is a Good Vega?	1096
26.3.3 Indirect Vega Calculations	1099
26.3.3.1 Definition and Analysis	1099
26.3.3.2 Numerical Example and Performance Analysis	1102
26.3.4 Hybrid Vega Calculations	1105
26.3.4.1 Definition and Analysis	1105
26.3.4.2 Numerical Example	1107
26.4 Skew and Smile Vegas	1107
26.5 Vegas and Correlations	1109
26.5.1 Term Correlation Effects	1109
26.5.2 What Correlations should be Kept Constant?	1110
26.5.3 Vegas with Fixed Term Correlations	1112
26.5.4 Numerical Example	1113
26.6 Deltas with Backbone	1114
26.7 Vega Projections	1116
26.8 Some Notes on Computing Model Vegas	1118

Appendix

A Markovian Projection	1123
A.1 Marginal Distributions of Ito Processes	1123
A.2 Approximations for Conditional Expected Values	1128
A.2.1 Gaussian Approximation	1128
A.2.2 Least-Squares Projection	1130
A.3 Applications to Local Stochastic Volatility Models	1131
A.3.1 Markovian Projection onto an SV Model	1131
A.3.2 Fitting the Market with an LSV Model	1133
A.3.3 On Calculating Proxy Local Volatility	1137
A.4 Basket Options in Local Volatility Models	1139
A.5 Basket Options in Stochastic Volatility Models	1143
A.A Appendix: $E(\sqrt{z_n(t)z_m(t)})$ and $E(\sqrt{z_n(t)})$	1146
A.A.1 Proof of Proposition A.A.1	1147
A.A.1.1 Step 1. Reduction to Covariance	1147
A.A.1.2 Step 2. Linear Approximation	1148
A.A.1.3 Step 3. Coefficients	1148
A.A.1.4 Step 4. Order of Approximation	1149
A.A.2 Proof of Lemma A.A.2	1149

Introduction to Arbitrage Pricing Theory

For reference, this chapter reviews selected results from stochastic calculus and from the modern theory of asset pricing. The material in this chapter is well covered in existing literature, so we keep the chapter brief and the mathematical treatment informal. For a more rigorous treatment we refer to Duffie [2001] or Musiela and Rutkowski [1997]. Most of the necessary mathematical foundation for the theory is available in Karatzas and Shreve [1997], Øksendal [1992], and Protter [2005].

The treatment in this chapter focuses on asset pricing in general; we shall specialize it to interest rate securities in Chapter 4. Chapter 5 introduces fixed income markets in detail.

1.1 The Setup

Unless otherwise noted, in this book we shall always consider an economy with continuous and frictionless trading taking place inside a finite horizon $[0, T]$. We assume the existence of traded dividend-free assets with prices characterized by a p -dimensional vector-valued stochastic process $X(t) = (X_1(t), \dots, X_p(t))^\top$. Uncertainty and information arrival is modeled by a probability space (Ω, \mathcal{F}, P) , with Ω being a sample space with outcome elements ω ; \mathcal{F} being a σ -algebra on Ω ; and P being a probability measure on the measure space (Ω, \mathcal{F}) . Information is revealed over time according to a filtration $\{\mathcal{F}_t, t \in [0, T]\}$, a family of sub- σ -algebras of \mathcal{F} satisfying $\mathcal{F}_s \subseteq \mathcal{F}_t$ whenever $s \leq t$. We can loosely think of \mathcal{F}_t as the information available at time t . We assume that the process $X(t)$ is adapted to $\{\mathcal{F}_t\}$, i.e. that $X(t)$ is fully observable at time t . For technical reasons, we require that the filtration satisfies the “usual conditions”¹. Let $E^P(\cdot)$ be the expectation

¹To satisfy the “usual conditions”, \mathcal{F}_t must be right-continuous for all t , and \mathcal{F}_0 must contain all the null-sets of \mathcal{F} , i.e. all subsets of sets of zero P -probability.

operator for the measure P ; when conditioning on information at time t , we will use the notation $E_t^P(\cdot) = E^P(\cdot | \mathcal{F}_t)$.

In all of the models in this book, we specialize the abstract setup above to the situation where information is generated by a d -dimensional vector-valued *Brownian motion* (or *Wiener process*) $W(t) = (W_1(t), \dots, W_d(t))^\top$, where W_i is independent of W_j for $i \neq j$. Brownian motions are treated in detail in Karatzas and Shreve [1997]; here, we just recall that a scalar Brownian motion W_i is a continuous stochastic process starting at 0 (i.e. $W_i(0) = 0$), having independent Gaussian increments: $W_i(t) - W_i(s) \sim \mathcal{N}(0, t-s)$, $t \geq s$. The filtration we consider is normally always the one *generated* by W , $\mathcal{F}_t = \sigma\{W(u), 0 \leq u \leq t\}$, possibly augmented to satisfy the usual conditions. We will generally assume that the price vector $X(t)$ is described by a vector-valued *Ito process*:

$$X(t) = X(0) + \int_0^t \mu(s, \omega) ds + \int_0^t \sigma(s, \omega) dW(s), \quad (1.1)$$

or, in differential notation,

$$dX(t) = \mu(t, \omega) dt + \sigma(t, \omega) dW(t), \quad (1.2)$$

where $\mu : \mathbb{R} \times \Omega \rightarrow \mathbb{R}^p$ and $\sigma : \mathbb{R} \times \Omega \rightarrow \mathbb{R}^{p \times d}$ are processes of dimension p and $p \times d$, respectively. We assume that both μ and σ are adapted to $\{\mathcal{F}_t\}$ and are in L^1 and L^2 respectively, in the sense that for all $t \in [0, T]$,

$$\int_0^t |\mu(s, \omega)| ds < \infty, \quad (1.3)$$

$$\int_0^t |\sigma(s, \omega)|^2 ds < \infty, \quad (1.4)$$

almost surely². In (1.4), we have defined

$$|\sigma(t, \omega)|^2 = \text{tr} (\sigma(t, \omega) \sigma(t, \omega)^\top). \quad (1.5)$$

We notice that the sample paths of X generated by (1.1) are almost surely continuous, with no jumps in asset prices.

A technical treatment of Ito processes and the Ito integral with respect to Brownian motion can be found in Karatzas and Shreve [1997]. For our needs, it suffices to think of the Ito integral as

$$\int_0^t \sigma(s, \omega) dW(s) = \lim_{n \rightarrow \infty} \sum_{i=1}^n \sigma((i-1)\delta, \omega) [W(i\delta) - W((i-1)\delta)], \quad (1.6)$$

²An event holds “almost surely” — often abbreviated by “a.s.” — if the probability of the event is one.

where $\delta \triangleq t/n$. We note that the integrand σ is here always evaluated at the *left* of each interval $[(i-1)\delta, i\delta]$. Other choices are possible³, but, as we shall see, the “non-anticipative” structure of the Ito integral gives rise to a number of useful results and makes it particularly useful as a model of trading gains (see Section 1.2).

We list a few relevant definitions and results below.

Definition 1.1.1 (Martingale). *Let $Y(t)$ be an adapted vector-valued process with $E^P(|Y(t)|) < \infty$ for all $t \in [0, T]$. We say that $Y(t)$ is a martingale under measure P if for all $s, t \in [0, T]$ with $t \leq s$,*

$$E_t^P(Y(s)) = Y(t), \quad a.s.$$

If we replace the equality sign in this equation with \leq or \geq , $Y(t)$ is said to be a *supermartingale* or a *submartingale*, respectively.

Definition 1.1.2 (Space H^2). *Let $|\sigma(t, \omega)|^2$ be as defined in (1.5). We say that σ is in H^2 , if for all $t \in [0, T]$ we have*

$$E^P \left(\int_0^t |\sigma(s, \omega)|^2 ds \right) < \infty.$$

The importance of Definition 1.1.2 becomes clear from the following result:

Theorem 1.1.3 (Properties of Ito Integral). *Define $I(t) = \int_0^t \sigma(s, \omega) dW(s)$ and assume that σ is in H^2 . Then*

1. $I(t)$ is \mathcal{F}_t -measurable.
2. $I(t)$ is a continuous martingale. In particular, $E^P(I(t)) = 0$ for all $t \in [0, T]$.
3. $E^P(|I(t)|^2) = E^P(\int_0^t |\sigma(s, \omega)|^2 ds) < \infty$.
4. $E^P(I(t)I(s)^\top) = E^P(\int_0^{\min(t,s)} \sigma(u, \omega) \sigma(u, \omega)^\top du)$.

A proof of Theorem 1.1.3 can be found in, e.g., Karatzas and Shreve [1997]. The equality in the third item of Theorem 1.1.3 is known as the *Ito isometry*. Due to the inequality in the third item, we say that the martingale defined in the process is a *square-integrable martingale*.

While it is common in applied work to simply assume that Ito integrals are martingales, without technical regularity conditions on $\sigma(t, \omega)$ (such as the H^2 restriction in Theorem 1.1.3), we should note that Ito integrals involving general processes in L^2 can, in fact, only be guaranteed to be *local martingales*. A process X is said to be a local martingale if there exists a

³The Stratonovich stochastic integral evaluates σ at the mid-point of each interval.

sequence of stopping times⁴ $\{\tau_n\}_{n=1}^\infty$, with $\tau_n \rightarrow \infty$ as $n \rightarrow \infty$, such that $X(\min(t, \tau_n))$, $t \geq 0$, is a martingale for all n . In other words, all “driftless” Ito processes of the type

$$dY(t) = \sigma(t, \omega) dW(t) \quad (1.7)$$

are local martingales, but not necessarily martingales. Interestingly, a converse result holds as well; all local martingales adapted to the filtration generated by the Brownian motion W can be represented as Ito processes of the form (1.7):

Theorem 1.1.4 (Martingale Representation Theorem). *If Y is a local martingale adapted to the filtration generated by a Brownian motion W , then there exists a process σ such that (1.7) holds. If Y is a square-integrable martingale, then σ is in H^2 .*

The proof of Theorem 1.1.4 can be found in Karatzas and Shreve [1997].

In the manipulation of functionals of Ito processes, the key result is a famous result by K. Ito:

Theorem 1.1.5 (Ito's Lemma). *Let $f(t, x)$, $x = (x_1, \dots, x_p)^\top$, denote a continuous function, $f : [0, T] \times \mathbb{R}^p \rightarrow \mathbb{R}$, with continuous partial derivatives $\partial f / \partial t = f_t$, $\partial f / \partial x_i = f_{x,i}$, $\partial^2 f / \partial x_i \partial x_j = f_{x,x,j}$. Let $X(t)$ be given by the Ito process (1.2) and define a scalar process $Y(t) = f(t, X(t))$. Then $Y(t)$ is an Ito process with stochastic differential*

$$\begin{aligned} dY(t) &= f_t(t, X(t)) dt + f_{x,t}(t, X(t)) \mu(t, \omega) dt + f_x(t, X(t)) \sigma(t, \omega) dW(t) \\ &\quad + \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p f_{x,x,j}(t, X(t)) (\sigma(t, \omega) \sigma(t, \omega)^\top)_{i,j} dt, \end{aligned}$$

where $f_x = (f_{x,1}, \dots, f_{x,p})$.

For easy reference, the result below lists Ito's lemma for the special case where $p = d = 1$.

Corollary 1.1.6. *For the case $p = d = 1$, Ito's lemma becomes*

$$\begin{aligned} dY(t) &= \left(f_t(t, X(t)) + f_x(t, X(t)) \mu(t, \omega) + \frac{1}{2} f_{xx}(t, X(t)) \sigma(t, \omega)^2 \right) dt \\ &\quad + f_x(t, X(t)) \sigma(t, \omega) dW(t). \end{aligned}$$

Ito's lemma can be motivated heuristically from a Taylor expansion. For instance, for the scalar case in Corollary 1.1.6, we write informally

⁴Recall that a stopping time τ is simply a random time adapted to the given filtration, in the sense that the event $\{\tau \leq t\}$ belongs to \mathcal{F}_t .

$$f(t + dt, X(t + dt)) = f(t, X(t)) + f_t dt + f_x dX(t) + \frac{1}{2} f_{xx} (dX(t))^2 + \dots \quad (1.8)$$

Here, we have

$$(dX(t))^2 = \mu(t, \omega)^2 (dt)^2 + \sigma(t, \omega)^2 (dW(t))^2 + 2\mu(t, \omega)\sigma(t, \omega) dt dW(t).$$

As shown earlier, $(dW(t))^2 = dt$ in quadratic mean, whereas all other terms in the expression for $(dX(t))^2$ are of order $O(dt^{3/2})$ or higher and can be neglected for small dt . In the limit, we therefore have $(dX(t))^2 = \sigma(t, \omega)^2 dt$ which can be inserted into (1.8). The result in Corollary 1.1.6 then emerges.

Remark 1.1.7. The quantity $(dX(t))^2$ discussed above is the differential of the *quadratic variation* of $X(t)$, often denoted by $\langle X(t), X(t) \rangle$. That is,

$$d\langle X(t), X(t) \rangle = (dX(t))^2 \Rightarrow \langle X(t), X(t) \rangle = \int_0^t (dX(u))^2.$$

For two different (scalar) Ito processes $X(t)$ and $Y(t)$, we may equivalently define the *quadratic covariation* process $\langle X(t), Y(t) \rangle$ by

$$d\langle X(t), Y(t) \rangle = dX(t) dY(t).$$

Sometimes we also write $d\langle X(t), Y(t) \rangle = \langle dX(t), dY(t) \rangle$. If $X(t)$ is a p -dimensional process and $Y(t)$ is a q -dimensional process, the quadratic covariation $\langle X(t), Y(t)^\top \rangle$ is a $(p \times q)$ -dimensional matrix process whose (i, j) -th element is $\langle X_i(t), Y_j(t) \rangle$, $i = 1, \dots, p$, $j = 1, \dots, q$.

The so-called *Tanaka extension* (see Karatzas and Shreve [1997]) extends Ito's lemma to continuous but non-differentiable functions. At points where the function has a kink, the Tanaka extension (loosely speaking) justifies using the Heaviside (step-) function for the first-order derivative and the Dirac delta function for the second-order derivative. An application of the Tanaka extension can be found in Section 1.9.2 and in Chapter 7, along with further discussion and references.

1.2 Trading Gains and Arbitrage

Working in the setting of Section 1.1 with assets driven by Ito processes, we now consider an investor engaging in a trading strategy involving the p assets X_1, \dots, X_p . Let the trading strategy be characterized by a predictable⁵ adapted process $\phi(t, \omega) = (\phi_1(t, \omega), \dots, \phi_p(t, \omega))^\top$, with $\phi_i(t, \omega)$ denoting

⁵A *predictable* process is one where we, loosely speaking, can “foretell” the value of the process at time t , given all information available up to, but not including, time t . All adapted continuous processes are thus predictable. For a technical definition of predictable processes, see Karatzas and Shreve [1997].

the holdings at time t in the i -th asset X_i . The value $\pi(t)$ of the trading strategy at time t is thus (dropping the dependence on ω in the notation)

$$\pi(t) = \phi(t)^\top X(t). \quad (1.9)$$

The gain from trading over a small time interval $[t, t + \delta]$ is (approximately) $\phi(t)^\top [X(t + \delta) - X(t)]$, suggesting (compare to (1.6)) that the Ito integral

$$\int_0^t \phi(s)^\top dX(s) = \int_0^t \phi(s)^\top \mu(s) ds + \int_0^t \phi(s)^\top \sigma(s) dW(s)$$

is a proper model for trading gains over $[0, t]$. An investment strategy is said to be *self-financing* if, for any $t \in [0, T]$,

$$\pi(t) - \pi(0) = \int_0^t \phi(s)^\top dX(s). \quad (1.10)$$

This relationship simply expresses that changes in portfolio value are solely caused by trading gains or losses, with no funds being added or withdrawn.

Self-financing trading strategies allow investors to turn a certain initial investment $\pi(0)$ into stochastic future wealth $\pi(t)$. Under natural assumptions on possible trading strategies (e.g., that there is finite supply of all assets) we would expect that there should be limitations to the profits that self-financing strategies can create. Most notably, it should be impossible to create “something for nothing”, that is, to turn a zero initial investment into future wealth that is certain to be non-negative and may be positive with non-zero probability. To express this formally, we introduce the concept of an *arbitrage opportunity*:

Definition 1.2.1 (Arbitrage). *An arbitrage opportunity is a self-financing strategy ϕ for which $\pi(0) = 0$ and, for some $t \in [0, T]$,*

$$\pi(t) \geq 0 \text{ a.s., and } P(\pi(t) > 0) > 0, \quad (1.11)$$

with π given in (1.9).

In economic equilibrium, arbitrage strategies cannot exist and precluding (1.11) constitutes a fundamental consistency requirement on the asset processes.

1.3 Equivalent Martingale Measures and Arbitrage

We turn to the question of characterizing the conditions under which the trading economy is free of arbitrage opportunities. A concise way to state these conditions involves *equivalent martingale measures*, a concept we shall work our way up to in a number of steps. First, we recall that two

probability measures P and \widehat{P} on the same measure space (Ω, \mathcal{F}) are said to be *equivalent* if $P(A) = 0 \Leftrightarrow \widehat{P}(A) = 0, \forall A \in \mathcal{F}$; that is, the two measures have the same null-sets. An important result from measure theory states that equivalent measures are uniquely associated through a quantity known as a *Radon-Nikodym derivative*:

Theorem 1.3.1 (Radon-Nikodym Theorem). *Let P and \widehat{P} be equivalent probability measures on the common measure space (Ω, \mathcal{F}) . There exists a unique (a.s.) non-negative random variable R with $E^P(R) = 1$, such that*

$$\widehat{P}(A) = E^P(R1_{\{A\}}), \quad \text{for all } A \in \mathcal{F}.$$

For a proof of Theorem 1.3.1, see e.g. Billingsley [1995]. The random variable R in the theorem is known as a *Radon-Nikodym derivative* and is denoted $d\widehat{P}/dP$. In the theorem we have used an *indicator* $1_{\{A\}}$; this quantity is 1 if the event A comes true, 0 if not.

For later use, we associate any probability measure \widehat{P} with a *density process*

$$\varsigma(t) = E_t^P \left(\frac{d\widehat{P}}{dP} \right), \quad \forall t \in [0, T]. \quad (1.12)$$

Clearly, $\varsigma(t)$ is a P -martingale with $\varsigma(0) = 1$ and $\varsigma(t) = E_t^P(\varsigma(T))$. A simple conditioning exercise demonstrates that for any \mathcal{F}_T -measurable random variable $Y(T)$, with $R = d\widehat{P}/dP$,

$$\begin{aligned} E^{\widehat{P}}(Y(T)|\mathcal{F}_t) &= \frac{1}{E^P(R|\mathcal{F}_t)} E^P(RY(T)|\mathcal{F}_t) \\ &= \varsigma(t)^{-1} E^P(E^P(R|\mathcal{F}_T)Y(T)|\mathcal{F}_t) \\ &= E^P \left(Y(T) \frac{\varsigma(T)}{\varsigma(t)} \middle| \mathcal{F}_t \right). \end{aligned} \quad (1.13)$$

We shall use this result on numerous occasions in this book.

We now introduce the important concept of a *deflator*, a strictly positive Ito process used to normalize the asset prices. Let the deflator be denoted $D(t)$, and define the normalized asset process $X^D(t) = (X_1(t)/D(t), \dots, X_p(t)/D(t))^{\top}$. We say that a measure Q^D is an *equivalent martingale measure induced by D* if $X^D(t)$ is a martingale with respect to Q^D . If Q^D is a martingale measure, we say that a self-financing trading strategy is *permissible* if

$$\int_0^t \phi(s)^{\top} dX^D(s)$$

is a martingale. For the Ito setup discussed earlier, a permissible strategy⁶ is obtained by, say, requiring that $\phi(t)^{\top} \sigma(t)$ is in H^2 ; see Theorem 1.1.3. An

⁶The technical restriction on trading positions imposed by only considering permissible trading strategies rules out certain pathological strategies, such as the

application of Ito's lemma combined with (1.9)–(1.10) implies that $\pi(t)/D(t)$ is a Q^D -martingale when the trading strategy is permissible.

For permissible trading strategies, the importance of equivalent martingale measures follows from the following theorem:

Theorem 1.3.2 (Sufficient Condition for No-Arbitrage). *Restrict attention to permissible trading strategies. If there is a deflator D such that the deflated asset price process allows for an equivalent martingale measure, then there is no arbitrage.*

For a proof we refer to Musiela and Rutkowski [1997]. We note that Theorem 1.3.2 only provides sufficient conditions for the absence of arbitrage, and known (and rather technical) counterexamples demonstrate that the existence of an equivalent martingale measure does not follow from the absence of arbitrage in a setting with permissible trading strategies. A body of results known as the *fundamental theorem of arbitrage* establishes the conditions under which the existence of an equivalent martingale measure is also a necessary condition for the absence of arbitrage. The results are rather technical, but generally state that absence of arbitrage and the existence of an equivalent martingale measure are “nearly” equivalent concepts. The exact notion of “nearly” equivalent is discussed in Duffie [2001] as well as in the authoritative reference⁷ Delbaen and Schachermayer [1994]. For our purposes in this book, we ignore many of these technicalities and often simply treat the absence of arbitrage and the existence of a martingale measure as equivalent concepts.

Finally, if the deflator is one of the p assets, we call the deflator a *numeraire*. Let us, say, assume that X_1 is strictly positive and can be used as a numeraire. Also assume that a deflator D has been identified such that Theorem 1.3.2 holds. As $X_1(t)/D(t)$ is a Q^D -martingale, we can use the Radon-Nikodym theorem to define a new measure Q^{X_1} by the density $\varsigma(t) = (X_1(t)/D(t))/(X_1(0)/D(0))$. For an \mathcal{F}_T -measurable variable $Y(T)$, we then have, from (1.13),

$$X_1(t)E_t^{Q^{X_1}}\left(\frac{Y(T)}{X_1(T)}\right) = D(t)E_t^{Q^D}\left(\frac{Y(T)}{D(T)}\right). \quad (1.14)$$

In particular, if $Y(t)/D(t)$ is a Q^D -martingale, $Y(t)/X_1(t)$ must also be a Q^{X_1} -martingale. In practice, it normally suffices to only consider deflators from the set of available numeraires.

Remark 1.3.3. Some sources define $1/D(t)$ (rather than $D(t)$) as the deflator. The convention used in this book is more natural for our applications.

doubling strategy considered in Harrison and Kreps [1979]. A realistic resource-constrained economy will always bound the size of the positions one can take in an asset, sufficing to ensure that predictable trading strategies are permissible.

⁷In a nutshell, Delbaen and Schachermayer [1994] show that absence of arbitrage implies only the existence of a *local* martingale measure.

1.4 Derivative Security Pricing and Complete Markets

A T -maturity *derivative security* (also known as a *contingent claim*) pays out at time T an \mathcal{F}_T -measurable random variable $V(T)$, and makes no payments before T . We assume that $V(T)$ has finite variance, and say that the derivative security is *attainable* (or sometimes *redundant*) if there exists a permissible trading strategy ϕ such that $V(T) = \phi(T)^\top X(T) = \pi(T)$ a.s. The trading strategy is said to *replicate* the derivative security. Importantly, the absence of arbitrage dictates that the time 0 price of an attainable derivative security $V(0)$ must be equal to the cost of setting up the self-financing strategy, i.e. $V(0) = \pi(0)$. More generally, $V(t) = \pi(t)$, $t \in [0, T]$. This observation is the foundation of *arbitrage pricing* and allows us to price derivative securities as expectations under an equivalent martingale measure. Specifically, consider a deflator D and assume the existence of an equivalent martingale measure Q^D induced by D ; the existence of Q^D guarantees that there are no arbitrages in the market, by Theorem 1.3.2. Now, from the martingale property of $\pi(t)/D(t)$ in the measure Q^D and the relation $V(t) = \pi(t)$ it immediately follows that

$$\frac{V(t)}{D(t)} = E_t^{Q^D} \left(\frac{V(T)}{D(T)} \right)$$

or

$$V(t) = D(t) E_t^{Q^D} \left(\frac{V(T)}{D(T)} \right). \quad (1.15)$$

If all finite-variance \mathcal{F}_T -measurable random variables can be replicated, the market is said to be *complete*. In a complete market, all derivatives are “spanned” and hence have unique prices. Interestingly, a similar uniqueness result holds for equivalent martingale measures:

Theorem 1.4.1. *In the absence of arbitrage, a market is complete if and only if there exists a deflator inducing a unique martingale measure.*

From (1.14) it follows that the martingale measures induced by all numeraires must then be unique as well.

In practical applications, we shall often manipulate the choice of numeraire asset to simplify computations. The following result is useful for this:

Theorem 1.4.2 (Change of Numeraire). *Consider two numeraires $N(t)$ and $M(t)$, inducing equivalent martingale measures Q^N and Q^M , respectively. If the market is complete, then the density of the Radon-Nikodym derivative relating the two measures is uniquely given by*

$$\varsigma(t) = E_t^{Q^N} \left(\frac{dQ^M}{dQ^N} \right) = \frac{M(t)/M(0)}{N(t)/N(0)}.$$

Proof. As the market is complete, all derivatives prices are unique. Consider an integrable \mathcal{F}_T -measurable payout $V(T) = Y(T)M(T)$, with time t price $V(t)$. From Theorem 1.4.1 and (1.15) we must have

$$V(t) = N(t)\mathbb{E}_t^{\mathbb{Q}^N} \left(\frac{M(T)Y(T)}{N(T)} \right) = M(t)\mathbb{E}_t^{\mathbb{Q}^M} \left(\frac{M(T)Y(T)}{M(T)} \right)$$

or

$$\mathbb{E}_t^{\mathbb{Q}^M} (Y(T)) = \mathbb{E}_t^{\mathbb{Q}^N} \left(Y(T) \frac{M(T)/N(T)}{M(t)/N(t)} \right).$$

Comparison with (1.13), and the fact that the density must be scaled to equal 1 at time 0, reveals that the Radon-Nikodym derivative for the measure shift is characterized by the density in the theorem. \square

1.5 Girsanov's Theorem

The last two sections have demonstrated a close link between the concept of arbitrage and the existence and uniqueness of equivalent martingale measures. In this section, we consider i) the conditions on the asset prices that allow for an equivalent martingale measure; and ii) the effect on asset dynamics from a change of probability measure. We consider two measures P and $P(\theta)$ related by a density $\varsigma^\theta(t) = \mathbb{E}_t^P(dP(\theta)/dP)$, where $\varsigma^\theta(t)$ is an *exponential martingale* given by the Ito process

$$d\varsigma^\theta(t)/\varsigma^\theta(t) = -\theta(t)^\top dW(t),$$

where $W(t)$ is a d -dimensional P -Brownian motion. The d -dimensional process θ is known as the *market price of risk*. By an application of Ito's lemma, we can write

$$\begin{aligned} \varsigma^\theta(t) &= \exp \left(- \int_0^t \theta(s)^\top dW(s) - \frac{1}{2} \int_0^t \theta(s)^\top \theta(s) ds \right) \\ &\triangleq \mathcal{E} \left(- \int_0^t \theta(s)^\top dW(s) \right) \end{aligned} \tag{1.16}$$

where $\mathcal{E}(\cdot)$ is the *Doleans exponential*. An often-quoted sufficient condition on $\theta(t)$ for (1.16) to define a proper martingale (and not just a local martingale) is the *Novikov condition*

$$\mathbb{E}^P \left[\exp \left(\frac{1}{2} \int_0^t \theta(s)^\top \theta(s) ds \right) \right] < \infty. \tag{1.17}$$

The Novikov condition can often be difficult to verify in practical applications.

Armed with the notation above, we are now ready to state the main result of this section.

Theorem 1.5.1 (Girsanov's Theorem). Suppose that $\varsigma^\theta(t)$ defined in (1.16) is a martingale. Then for all $t \in [0, T]$

$$W^\theta(t) = W(t) + \int_0^t \theta(s) ds$$

is a Brownian motion under the measure $P(\theta)$.

To discuss a strategy to prove Girsanov's theorem, assume for simplicity that the dimension of the Brownian motion is $d = 1$. One way to construct a proof for Theorem 1.5.1 is to demonstrate that the joint moment-generating function (mgf)⁸ (under $P(\theta)$) of the increments

$$W^\theta(t_1), W^\theta(t_2) - W^\theta(t_1), \dots, W^\theta(t_n) - W^\theta(t_{n-1}), \quad 0 < t_1 < \dots < t_n,$$

is the same as that of n independent Gaussian random variables with expectations 0 and variances $t_1, t_2 - t_1, \dots$. That is, for any positive integer value of n and any set of values $\alpha_i \in \mathbb{R}$, $i = 1, 2, \dots, n$, we need to show that,

$$E^{P^\theta} \left[\exp \left(\sum_{i=1}^n \alpha_i (W^\theta(t_i) - W^\theta(t_{i-1})) \right) \right] = \prod_{i=1}^n \exp (\alpha_i^2 (t_i - t_{i-1}) / 2),$$

where we have defined $t_0 = 0$. While carrying out such a proof is not difficult, we here merely justify the final result by examining the case $n = 1$ only. Specifically, we consider

$$E^{P(\theta)} [\exp (\alpha W^\theta(t))],$$

where $\alpha \in \mathbb{R}$ and $t > 0$. Shifting probability measure, we get

$$\begin{aligned} E^{P(\theta)} [\exp (\alpha W^\theta(t))] &= E^{P(\theta)} \left[\exp \left(\alpha W(t) + \alpha \int_0^t \theta(s) ds \right) \right] \\ &= E^P \left[\exp \left(\alpha W(t) + \alpha \int_0^t \theta(s) ds \right) \mathcal{E} \left(- \int_0^t \theta(s) dW(s) \right) \right] \\ &= e^{\alpha^2 t / 2} E^P \left[\exp \left(\int_0^t (\alpha - \theta(s)) dW(s) - \frac{1}{2} \int_0^t (\alpha - \theta(s))^2 ds \right) \right] \\ &= e^{\alpha^2 t / 2} E^P \left[\mathcal{E} \left(\int_0^t (\alpha - \theta(s)) dW(s) \right) \right] \\ &= e^{\alpha^2 t / 2}, \end{aligned}$$

⁸Recall that the moment-generating function of a random variable Y in some measure P is defined as the expectation $E^P(\exp(\alpha Y))$, $\alpha \in \mathbb{R}$. Unlike the characteristic function, the moment-generating function is not always well-defined for all values of the argument α .

as desired. In the last step, we used the fact that the Doleans exponential is a martingale with initial value 1.

Girsanov's theorem implies that we can shift probability measure to transform an Ito process with a given drift to an Ito process with nearly arbitrary drift. Specifically, we notice that our asset price process (under P)

$$dX(t) = \mu(t) dt + \sigma(t) dW(t)$$

can be written

$$dX(t) = (\mu(t) - \sigma(t)\theta(t)) dt + \sigma(t)dW^\theta(t),$$

where $W^\theta(t)$ is a Brownian measure under the measure $P(\theta)$. This process will be driftless provided that θ satisfies the “spanning condition” $\mu(t) = \sigma(t)\theta(t)$ for all $t \in [0, T]$. This gives us a convenient way to check for the existence of equivalent martingale measures:

Corollary 1.5.2. *For a given numeraire D, assume that the deflated asset process satisfies*

$$dX^D(t) = \mu^D(t) dt + \sigma^D(t) dW(t),$$

where $\sigma^D(t)$ is sufficiently regular to make $\int_0^t \sigma^D(s) dW(s)$ a martingale. Assume also that there exists a θ such that the density ς^θ is a martingale and (a.s.)

$$\sigma^D(t)\theta(t) = \mu^D(t), \quad t \in [0, T], \quad (1.18)$$

then D induces an equivalent martingale measure and there is no arbitrage.

Equation (1.18) is a system of linear equations and we can use rank results from linear algebra to determine the circumstances under which (1.18) will have solutions (no arbitrage) and when these are unique (complete market). For instance, a necessary condition for the market to be complete is that $\text{rank}(\sigma) = d$. Further results along these lines can be found in Musiela and Rutkowski [1997] and Duffie [2001].

We conclude this section by noting that while a change of probability measure affects the drift μ of an Ito process, it does not change the diffusion coefficient σ . This is sometimes known as the *diffusion invariance principle*.

1.6 Stochastic Differential Equations

So far we have defined the asset process vector to be an Ito process with general measurable coefficients $\mu(t, \omega)$ and $\sigma(t, \omega)$. In virtually all applications, however, we restrict our attention to the case where these coefficients

are deterministic functions of time and the state of the asset process⁹. In other words, we consider a *stochastic differential equation* (SDE) of the form

$$dX(t) = \mu(t, X(t)) dt + \sigma(t, X(t)) dW(t), \quad X(0) = X_0, \quad (1.19)$$

with $\mu : [0, T] \times \mathbb{R}^p \rightarrow \mathbb{R}^p$; $\sigma : [0, T] \times \mathbb{R}^p \rightarrow \mathbb{R}^{p \times d}$; and X_0 an initial condition. A *strong solution*¹⁰ to (1.19) is an Ito process

$$X(t) = X_0 + \int_0^t \mu(s, X(s)) ds + \int_0^t \sigma(s, X(s)) dW(s).$$

A number of restrictions on μ and σ are needed to ensure that the solution to (1.19) exists and is unique. A standard result is listed below.

Theorem 1.6.1. *In (1.19) assume that there exists a constant K such that for all $t \in [0, T]$ and all $x, y \in \mathbb{R}^p$,*

$$\begin{aligned} |\mu(t, x) - \mu(t, y)| + |\sigma(t, x) - \sigma(t, y)| &\leq K|x - y|, & (\text{Lipschitz condition}), \\ |\mu(t, x)|^2 + |\sigma(t, x)|^2 &\leq K^2(1 + |x|^2), & (\text{growth condition}). \end{aligned}$$

Then there exists a unique solution to (1.19).

We notice that the dynamics of (1.19) do not depend on the past evolution of $X(t)$ beyond the state of X at time t . This lack of path-dependence suggests that X is a *Markov* process. We formalize this as follows.

Definition 1.6.2 (Markov Process). *The \mathbb{R}^p -valued stochastic process $X(t)$ is called a *Markov process* if for all $s, t \in [0, T]$ with $t \leq s$,*

$$\mathbb{P}(X(s) \in B | \mathcal{F}_t) = \mathbb{P}(X(s) \in B | X(t)) \quad (1.20)$$

for all sets B in the p -dimensional σ -algebra of Borel set \mathfrak{B}^p . If (1.20) holds with s replaced by a stopping time, the process is a strong Markov process.

Expressed verbally, the Markov property implies that the past and future become statistically independent when we condition on the present.

Theorem 1.6.3 (Markov Property of SDEs). *Let the coefficient of the SDE for $X(t)$ satisfy the conditions in Theorem 1.6.1. Then $X(t)$ is a strong Markov process.*

⁹In this section, the process X is generic and need not represent financial assets.

¹⁰In a strong solution, the Brownian motion is given and the solution is adapted to the filtration generated by it. If we are free to pick our own Brownian motion on some different probability space, we say that (1.19) holds in a *weak sense*. For financial applications where we normally only need the law of the underlying process, weak solutions are typically sufficient. The distinction between weak and strong solutions is of little importance for our purposes and we shall ignore it going forward.

Let us consider the explicit solutions of a few simple SDEs. First, consider a *linear* SDE

$$dX(t) = (AX(t) + B(t)) dt + C(t) dW(t),$$

where A is a constant $p \times p$ matrix, and B and C are deterministic matrices of dimension $p \times 1$ and $p \times d$, respectively. The solution to this equation can, by Ito's lemma, be verified to be

$$X(t) = e^{At} X(0) + \int_0^t e^{A(t-s)} (B(s) ds + C(s) dW(s)).$$

The term

$$\int_0^t e^{A(t-s)} C(s) dW(s)$$

is distributed as a p -dimensional Gaussian random variable with mean 0 and, from Theorem 1.1.3, covariance matrix

$$\Sigma = \int_0^t e^{A(t-s)} C(s) C(s)^\top e^{A^\top(t-s)} ds.$$

Extensions to time-varying A are straightforward, and basically involve replacing the exponential matrix e^{At} with the solution of a homogeneous ODE with time-dependent coefficients. Details can be found in, e.g., Arnold [1974] and key results are listed in Chapter 12.

Now let us specialize to the scalar case with $p = 1$. An SDE of great importance is the *geometric Brownian motion with drift* (GBMD):

$$dX(t)/X(t) = \mu(t) dt + \sigma(t) dW(t),$$

where $\mu(t)$ and $\sigma(t)$ are *deterministic* (with $\sigma(t)$ having dimension $1 \times d$). An application of Ito's lemma to $\ln(X(t))$ reveals that

$$\begin{aligned} X(t) &= X(0) \exp \left(\int_0^t \left(\mu(s) - \frac{1}{2} \sigma(s) \sigma(s)^\top \right) ds + \int_0^t \sigma(s) dW(s) \right) \\ &= X(0) \exp \left(\int_0^t \mu(s) ds \right) \mathcal{E} \left(\int_0^t \sigma(s) dW(s) \right). \end{aligned} \quad (1.21)$$

Being an exponential of a Gaussian random variable, $X(t)$ follows a *log-normal* distribution, with moments (see Karatzas and Shreve [1997])

$$\mathbb{E}^P(X(t)) = X(0) \exp \left(\int_0^t \mu(s) ds \right), \quad (1.22)$$

$$\mathbb{E}^P(X(t)^2) = \mathbb{E}^P(X(t))^2 \exp \left(\int_0^t \sigma(s) \sigma(s)^\top ds \right). \quad (1.23)$$

1.7 Explicit Trading Strategies and PDEs

After the mathematical interlude of Section 1.6, we now return to financial markets and a more careful analysis of the trading strategies that replicate derivative securities. We have already established that in a complete market such strategies must exist for any given derivative, but it still remains to determine these strategies explicitly. Consider a Markovian setup where the asset vector X satisfies an SDE of the form (1.19). Let there be given a derivative security V paying out at time T an amount $V(T) = g(X(T))$, for some smooth payout function $g : \mathbb{R}^p \rightarrow \mathbb{R}$. The Markovian form of the asset dynamics suggests that the time t derivative price is a function of t and $X(t)$ only, $V(t) = V(t, X(t))$ for some deterministic function $V(t, x)$, $x \in \mathbb{R}^p$. Conjecturing that this function is smooth enough to allow for an application of Ito's lemma for all $t \in [0, T]$, Theorem 1.1.5 implies (suppressing dependence on $X(t)$ for brevity)

$$\begin{aligned} dV(t) &= V_t(t) dt + \sum_{i=1}^p V_{x_i}(t) \mu_i(t) dt \\ &\quad + \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p V_{x_i x_j}(t) \Sigma_{i,j}(t) dt + \sum_{i=1}^p V_{x_i}(t) \sigma_i(t) dW(t), \end{aligned} \quad (1.24)$$

where σ_i is the i -th row of the $p \times d$ matrix σ and $\Sigma_{i,j}$ is the (i, j) -th element in $\sigma \sigma^\top$. We recall that subscripts like V_X , denote partial differentiation, see Theorem 1.1.5.

If $V(t)$ can be replicated by a self-financing trading strategy ϕ in the p assets, we must also have, from (1.10),

$$dV(t) = \phi(t)^\top dX(t) = \sum_{i=1}^p \phi_i(t) \mu_i(t) dt + \sum_{i=1}^p \phi_i(t) \sigma_i(t) dW(t). \quad (1.25)$$

Comparing terms in (1.24) and (1.25) we see that both equations will hold, provided that for all $t \in [0, T]$

$$\phi_i(t) = \frac{\partial V(t, X(t))}{\partial x_i}, \quad i = 1, \dots, p, \quad (1.26)$$

and

$$\frac{\partial V(t, x)}{\partial t} + \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p \frac{\partial^2 V(t, x)}{\partial x_i \partial x_j} \Sigma_{i,j}(t, x) = 0. \quad (1.27)$$

To the extent that the system above allows for a solution (it may not if the market is not complete, from (1.26) we see that the trading strategy that replicates the derivative V holds $\partial V(t, X(t))/\partial x_i$ units of asset X_i at time

t. The quantity $\partial V / \partial x_i$ is often known as the *delta* with respect to X_i ¹¹. Note that, from (1.9) and (1.26) we have that

$$V(t, X(t)) = \sum_{i=1}^p \frac{\partial V(t, X(t))}{\partial x_i} X_i(t). \quad (1.28)$$

Besides identifying an explicit replication strategy, the arguments above have also produced (1.27), a partial differential equation (PDE) for the value function $V(t, x)$. The PDE is a second-order parabolic equation in p spatial variables, with known terminal condition $V(T, x) = g(x)$ (a so-called *Cauchy problem*). Solving this PDE provides an alternative way to price the derivative, as compared to the purely probabilistic expectations-based methods outlined earlier (see (1.15)). We shall investigate the link between expectations and PDEs in more detail in Section 1.8.

Inspection of the valuation PDE (1.27) reveals that the drifts μ_i of the asset price SDE (1.19) are notably absent, making the price of the derivative security independent of drifts. This is typical of derivatives in complete markets and follows from the fact that derivatives can be priced preference-free, by arbitrage arguments. In contrast, for the elements of the fundamental asset price vector, risk-averse investors would demand that assets with high volatilities $|\sigma_i|$ be rewarded with higher drifts (more precisely, higher rates of return) as compensation for the additional uncertainty.

1.8 Kolmogorov's Equations and the Feynman-Kac Theorem

In earlier sections, we have seen that derivatives prices can be expressed as expectations under certain probability measures or as solutions to PDEs. This hints at a deeper connection between expectations and PDEs, a connection we shall explore in this section. As part of this exploration, we list results for transition densities that will be useful later in model calibration.

As in Section 1.6, we consider a Markov vector SDE of the type (see (1.19))

$$dX(t) = \mu(t, X(t)) dt + \sigma(t, X(t)) dW(t), \quad X(0) = X_0, \quad (1.29)$$

where the coefficients are assumed smooth enough to allow for a unique solution (see Theorem 1.6.1). Now define a functional

$$u(t, x) = \mathbb{E}^P(g(X(T)) | X(t) = x),$$

¹¹Note that taking a position in V and following a trading strategy with $\phi_i = -\partial V / \partial x_i$, $i = 1, \dots, p$ will effectively remove any exposure to V (as we simultaneously take a long position in V and, through a trading strategy, a short position in V). This strategy is known as a *delta hedge*.

for a function $g : \mathbb{R}^p \rightarrow \mathbb{R}$. Under regularity conditions on g , it is easy to see that the process $u(t, X(t))$, being a conditional expectation, must be a martingale. Proceeding informally, an application of Ito's lemma gives, for $t \in [0, T)$ (suppressing dependence on $X(t)$),

$$du(t) = u_t(t) dt + \sum_{i=1}^p u_{x_i}(t) \mu_i(t) dt + \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p u_{x_i x_j}(t) \Sigma_{i,j}(t) dt + O(dW(t)),$$

where as before $\Sigma_{i,j}$ is the (i, j) -th element of $\sigma\sigma^\top$. From earlier results, we know that for $u(t, X(t))$ to be a martingale, the term multiplying dt in the equation above must be zero. Defining the operator

$$\mathcal{A} = \sum_{i=1}^p \mu_i(t, x) \frac{\partial}{\partial x_i} + \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p \Sigma_{i,j}(t, x) \frac{\partial^2}{\partial x_i \partial x_j},$$

we deduce that $u(t, x)$ satisfies the PDE

$$\frac{\partial u(t, x)}{\partial t} + \mathcal{A}u(t, x) = 0, \quad (1.30)$$

with terminal condition $u(T, x) = g(x)$. The equation above is known as the *Kolmogorov backward equation* for the SDE (1.29). The operator \mathcal{A} is known as the *generator* or *infinitesimal operator* of the SDE, and can be identified as

$$\mathcal{A}u(t, x) = \lim_{h \downarrow 0} \frac{\mathbb{E}^P(u(t+h, X(t+h)) | X(t) = x) - u(t, x)}{h}.$$

In arriving at (1.30) we made several implicit assumptions, most notably that the function $u(t, x)$ exists and is twice differentiable. Sufficient conditions for the validity of (1.30) can be found in Karatzas and Shreve [1997], for instance. A relevant result is listed below.

Theorem 1.8.1. *Let the process $X(t)$ be given by the SDE (1.29), where the coefficients μ and σ are continuous in x and satisfy the Lipschitz and growth conditions of Theorem 1.6.1. Consider a continuous function $g(x)$ that is either non-negative or satisfies a polynomial growth condition, meaning that for some positive constants K and q*

$$g(x) \leq K(1 + |x|^q), \quad x \in \mathbb{R}^p.$$

If $u(t, x)$ solves (1.30) with boundary condition $u(T, x) = g(x)$, and $u(t, x)$ satisfies a polynomial growth condition in x , then

$$u(t, x) = \mathbb{E}^P(g(X(T)) | X(t) = x), \quad t \in [0, T]. \quad (1.31)$$

Conditions required to ensure existence of a solution to (1.30) are more involved, and we just refer to Karatzas and Shreve [1997] and the references therein.

A family of functions g of particular importance to many of our applications is

$$g(x) = e^{ik^T x}, \quad k \in \mathbb{R}^p,$$

where $i = \sqrt{-1}$ is the imaginary unit. In this case $u(t, x)$ becomes the *characteristic function* of $X(T)$, conditional on $X(t) = x$. We refer to any standard statistics textbook (e.g. Ochi [1990]) for the many useful properties of characteristic functions.

For the Markov process $X(t)$ in (1.29), let us now introduce a *transition density*, given heuristically by

$$p(t, x; s, y) dy \triangleq \mathbb{P}(X(s) \in [y, y + dy] | X(t) = x), \quad 0 \leq t \leq s \leq T.$$

We can loosely think of the transition density as a special case of the functional $u(t, x)$ above, with boundary condition $u(s, x) = \delta(x - y)$, where $\delta(\cdot)$ is the Dirac delta function. Sometimes $p(\cdot, \cdot; \cdot, \cdot)$ is called a *Green's function* or a *fundamental solution* to (1.30). Under certain regularity conditions discussed in Karatzas and Shreve [1997], the transition density solves the Kolmogorov backward equation

$$\frac{\partial p(t, x)}{\partial t} + \mathcal{A}p(t, x) = 0, \quad (s, y) \text{ fixed},$$

subject to the boundary condition $p(s, x; s, y) = \delta(x - y)$. Further, the general expectation $u(t, x) = \mathbb{E}^P(g(X(T)) | X(t) = x)$ in Theorem 1.8.1 can be written

$$u(t, x) = \int_{\mathbb{R}^p} g(y)p(t, x; T, y) dy, \quad t \in [0, T]. \quad (1.32)$$

In many applications, it is useful to have a result that produces transition densities at future times $s \geq t$ from a known state at time t , rather than vice-versa. For this, we first define an operator \mathcal{A}^* by

$$\mathcal{A}^* f(s, y) = - \sum_{i=1}^p \frac{\partial [\mu_i(s, y) f(s, y)]}{\partial y_i} + \frac{1}{2} \sum_{i=1}^p \sum_{j=1}^p \frac{\partial^2 [\Sigma_{i,j}(s, y) f(s, y)]}{\partial y_i \partial y_j}.$$

In the transition density $p(t, x; s, y)$ now consider (t, x) fixed and let \mathcal{A}^* operate on the resulting function of s and y . Under additional regularity conditions, we then have the *forward Kolmogorov equation*

$$-\frac{\partial p(s, y)}{\partial s} + \mathcal{A}^* p(s, y) = 0, \quad (t, x) \text{ fixed}, \quad (1.33)$$

subject to the boundary condition $p(t, x; t, y) = \delta(x - y)$.

The forward Kolmogorov equation is sometimes known as the *Fokker-Planck* equation. We stress that the backward equation is more general than the forward equation, in the sense that the former holds for general terminal conditions $g(x)$, whereas the latter only holds for δ -type initial conditions.

We round off this section by a useful extension to the Kolmogorov backward equation. Specifically, consider extending the PDE (1.30) to

$$\frac{\partial u(t, x)}{\partial t} + \mathcal{A}u(t, x) + h(t, x) = r(t, x)u(t, x), \quad (1.34)$$

where $h, r : [0, T] \times \mathbb{R}^p \rightarrow \mathbb{R}$. Given the boundary condition $u(T, x) = g(x)$, the *Feynman-Kac solution* to (1.34), should it exist, is given by

$$u(t, x) = \mathbb{E}^P \left(\psi(t, T)g(X(T)) + \int_t^T \psi(s, X(s))h(s, X(s)) ds \middle| X(t) = x \right), \quad (1.35)$$

where

$$\psi(t, T) = \exp \left(- \int_t^T r(s, X(s)) ds \right), \quad t \in [0, T].$$

The result is easily understood from an application of Ito's lemma, similar to the one used above to motivate the backward Kolmogorov equation. Sufficient regularity conditions for the Feynman-Kac result to hold are identical to those of Theorem 1.8.1, supplemented with the requirement that r be nonnegative and continuous in x ; and the requirement that h be continuous in x and either be nonnegative or satisfy a polynomial growth requirement in x . See Duffie [2001] for further details about the often delicate regularity issues surrounding the Feynman-Kac result.

For later use, let us finally note that when $g(x) = \delta(x - y)$ and $h(t, x) = 0$, $u(t, x)$ in (1.35) will equal

$$G(t, x; T, y) \triangleq \mathbb{E}^P \left(e^{-\int_t^T r(s, X(s)) ds} \delta(X(T) - y) | X(t) = x \right).$$

The function G is known as a *state-price density* or as an *Arrow-Debreu security price* function. In particular, notice that for an arbitrary $g(x)$, we then have

$$\mathbb{E}^P \left(e^{-\int_t^T r(s, X(s)) ds} g(X(T)) | X(t) = x \right) = \int_{\mathbb{R}} G(t, x; T, y) g(y) dy. \quad (1.36)$$

Comparison with (1.32) shows that the state-price density is, essentially, equivalent to a Green's function with built-in discounting.

1.9 Black-Scholes and Extensions

In reviews of asset pricing theory, a discussion of the seminal *Black-Scholes-Merton* model (sometimes just known as the *Black-Scholes* model) of Black and Scholes [1973] and Merton [1973] is nearly mandatory. As the Black-Scholes-Merton (BSM) model constitutes a well-behaved setting in which to tie elements of previous sections together, our text is no exception. To provide a smoother transition to material that follows, we do, however, extend the usual analysis to include a simple case of stochastic interest rates.

1.9.1 Basics

In the basic BSM economy, two assets are traded: a money market account β and a stock S . In previous notations, $X(t) = (\beta(t), S(t))^\top$ and $p = 2$. The money market account value is 1 at time 0 and accrues risk-free interest at a continuously compounded, non-negative rate of r , initially assumed constant. The dynamics for β are thus given by an ordinary differential equation (ODE)

$$d\beta(t)/\beta(t) = r dt, \quad \beta(0) = 1,$$

implying that simply $\beta(t) = \beta(0)e^{rt}$.

The stock dynamics are assumed to satisfy GBMD under measure P:

$$dS(t)/S(t) = \mu dt + \sigma dW(t), \quad (1.37)$$

where W is a Brownian motion of dimension $d = 1$, and μ and σ are constants.

Taking first a probabilistic approach, we notice that β is positive and can be used as a numeraire. Let $S^\beta(t) = S(t)/\beta(t)$ be the stock price deflated by β . By Ito's lemma,

$$dS^\beta(t)/S^\beta(t) = (\mu - r) dt + \sigma dW(t).$$

Applying Girsanov's theorem (see Theorem 1.5.1) and Corollary 1.5.2, we see that if $\sigma \neq 0$, β will induce a unique equivalent martingale measure, with the measure shift characterized by the density process¹²

$$d\varsigma(t)/\varsigma(t) = -\theta dW(t), \quad \theta = \frac{\mu - r}{\sigma}.$$

Clearly, $\varsigma(t)$ defines an exponential martingale. The probability measure induced by the money market account β is called the *risk-neutral martingale measure* and is traditionally denoted Q. Under Q, $W^\beta(t) = W(t) + \theta t$ is a Brownian motion, and

$$\begin{aligned} dS^\beta(t)/S^\beta(t) &= \sigma dW^\beta(t), \\ dS(t)/S(t) &= r dt + \sigma dW^\beta(t), \end{aligned} \quad (1.38)$$

or, from (1.21),

$$S(T) = S(t)e^{(r - \frac{1}{2}\sigma^2)(T-t) + \sigma(W^\beta(T) - W^\beta(t))}, \quad t \in [0, T]. \quad (1.39)$$

We note that under Q, the drift μ of the stock process is replaced by the risk-free interest rate r . That is, under Q agents in the economy will

¹²The reader may recognize the market price of risk θ as the *Sharpe ratio* of the stock S , a measure of how well the risk of stock (represented by σ) is compensated by excess return (represented by $\mu - r$).

appear to be indifferent (“neutral”) to the risk of the stock, content with an average growth rate of the stock equal to that of the money market account.

Before proceeding with the BSM analysis, we wish to emphasize that the drift restriction imposed on the stock in the risk-neutral measure Q is a general result. In a larger setting with a p -dimensional vector asset process X , if the Q -dynamics of the components of X are all of the form

$$dX_i(t) = rX_i dt + O(dW(t)), \quad i = 1, \dots, p,$$

there is no arbitrage. This result holds unchanged if the interest rate is random (see Section 1.9.3).

Returning to the BSM setting, we note that the risk-neutral measure is unique, whereby the market is complete and all derivative securities on S (and β) are attainable. Let us consider a few such securities. First, we consider a security paying at time T \$1 for certain. Such a security is a *discount bond* and we shall denote its time t price by $P(t, T)$, $t \in [0, T]$. If the interest rate is positive, we would expect $P(t, T) \leq 1$ as a reflection of the time value of money, with equality only holding for $t = T$. Application of the basic derivative pricing equation (1.15) immediately gives

$$P(t, T) = \beta(t) E_t^Q \left(\frac{1}{\beta(T)} \right) = E_t^Q \left(e^{-r(T-t)} \right) = e^{-r(T-t)}.$$

This result is trivial, as it is easily seen that the amount $e^{-r(T-t)}$ invested in the money market account at time t will grow to exactly \$1 at time T .

Second, consider a derivative V paying $V(T) = S(T) - K$ at time T , with K being an arbitrary constant. Proceeding as above, at time $t \leq T$ the arbitrage-free price must be

$$\begin{aligned} V(t) &= E_t^Q \left(e^{-r(T-t)} (S(T) - K) \right) \\ &= e^{-r(T-t)} \left(E_t^Q (S(T)) - K \right) = S(t) - K P(t, T), \end{aligned} \quad (1.40)$$

where the last equality follows from property (1.22) of GBMD. We notice that $V(t) = 0$ if $K = S(t)/P(t, T)$. This value of K is known as the time t *forward price of $S(T)$* ¹³.

Third, consider the derivative that was the main focus of the original BSM analysis, a *European call option* paying¹⁴ $c(T) = (S(T) - K)^+$, with K being a positive *strike price*. Following (1.40), we can write

$$c(t) = P(t, T) E_t^Q \left((S(T) - K)^+ \right). \quad (1.41)$$

From the representation (1.39), basic probability theory allows us to write this expectation as

¹³We shall touch on the closely related concept of a *futures price* in Section 4.1.2.

¹⁴We use the notations $x^+ = \max(x, 0)$, $x^- = \min(x, 0)$ throughout this book.

$$c(t) = P(t, T) \int_{-\infty}^{\infty} \left(S(t) e^{(r - \frac{1}{2}\sigma^2)(T-t) + z\sigma\sqrt{T-t}} - K \right)^+ \phi(z) dz, \quad (1.42)$$

where $\phi(z) = (2\pi)^{-1/2} \exp(-z^2/2)$ is the standard Gaussian density. A straightforward evaluation of the integral leads to the famous *Black-Scholes-Merton call pricing formula*:

Theorem 1.9.1. *In the BSM economy, the arbitrage-free time t price of a K -strike call option maturing at time T is*

$$c(t) = S(t)\Phi(d_+) - K P(t, T)\Phi(d_-), \quad (1.43)$$

$$d_{\pm} \triangleq \frac{\ln(S(t)/K) + (r \pm \sigma^2/2)(T-t)}{\sigma\sqrt{T-t}}, \quad t < T,$$

where $\Phi(\cdot)$ is the Gaussian cumulative distribution function.

A formula for a *European put option* $p(t)$ paying $(K - S(T))^+$ can be obtained from (1.43) by *put-call parity*:

$$c(t) - p(t) = V(t),$$

where $V(t)$ is the forward contract defined above.

Remark 1.9.2. At time t , call and put options with strikes equal to $S(t)$ are said to be *at-the-money* (ATM). If $S(t) > K$, the call option is *in-the-money* (ITM) and the put option is *out-of-the-money* (OTM). If $S(t) < K$, the call is OTM and the put is ITM. The ATM, ITM, and OTM monikers are sometimes used to refer to the ordering of the *forward value* $E_t(S(T)) = S(t)e^{r(T-t)}$ (for a T -maturity option) rather than the spot $S(t)$, relative to the strike K .

In deriving (1.43), the choice of β as numeraire was arbitrary. If we instead use S (which is also strictly positive) as numeraire, we can write

$$c(t) = S(t)E_t^{Q^S} \left(\frac{(S(T) - K)^+}{S(T)} \right) = S(t)E_t^{Q^S} \left((1 - K/S(T))^+ \right), \quad (1.44)$$

where Q^S is the martingale measure induced by S . To identify the measure shift involved in moving from P to Q^S , consider that $\beta^S(t) = \beta(t)/S(t)$ must be a martingale in Q^S . By Ito's lemma, in measure P we have

$$d\beta^S(t)/\beta^S(t) = (r - \mu + \sigma^2) dt - \sigma dW(t),$$

such that $dW^S(t) = ((r - \mu)/\sigma + \sigma) dt + dW(t)$ is a Brownian motion under Q^S . Application of Ito's lemma on $1/S(t)$ yields, after a few rearrangements,

$$dS(t)^{-1}/S(t)^{-1} = -rdt - \sigma dW^S(t),$$

which is a GBMD as before. Evaluation of the expectation (1.44) can be verified to recover the BSM formula (1.43).

Our derivation of the BSM formula was so far entirely probabilistic. Writing $c(t) = c(t, \beta, S)$, the arguments in Section 1.7 allow us to write c as a solution to the PDE (see (1.27))

$$\frac{\partial c}{\partial t} + \frac{1}{2}\sigma^2 S^2 \frac{\partial^2 c}{\partial S^2} = 0, \quad (1.45)$$

subject to the boundary condition $c(T, \beta, S) = (S - K)^+$. From (1.28) we also have that the replication positions in β and S are $\frac{\partial c}{\partial \beta}$ and $\frac{\partial c}{\partial S}$, respectively. That is,

$$c(t, \beta, S) = \frac{\partial c}{\partial \beta} \beta + \frac{\partial c}{\partial S} S. \quad (1.46)$$

As β is deterministic, we can actually eliminate c -dependence on this variable by a change of variables $\tilde{c}(t, S) = c(t, \beta, S)$. By the chain rule

$$\frac{\partial \tilde{c}}{\partial t} = \frac{\partial c}{\partial t} + \frac{\partial c}{\partial \beta} \frac{\partial \beta}{\partial t} = \frac{\partial c}{\partial t} + \frac{\partial c}{\partial \beta} r\beta = \frac{\partial c}{\partial t} + rc - \frac{\partial c}{\partial S} rS$$

where the last equation follows from (1.46). Inserting this into (1.45) yields the original *Black-Scholes PDE*

$$\frac{\partial \tilde{c}}{\partial t} + rS \frac{\partial \tilde{c}}{\partial S} + \frac{1}{2}\sigma^2 S^2 \frac{\partial^2 \tilde{c}}{\partial S^2} = r\tilde{c}, \quad (1.47)$$

with $\tilde{c}(T, S) = (S - K)^+$. We can solve this equation by classical methods (see Lipton [2001] for several techniques), or we can use the Feynman-Kac result to write it as an expectation. We leave it as an exercise to the reader to verify that Feynman-Kac leads to the same expectation as derived earlier by probabilistic means (see (1.41)).

A final note: the derivation of the Black-Scholes PDE above was somewhat non-standard due to the initial assumption of option price being a function of the deterministic numeraire β . A more conventional (but entirely equivalent) argument sets up a portfolio of the call option and a position in the stock, and demonstrates that the stock position can be set such that the total portfolio growth is deterministic (risk-free) on $[t, t+dt]$. Equating the portfolio growth with the risk-free rate yields the Black-Scholes PDE (1.47). See Hull [2006] for details of this approach.

1.9.2 Alternative Derivation

We have already demonstrated several different ways of proving the BSM call pricing formula, but as shown in Andreasen et al. [1998] there are many more. One particularly enlightening proof is based on the concept of *local time* and shall briefly be discussed in this section. The proof, which borrows

from the results in Carr and Jarrow [1990], will also allow us to demonstrate the Tanaka extension of Ito's lemma, mentioned earlier in Section 1.1.

As above, we assume that the stock price process is as in (1.38), and define the forward stock price $F(t) \triangleq S(t)/P(t, T)$. Clearly,

$$dF(t)/F(t) = \sigma dW^\beta(t), \quad t \leq T, \quad (1.48)$$

where W^β is a Brownian motion in the risk-neutral measure. Define the random variable $I(t) = (F(t) - K)^+$. The first derivative of I with respect to F is an indicator function $1_{\{F(t)>K\}}$ and the second derivative can be interpreted as the Dirac delta function, $\delta(F(t) - K)$. As I is clearly not twice differentiable, Ito's lemma formally does not apply, but the Tanaka extension nevertheless gives us permission to write

$$\begin{aligned} dI(t) &= 1_{\{F(t)>K\}} dF(t) + \frac{1}{2}\sigma^2 F(t)^2 \delta(F(t) - K) dt \\ &= 1_{\{F(t)>K\}} \sigma F(t) dW^\beta(t) + \frac{1}{2}\sigma^2 K^2 \delta(F(t) - K) dt. \end{aligned}$$

In integrated form,

$$I(T) = I(t) + \int_t^T 1_{\{F(u)>K\}} \sigma F(u) dW^\beta(u) + \frac{1}{2}\sigma^2 K^2 \int_t^T \delta(F(u) - K) du.$$

The second integral in this expression is a random variable known as the *local time of F spent at the level K* , on the interval $[t, T]$. Taking expectations, it follows that

$$E_t^Q(I(T)) = I(t) + \frac{1}{2}\sigma^2 K^2 \int_t^T E_t^Q(\delta(F(u) - K)) du.$$

Here, if $p(t, y; u, x)$ is the density of $F(u)$ given $F(t) = y$, $u \geq t$, then obviously

$$E_t^Q(\delta(F(u) - K)) = p(t, F(t); u, K).$$

By the definition of $F(T)$ we have $F(T) = S(T)$, such that $I(T) = (S(T) - K)^+$. From (1.41), we may therefore write the time t European call option price as

$$\begin{aligned} c(t) &= P(t, T) E_t^Q(I(T)) \\ &= (S(t) - KP(t, T))^+ + \frac{P(t, T)}{2}\sigma^2 K^2 \int_t^T p(t, F(t); u, K) du. \quad (1.49) \end{aligned}$$

The formula (1.49) decomposes the call option into a sum of two terms, the *intrinsic value* and the *time value*, respectively. The time value can be made more explicit by observing from the representation (1.39) that¹⁵

¹⁵This also follows directly from the fact that $F(u)$ is a log-normal random variable with moments given by (1.22) and (1.23).

$$p(t, F(t); u, K) = \frac{1}{K\sigma\sqrt{u-t}\sqrt{2\pi}} \exp\left(-\frac{1}{2}d(u)^2\right),$$

$$d(u) \triangleq \frac{\ln(F(t)/K) - \frac{1}{2}\sigma^2(u-t)}{\sigma\sqrt{u-t}}.$$

In other words, we have arrived at the following result.

Proposition 1.9.3. *The European call option price $c(t)$ on the process (1.38) can be written as*

$$c(t) = (S(t) - KP(t, T))^+ + \frac{P(t, T)\sigma K}{2} \int_t^T \frac{\phi(d(u))}{\sqrt{u-t}} du, \quad (1.50)$$

where $\phi(x)$ is the Gaussian density.

Explicit evaluation of the integral in (1.50) can be verified to produce the BSM formula in Theorem 1.9.1. We leave this as an exercise to the reader.

1.9.3 Extensions

1.9.3.1 Deterministic Parameters and Dividends

In our basic BSM setup, consider now first a simple extension to a deterministic interest rate $r(t)$ and a deterministic volatility $\sigma(t)$. Carrying out the analysis as before, we see that discount bond prices now become

$$P(t, T) = e^{-\int_t^T r(s) ds}. \quad (1.51)$$

The BSM call pricing formula (1.43) holds unchanged provided $P(t, T)$ is changed according to (1.51), and we redefine

$$d_{\pm} \triangleq \frac{\ln(S(t)/K) + \int_t^T (r(s) \pm \sigma(s)^2/2) ds}{\sqrt{\int_t^T \sigma(s)^2 ds}}.$$

Let us further assume that the stock pays dividends at a deterministic rate of $q(t)$. Our framework so far, however, has assumed that assets pay no cash over $[0, T]$. To salvage the situation, consider a fictitious asset S^* obtained by reinvesting all dividends into the stock S itself. It is easily seen that

$$S^*(t) = S(t)e^{\int_0^t q(s) ds},$$

and clearly $S^*(t)$ satisfies the requirements of generating no cash flows on $[0, T]$. Stating the call option payout as

$$c(T) = (S(T) - K)^+ = \left(S^*(T)e^{-\int_0^T q(s) ds} - K \right)^+$$

and performing the pricing analysis of Section 1.9.1 on $S^*(t)$, rather than $S(t)$, results in a dividend-extended BSM call option formula:

$$c(t) = S(t)e^{-\int_t^T q(s) ds} \Phi(d_+) - KP(t, T)\Phi(d_-),$$

$$d_{\pm} \triangleq \frac{\ln(S(t)/K) + \int_t^T (r(s) - q(s)) \pm \sigma(s)^2/2 ds}{\sqrt{\int_t^T \sigma(s)^2 ds}}.$$

When the stock pays a dividend rate of $q(t)$, note that the risk-neutral process for $S(t)$ is

$$dS(t)/S(t) = (r(t) - q(t)) dt + \sigma(t) dW^{\beta}(t),$$

which extends (1.38). Note that for the special case where $r(t) = q(t)$, $S(t)$ becomes a martingale and the call option price formula simplifies to

$$c(t) = P(t, T)(S(t)\Phi(d_+) - K\Phi(d_-)), \quad (1.52)$$

where now

$$d_{\pm} \triangleq \frac{\ln(S(t)/K) \pm \frac{1}{2} \int_t^T \sigma(s)^2 ds}{\sqrt{\int_t^T \sigma(s)^2 ds}}.$$

Remark 1.9.4. The martingale call formula (1.52) typically emerges when pricing options on futures and forward prices (see (1.48)) and is often called the *Black formula*, in honor of the work in Black [1976].

1.9.3.2 Stochastic Interest Rates

We now get even more ambitious and wish to consider call option pricing in the case where the interest rate r is stochastic. The money market account β becomes

$$\beta(t) = e^{\int_0^t r(s) ds},$$

and is now assumed an \mathcal{F}_t -measurable random variable. Proceeding as in Section 1.9.1, we find that under the risk-neutral measure Q , the call option price expression is (assuming that the stock pays no dividends)

$$c(t) = \beta(t) \mathbb{E}_t^Q \left(\frac{1}{\beta(T)} (S(T) - K)^+ \right) = \mathbb{E}_t^Q \left(e^{-\int_t^T r(s) ds} (S(T) - K)^+ \right). \quad (1.53)$$

In (1.53), we emphasize that the numeraire no longer can be pulled out from the expectation. Still, to simplify call option computations, it would be convenient to somehow remove the term $\exp(-\int_t^T r(s) ds)$ from the expectation in (1.53). By substituting 1 for $(S(T) - K)^+$ in the expression above, we first notice that

$$P(t, T) = \mathbb{E}_t^Q \left(e^{-\int_t^T r(s) ds} \right).$$

This inspires us to perform a new measure shift, where we use the discount bond $P(t, T)$, rather than $\beta(t)$, as our numeraire. Let the martingale measure induced by $P(t, T)$ be denoted Q^T , often termed the *T-forward measure*. By the standard result (1.15) we have

$$\begin{aligned} c(t) &= P(t, T) \mathbb{E}_t^{Q^T} \left(P(T, T)^{-1} (S(T) - K)^+ \right) \\ &= P(t, T) \mathbb{E}_t^{Q^T} \left((S(T) - K)^+ \right), \end{aligned}$$

where we have used that $P(T, T) = 1$. From Theorem 1.4.2, Q^T and Q are related by the density

$$\varsigma(t) = \mathbb{E}_t^Q \left(\frac{dQ^T}{dQ} \right) = \frac{P(t, T)/P(0, T)}{\beta(t)}. \quad (1.54)$$

To proceed, we need to add more structure to the model by making assumptions about the stochastic process for $P(t, T)$. We shall spend considerable effort in subsequent chapters on this issue, but for this initial application we simply assume that $P(t, T)$ has Q dynamics

$$dP(t, T)/P(t, T) = r(t) dt - \sigma_P(t, T) dW_P(t), \quad (1.55)$$

where $\sigma_P(t, T)$ is deterministic and $W_P(t)$ is a Brownian motion correlated to the stock Brownian motion. Notice that the drift of $P(t, T)$ under Q is not freely specifiable and must be equal to the risk-free rate; see the discussion following (1.38). For clarity, let the stock Brownian motion be renamed $W_S(t)$, and assume that the correlation between $W_P(t)$ and $W_S(t)$ is a constant ρ . In the setting of vector-valued Brownian motion with independent components used in earlier sections, we can introduce correlation by writing $W(t) = (W_1(t), W_2(t))^\top$ and setting, say,

$$\begin{aligned} W_P(t) &= W_1(t), \\ W_S(t) &= \rho W_1(t) + \sqrt{1 - \rho^2} W_2(t). \end{aligned}$$

The filtration $\{\mathcal{F}_t\}$ of our extended BSM setting is the one generated by the 2-dimensional $W(t)$.

Under Q^T , the deflated process $S^P(t) = S(t)/P(t, T)$ is a martingale. An application of Ito's lemma combined with the Diffusion Invariance Principle shows that the Q^T process for $S^P(t)$ is

$$dS^P(t)/S^P(t) = \sigma_P(t, T) dW_1(t) + \sigma(t) \left(\rho dW_1(t) + \sqrt{1 - \rho^2} dW_2(t) \right), \quad (1.56)$$

where $\sigma(t)$ as before is the deterministic volatility of the stock S . We recognize $S^P(t)$ as a drift-free geometric Brownian motion with instantaneous variance of

$$\sigma_P(t, T)^2 + \sigma(t)^2 + 2\rho\sigma(t)\sigma_P(t, T).$$

Exploiting the convenient fact that $S^P(T) = S(T)$ and $c(T) = (S^P(T) - K)^+$ (as $P(T, T) = 1$), we get

$$\begin{aligned} c(t) &= P(t, T) E_t^{Q^T} \left((S^P(T) - K)^+ \right) \\ &= P(t, T) \int_{-\infty}^{\infty} \left(S^P(t) e^{-\frac{1}{2}v(t, T) + z\sqrt{v(t, T)}} - K \right)^+ \phi(z) dz, \end{aligned} \quad (1.57)$$

where we have defined the “term”, or total, variance

$$v(t, T) \triangleq \int_t^T (\sigma_P(s, T)^2 + \sigma(s)^2 + 2\rho\sigma(s)\sigma_P(s, T)) ds. \quad (1.58)$$

Completing the integration (compare with (1.42)) and using $S^P(t) = S(t)/P(t, T)$, we arrive at a modified BSM-type call option formula:

Proposition 1.9.5. *Consider a BSM economy with stochastic interest rates evolving according to (1.55). Define term variance $v(t, T)$ as in (1.58). Then, the T -maturity European call option price is*

$$\begin{aligned} c(t) &= S(t)\Phi(d_+) - KP(t, T)\Phi(d_-), \\ d_{\pm} &= \frac{\ln(S(t)/(KP(t, T))) \pm \frac{1}{2}v(t, T)}{\sqrt{v(t, T)}}. \end{aligned}$$

Proposition 1.9.5 was originally derived in Merton [1973], using PDE methods. Extensions to dividend-paying stocks are straightforward and follow the arguments shown in Section 1.9.3.1.

1.10 Options with Early Exercise Rights

In our previous definition of a contingent claim, we assumed that the claim involved a single \mathcal{F}_T -measurable payout at time T . In reality, a number of derivative contracts may have intermediate cash payments from, say, scheduled coupons or through “rebates” for barrier-style options. Mostly, such complications are straightforwardly incorporated; see for instance Section 2.7.3. Of particular interest from a theoretical perspective are the claims that allow the holder to accelerate payments through *early exercise*. Derivative securities with early exercise are characterized by an adapted payout process $U(t)$, payable to the option holder at a stopping time (or *exercise policy*) $\tau \leq T$, chosen by the holder. If early exercise can take place at any time in some interval, we say that the derivative security is an *American option*; if exercise can only take place on a discrete set of dates, we say that it is a *Bermudan option*.

Let the allowed (and deterministic) set of exercise dates larger than or equal to t be denoted $\mathcal{D}(t)$, and suppose that we are given at time 0 a particular exercise policy τ taking values in $\mathcal{D}(0)$, as well as a pricing numeraire N inducing a unique martingale measure Q^N . Let $V^\tau(0)$ be the time 0 value of a derivative security that pays $U(\tau)$. Under some technical conditions on $U(t)$, we can write for the value of the derivative security

$$V^\tau(0) = E^{Q^N} \left(\frac{U(\tau)}{N(\tau)} \right), \quad (1.59)$$

where we have assumed, with no loss of generality, that $N(0) = 1$. Let $\mathcal{T}(t)$ be the time t set of (future) stopping times taking value in $\mathcal{D}(t)$. In the absence of arbitrage, the time 0 value of a security with early exercise into U must then be given by the *optimal stopping problem*

$$V(0) = \sup_{\tau \in \mathcal{T}(0)} V^\tau(0) = \sup_{\tau \in \mathcal{T}(0)} E^{Q^N} \left(\frac{U(\tau)}{N(\tau)} \right), \quad (1.60)$$

reflecting the fact that a rational investor would choose an exercise policy to optimize the value of his claim.

We can extend (1.60) to future times t by

$$V(t) = N(t) \sup_{\tau \in \mathcal{T}(t)} E_t^{Q^N} \left(\frac{U(\tau)}{N(\tau)} \right), \quad (1.61)$$

where $\sup_{\tau \in \mathcal{T}(t)} E_t^{Q^N} (U(\tau)/N(\tau))$ is known as the *Snell envelope* of U/N under Q^N . The process $V(t)$ must here be interpreted as the value of the option with early exercise, *conditional* on exercise not having taken place before time t . To make this explicit, let $\tau^* \in \mathcal{T}(0)$ be the optimal exercise policy, as seen from time 0. We can then write, for $0 < t \leq T$,

$$V(0) = E^{Q^N} (1_{\{\tau^* \geq t\}} V(t)/N(t)) + E^{Q^N} (1_{\{\tau^* < t\}} U(\tau^*)/N(\tau^*)), \quad (1.62)$$

where we break the time 0 value into two components: one from the time t value of the option, should it not have been exercised before time t ; and one from the right to exercise on $[0, t]$. As we can always elect — possibly suboptimally — to never exercise on $[0, t]$, from (1.62) we see that

$$V(0) \geq E^{Q^N} (V(t)/N(t)),$$

which establishes that $V(t)/N(t)$ is a *supermartingale* under Q^N . This result also follows directly from known properties of the Snell envelope; see, e.g., Musiela and Rutkowski [1997].

For later use, focus now on the Bermudan case and assume that $\mathcal{D}(0) = \{T_1, T_2, \dots, T_B\}$, where $T_1 > 0$ and $T_B = T$. For $t < T_{i+1}$, define $H_i(t)$ as the time t value of the Bermudan option when exercise is restricted to the dates $\mathcal{D}(T_{i+1}) = \{T_{i+1}, T_{i+2}, \dots, T_B\}$. That is

$$H_i(t) = N(t) \mathbb{E}_t^{\mathbb{Q}^N} (V(T_{i+1})/N(T_{i+1})) , \quad i = 1, \dots, B-1.$$

At time T_i , $H_i(T_i)$ can be interpreted as the *hold value* of the Bermudan option, that is, the value of the Bermudan option if not exercised at time T_i . If an optimal exercise policy is followed, clearly we must have at time T_i

$$V(T_i) = \max (U(T_i), H_i(T_i)) , \quad i = 1, \dots, B, \quad (1.63)$$

such that

$$H_i(t) = N(t) \mathbb{E}_t^{\mathbb{Q}^N} (\max (U(T_{i+1}), H_{i+1}(T_{i+1})) / N(T_{i+1})) , \quad i = 1, \dots, B-1. \quad (1.64)$$

Starting with the terminal condition $H_B(T) = 0$, (1.64) defines a useful iteration backwards in time for the value $V(0) = H_0(0)$. We shall use this later for the purposes of designing valuation algorithms in Chapter 18, and for computing price sensitivities (deltas) in Chapter 24.

We note that the idea behind (1.63) is often known as *dynamic programming* or the *Bellman principle*. Loosely speaking, we here work “from the back” to price the Bermudan option. As we shall see later (in Chapter 2), this idea is particularly well-suited for numerical methods that proceed backwards in time, such as finite difference methods.

1.10.1 The Markovian Case

We now specialize to the Markovian case where $U(t) = g(t, x(t))$, where $g : [0, T] \times \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous and

$$dx(t) = \mu(t, x(t)) dt + \sigma(t, x(t)) dW(t) \quad (1.65)$$

is an n -dimensional Markovian process, where μ and σ satisfy the regularity conditions of Theorem 1.6.1. The n -dimensional process¹⁶ $x(t)$ here defines the state of the exercise value $U(t)$, so we say that $x(t)$ is a *state variable process*. For concreteness let our numeraire $N(t)$ be the money market account

$$N(t) = \beta(t) = e^{\int_0^t r(u, x(u)) du},$$

where the short interest rate $r : [0, T] \times \mathbb{R}^n \rightarrow \mathbb{R}$ is here assumed a function of time and the state variable vector x . In (1.65), $W(t)$ is understood to be a d -dimensional Brownian motion in the risk-neutral measure \mathbb{Q} .

Writing $V(t) = V(t, x(t))$, we have from (1.61)

$$V(t, x) = \sup_{\tau \in \mathcal{T}(t)} \mathbb{E}^{\mathbb{Q}} \left(e^{-\int_t^\tau r(u, x(u)) du} g(\tau, x(\tau)) \middle| x(t) = x \right). \quad (1.66)$$

¹⁶Note that $x(t)$ is an abstract construct, and does not necessarily coincide with any asset price process.

For dates $t \in \mathcal{D}(0)$, clearly $V(t, x) \geq g(t, x)$, with equality holding only when time t exercise is optimal. This leads us to define the concept of an *exercise region* as

$$\mathcal{X} = \{(t, x) \in \mathcal{D}(0) \times \mathbb{R}^n : V(t, x) = g(t, x)\}.$$

Similarly, we define the complement of \mathcal{X} ,

$$\mathcal{C} = \{(t, x) \in [0, T] \times \mathbb{R}^n : (t, x) \notin \mathcal{X}\},$$

to be the *continuation region*, i.e. the region where we wait (either because exercise is not optimal or because it is not allowed, $t \notin \mathcal{D}(0)$) rather than exercise the option.

For Markovian systems, rather than solving the optimization problem (1.66) directly, it is often particularly convenient to invoke the Bellman principle. Extending the ideas presented earlier, let us, somewhat loosely, state the Bellman principle as follows: for any $t \in \mathcal{D}(0)$,

$$V(t, x) = \lim_{\Delta \downarrow 0} \max \left(g(t, x), \mathbb{E}_t^Q \left(e^{-\int_t^{t+\Delta} r(u, x(u)) du} V(t + \Delta, x(t + \Delta)) \right) \right). \quad (1.67)$$

Again, this simply says that the option value at time t is the maximum of the exercise value and the hold value, that is, the present value of continuing to hold on to the option for a small period of time. As we have seen above, for a Bermudan option, (1.67) also holds for finite Δ (namely up to the next exercise date).

The Bellman principle provides us with a link between present (time t) and future (time $t + \Delta$) option values that we can often exploit in a numerical scheme. For this, however, we need further characterization of $V(t, x)$ in the continuation region. By earlier arguments, we realize that $V(t, x)/\beta(t)$ must be a Q -martingale on the continuation region. Assuming sufficient smoothness for an application of Ito's lemma, this leads to a PDE formulation, to hold for $(t, x) \in \mathcal{C}$,

$$\mathcal{J}V(t, x) = 0, \quad (1.68)$$

where

$$\mathcal{J} = \frac{\partial}{\partial t} + \mu(t, x) \frac{\partial}{\partial x} + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\sigma(t, x) \sigma(t, x)^\top)_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} - r(t, x).$$

Assume first that our option is of the Bermudan type, and let T_i and T_{i+1} be subsequent exercise dates in the exercise schedule. For any function f of time, define $f(t \pm)$ to be the limits $\lim_{\epsilon \downarrow 0} f(t \pm \epsilon)$, and assume that $V(T_{i+1}-, x)$ is known for all x . As all values of $t \in (T_i, T_{i+1})$ by definition must be in the continuation region, we can use (1.68) to solve for $V(T_i+, x)$. Applying the Bellman principle (1.67) at time T_i then leads to the condition

$$V(T_i-, x) = \max(g(T_i, x), V(T_i+, x)).$$

In PDE parlance, this is a so-called *jump condition* which is straightforward to incorporate into a numerical solution; see Section 2.7.4 for details.

For American-style options, (1.68) continues to apply on \mathcal{C} . The Bellman principle here leads to the characterization that

$$\mathcal{J}V(t, x) < 0,$$

for $(t, x) \in \mathcal{X}$, i.e. we exercise when the rate of return from holding the option strictly fails to match $r(t, x)$. The American option pricing problem is often conveniently summarized in a *variational inequality*, to hold on $\mathcal{X} \cup \mathcal{C}$,

$$V(t, x) \geq g(t, x), \quad \mathcal{J}V(t, x) \leq 0, \quad (V(t, x) - g(t, x)) \mathcal{J}V(t, x) = 0, \quad (1.69)$$

and subject to the boundary condition $V(T, x) = g(T, x)$. The first of these three conditions expresses that the option is always worth at least its exercise value; the second expresses the supermartingale property of $V(t, x)$; and the third implies (after a little thought) that $\mathcal{J}V(t, x) = 0$ on \mathcal{C} and $\mathcal{J}V(t, x) < 0$ on \mathcal{X} . The system (1.69) is discussed more carefully in Duffie [2001], where additional discussion of regularity issues may also be found.

1.10.2 Some General Bounds

In many cases of practical interest, solving PDEs and/or variational inequalities is not computationally feasible. In such situations, we may be interested in at least bounding the value of an option with early exercise rights. Providing a lower bound is straightforward: postulate an exercise policy τ and compute the price $V^\tau(0)$ by direct methods. From (1.60), clearly this provides a lower bound

$$V^\tau(0) \leq V(0). \quad (1.70)$$

The closer the postulated exercise policy τ is to the optimal exercise policy τ^* , the tighter this bound will be. We shall later study a number of numerical techniques to generate good exercise strategies for fixed income options with early exercise rights, see Chapter 18.

To produce an upper bound, we can rely on duality results established in Rogers [2001], Haugh and Kogan [2004] and Andersen and Broadie [2004]. Let \mathcal{K} denote the space of adapted martingales M for which $\sup_{\tau \in [0, T]} \mathbb{E}^{Q^N} |M(\tau)| < \infty$. For a martingale $M \in \mathcal{K}$, we then write

$$\begin{aligned} V(0) &= \sup_{\tau \in \mathcal{T}(0)} \mathbb{E}^{Q^N} \left(\frac{U(\tau)}{N(\tau)} \right) \\ &= \sup_{\tau \in \mathcal{T}(0)} \mathbb{E}^{Q^N} \left(\frac{U(\tau)}{N(\tau)} + M(\tau) - M(\tau) \right) \\ &= M(0) + \sup_{\tau \in \mathcal{T}(0)} \mathbb{E}^{Q^N} \left(\frac{U(\tau)}{N(\tau)} - M(\tau) \right). \end{aligned}$$

In the second equality, we have relied on the *optional sampling theorem*, a result that states that the martingale property is satisfied up to a bounded random stopping time, i.e. that $E^{Q^N}(M(\tau)) = M(0)$; see Karatzas and Shreve [1997] for details. We now turn the above result into an upper bound by forming a pathwise maximum at all possible future exercise dates $\mathcal{D}(0)$:

$$\begin{aligned} V(0) &= M(0) + \sup_{\tau \in \mathcal{T}(0)} E^{Q^N} \left(\frac{U(\tau)}{N(\tau)} - M(\tau) \right) \\ &\leq M(0) + E^{Q^N} \left(\max_{t \in \mathcal{D}(0)} \left(\frac{U(t)}{N(t)} - M(t) \right) \right). \end{aligned} \quad (1.71)$$

With (1.70) and (1.71) we have, as desired, established upper and lower bounds for values of options with early exercise rights. Let us consider how to make these bounds tight. As mentioned earlier, to tighten the lower bound we need to pick exercise strategies close to the optimal one. Tightening the upper bound is a bit more involved and requires the following basic theorem, proven in Karatzas and Shreve [1997]:

Theorem 1.10.1 (Doob-Meyer Decomposition). *Let $\{Y(t), t \in [0, T]\}$ be a positive \mathcal{F}_t -adapted supermartingale process with right-continuous sample paths. Then we can write*

$$Y(t) = m(t) - A(t),$$

where $m(t)$ is a martingale process with $m(0) = Y(0)$ and $A(t)$ is an increasing predictable process with $A(0) = 0$.

Applying the Doob-Meyer decomposition on the supermartingale process $V(t)/N(t)$ under Q^N shows that

$$V(t)/N(t) = m(t) - A(t),$$

and $V(0) = m(0)$. Consider taking $M(t) = m(t)$ in equation (1.71), to get

$$\begin{aligned} V(0) &\leq V(0) + E^{Q^N} \left(\max_{t \in \mathcal{D}(0)} \left(\frac{U(t)}{N(t)} - m(t) \right) \right) \\ &= V(0) + E^{Q^N} \left(\max_{t \in \mathcal{D}(0)} \left(\frac{U(t)}{N(t)} - \frac{V(t)}{N(t)} - A(t) \right) \right) \\ &\leq V(0). \end{aligned}$$

The last inequality follows from the fact that $V(t) \geq U(t)$ and $A(t) \geq 0$. In conclusion, we have arrived at a *dual* formulation of the option price

$$V(0) = \inf_{M \in \mathcal{K}} \left\{ M(0) + E^{Q^N} \left(\max_{t \in \mathcal{D}(0)} \left(\frac{U(t)}{N(t)} - M(t) \right) \right) \right\}, \quad (1.72)$$

and have demonstrated that the infimum is attained when the martingale M is set equal to the martingale component of the deflated price process

$V(t)/N(t)$. In practice, we are obviously not privy to $V(t)/N(t)$ (which is a quantity that we are trying to estimate), but we are nevertheless provided with a strategy to make the upper bound (1.71) tight: use a martingale that is “close” to the martingale component of the true deflated option price process. In Chapter 18 we shall demonstrate how to make this strategy operational.

1.10.3 Early Exercise Premia

We finish our discussion of options with early exercise rights by listing some known results for puts and calls, including an interesting decomposition of American and Bermudan option prices into the sum of a European option price and an *early exercise premium*. For convenience, we work in a Markovian setting where the single state variable, denoted $S(t)$, follows one-dimensional GBMD. Specifically, we assume that

$$dS(t)/S(t) = (r - q) dt + \sigma dW^\beta(t), \quad (1.73)$$

with $W^\beta(t)$ being a one-dimensional Brownian motion in the risk-neutral measure, i.e. the measure induced by the money market account $\beta(t) = e^{rt}$. For simplicity we assume that the interest rate r , the dividend yield q , and the volatility σ are all constants; the extension to time-dependent parameters is straightforward.

Let $c(t)$, $C_A(t)$, and $C_B(t)$ be the time t European, American, and Bermudan prices of the call option with terminal maturity T , conditional on no exercise prior to time t . While obviously $c(t) \leq C_B(t) \leq C_A(t)$, in some cases these inequalities are equalities, as the following straightforward lemma shows.

Lemma 1.10.2. *Suppose that $r \geq 0$ and $q \leq 0$ in (1.73). It is then never optimal to exercise a call option early, and*

$$c(t) = C_A(t) = C_B(t).$$

Proof. Notice that, by Jensen’s inequality,

$$\begin{aligned} c(t) &= e^{-r(T-t)} \mathbb{E}_t^Q ((S(T) - K)^+) \\ &\geq e^{-r(T-t)} \left((\mathbb{E}_t^Q (S(T)) - K)^+ \right) = (e^{-q(T-t)} S(t) - e^{-r(T-t)} K)^+. \end{aligned}$$

It is therefore clear that if $r \geq 0$ and $q \leq 0$, then for any value of $T - t$,

$$c(t) \geq (S(t) - K)^+,$$

i.e. the European call option price dominates the exercise value. As the hold value of American and Bermudan options must be at least as large as the European option price, it follows that the option to exercise early is worthless. \square

Remark 1.10.3. For the put option, early exercise is never optimal if $r \leq 0$ and $q \geq 0$. As this situation rarely happens in practice, American put options nearly always trade at a premium to their European counterparts.

Lemma 1.10.2 demonstrates the well-known fact that American or Bermudan call options on stocks that pay no dividends ($q = 0$) should never be exercised early. On the other hand, if the stock does pay dividends, for an American call option there will, at time t , be a critical value of the stock, $S_A(t)$, at which the value of the stream of dividends paid by the stock will compensate for the cost of accelerating the payment of the strike K . In other words, an American option should be exercised at time t , provided that $S(t) \geq S_A(t)$. The deterministic curve $S_A(t)$ is known as the *early exercise boundary* and marks the boundary between the exercise and continuation regions, \mathcal{X} and \mathcal{C} . Writing $C_A(t) = C_A(t, S(t))$, we formally have

$$S_A(t) = \inf \left\{ S : C_A(t, S) = (S - K)^+ \right\}, \quad t \leq T.$$

For a Bermudan option, we may similarly define

$$S_B(t) = \inf \left\{ S : C_B(t, S) = (S - K)^+ \right\}, \quad t \in \mathcal{D}(0),$$

where we recall that $\mathcal{D}(0)$ is the (discrete) set of allowed exercise dates for the Bermudan option.

The following important result characterizes the exercise boundary of American call options.

Proposition 1.10.4. *For the American call option on a stock that follows (1.73), we have*

$$\frac{\partial S_A(t)}{\partial t} \leq 0, \quad t < T, \tag{1.74}$$

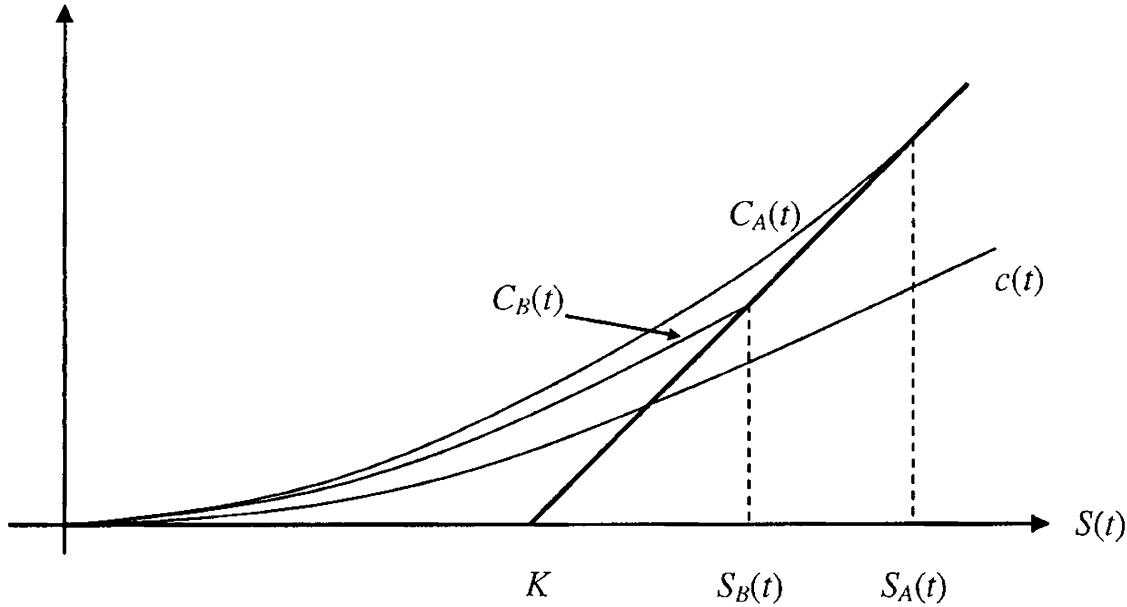
and

$$\left. \frac{\partial C_A(t, S)}{\partial S} \right|_{S=S_A(t)} = 1, \quad t < T. \tag{1.75}$$

Equation (1.74) states that the exercise boundary decreases as we approach maturity, a result that is easily understood. Statement (1.75) is more subtle, however, and amounts to a tangency condition that ensures that the American call option value transitions smoothly from hold value to exercise value across the early exercise boundary. As a consequence, (1.75) is often known as the *smooth pasting condition* or the *high contact condition*. A similar tangency condition does *not* hold for the Bermudan option value, which is not differentiable at the boundary but instead transitions into the exercise region at a “kink”:

$$\lim_{\varepsilon \downarrow 0} \left. \frac{C_B(t, S) - C_B(t, S - \varepsilon)}{\varepsilon} \right|_{S=S_B(t)} < 1, \quad t \in \mathcal{D}(0).$$

Fig. 1.1. Call Option Prices



Notes: Time t prices of American, Bermudan, and European call options, as a function of the asset price. The Bermudan option is assumed to be exercisable at time t .

Figure 1.1 shows a typical value profile for a Bermudan call, along with the corresponding profiles for the European and American options.

Smooth pasting is essentially an optimality condition, which is how Proposition 1.10.4 is traditionally derived (see, e.g., Merton [1973] or the more recent Brekke and Øksendal [1991]). A more descriptive proof based on hedging arguments is given in Tavella and Randall [2000] and Wilmott et al. [1993]. Loosely speaking, the idea is here that a delta hedger should not be able to make riskless profits when the underlying asset crosses into the exercise region. This requires that the delta is continuous across the boundary, which is (1.75).

Remark 1.10.5. For the American put option, $\partial S_A(t)/\partial t \geq 0$ and the high contact condition states that the delta equals -1 at the exercise boundary.

Establishing the boundary $S_A(t)$ will virtually always require numerical methods, although asymptotic results are known for t close to T (see for instance Lipton [2001]). One simple result is listed below.

Lemma 1.10.6. *Assume that $r \geq 0$ and $q \geq 0$, such that the early exercise boundary exists for the American call option. The exercise boundary just prior to maturity is then*

$$\lim_{\varepsilon \downarrow 0} S_A(T - \varepsilon) = K \max \left(1, \frac{r}{q} \right).$$

Proof. An informal proof of Lemma 1.10.6 proceeds as follows. At time $T - dt$, assume that $S(T - dt) > K$; otherwise it clearly makes no sense to exercise the option. If we exercise the option, we receive $S(T - dt) - K$ at time $T - dt$. On the other hand, if we postpone exercise, at time $T - dt$ our hold value is

$$\begin{aligned} e^{-r dt} \mathbb{E}_{T-dt}^Q (S(T) - K) &= S(T - dt) e^{-q dt} - K e^{-r dt} \\ &= S(T - dt) - K - S(T - t) q dt + K r dt. \end{aligned}$$

Clearly, we should then only exercise if

$$S(T - dt) - K > S(T - dt) - K - S(T - t) q dt + K r dt$$

or if

$$S(T - dt) q > K r.$$

□

Notice that since clearly $S_A(T) = K$, the call option exercise boundary will have a *discontinuity* at time T , if $q < r$.

One might guess that complete knowledge of the curve $S_A(t)$ should suffice to price the American option analytically. This intuition is confirmed by the following result due to Jamshidian [1992], Carr et al. [1992], Kim [1990], and Jacka [1991].

Proposition 1.10.7. *The American option price $C_A(t)$ satisfies*

$$C_A(t) = c(t) + E_A(t), \quad t \leq T, \quad (1.76)$$

where the (American) early exercise premium $E_A(t)$ is defined as

$$E_A(t) = \int_t^T e^{-r(u-t)} \mathbb{E}_t^Q (1_{\{S(u) \geq S_A(u)\}} (qS(u) - rK)) du \quad (1.77)$$

$$= \int_t^T \left(qS(t) e^{-q(u-t)} \Phi(d_+(u)) - rK e^{-r(u-t)} \Phi(d_-(u)) \right) du, \quad (1.78)$$

where

$$d_{\pm}(u) = \frac{\ln(S(t)/S_A(u)) + (r - q \pm \frac{1}{2}\sigma^2)(u - t)}{\sigma\sqrt{u - t}}.$$

Proof. Due to the smooth pasting condition in Proposition 1.10.4, we are justified¹⁷ in applying Ito's lemma. In informal notation,

$$\begin{aligned} dC_A(t) &= 1_{\{S(t) \geq S_A(t)\}} dS(t) \\ &+ 1_{\{S(t) < S_A(t)\}} \left\{ \frac{\partial C_A(t)}{\partial t} dt + \frac{\partial C_A(t)}{\partial S} dS(t) + \frac{1}{2} \sigma^2 S^2 \frac{\partial^2 C_A(t)}{\partial S^2} dt \right\}, \end{aligned} \quad (1.79)$$

¹⁷In particular, there is no local time contribution to $dC_A(t)$ at the boundary.

where we have used the fact that

$$1_{\{S(t) \geq S_A(t)\}} C_A(t) = 1_{\{S(t) \geq S_A(t)\}} (S(t) - K).$$

In the continuation region, $C_A(t, S)$ satisfies the PDE (1.47), i.e.

$$\frac{\partial C_A(t, S)}{\partial t} + (r - q)S \frac{\partial C_A(t, S)}{\partial S} + \frac{1}{2}\sigma^2 S^2 \frac{\partial^2 C_A(t, S)}{\partial S^2} = rC_A(t, S).$$

Inserting this into (1.79) we get, after a few rearrangements,

$$\begin{aligned} dC_A(t) &= rC_A(t) dt + 1_{\{S(t) < S_A(t)\}} (r - q)S(t) \frac{\partial C_A(t)}{\partial S} dW^\beta(t) \\ &\quad + 1_{\{S(t) \geq S_A(t)\}} \{((r - q)S(t) - rC_A(t)) dt + \sigma S(t) dW^\beta(t)\} \\ &= rC_A(t) dt + 1_{\{S(t) < S_A(t)\}} (r - q)S(t) \frac{\partial C_A(t)}{\partial S} dW^\beta(t) \\ &\quad + 1_{\{S(t) \geq S_A(t)\}} \{(rK - qS(t)) dt + \sigma S(t) dW^\beta(t)\} \end{aligned}$$

Setting $y(t) = C_A(t)/\beta(t)$, it follows from Ito's lemma that

$$\begin{aligned} dy(t) &= e^{-rt} 1_{\{S(t) < S_A(t)\}} (r - q)S(t) \frac{\partial C_A(t)}{\partial S} dW^\beta(t) \\ &\quad + 1_{\{S(t) \geq S_A(t)\}} e^{-rt} \{(rK - qS(t)) dt + \sigma S(t) dW^\beta(t)\}. \end{aligned}$$

Integrating and taking expectations leads to

$$E_t^Q(y(T)) = y(t) + \int_t^T e^{-ru} E_t^Q(1_{\{S(u) \geq S_A(u)\}} (rK - qS(u))) du.$$

Applying the definition of $y(t)$ and the fact that $y(T) = e^{-rT}(S(T) - K)^+$ proves (1.76). The explicit form of the early exercise premium in (1.78) follows from the properties of GBMD. \square

Remark 1.10.8. Combining results from Lemma 1.10.6 and Proposition 1.10.4, it follows that $E_A(t) \geq 0$, so $C_A(t) \geq c(t)$ as expected.

The integral representation of the American call option in Proposition 1.10.7 forms the basis for a number of proposed computational methods for American option pricing. Loosely speaking, these methods are based on the idea of iteratively estimating the exercise boundary $S_A(t)$, often working backwards from $t = T$, after which an application of Proposition 1.10.7 will yield the American option price. A representative example of these methods can be found in Ju [1998]. See Chiarella et al. [2004] for a survey of the literature, and Section 19.7.3 for an application in interest rate derivative pricing.

For a Bermudan option, an integral representation such as that in Proposition 1.10.7 is not possible. Nevertheless, it is still possible to break the

Bermudan call option into the sum of a European option and an early exercise premium. To show this, assume that the allowed exercise dates are $\mathcal{D}(0) = \{T_1, T_2, \dots, T_B\}$, and let $S_B(T_i)$ be the exercise level above which the Bermudan option should be exercised at time T_i , $i = 1, \dots, B$. Notice that if at time T_i we have $S(T_i) > S_B(T_i)$, then C_B will *jump down* in value when time progresses past time T_i , as a reflection of the missed exercise opportunity. Indeed, in the earlier notation of hold and exercise values, we have

$$\begin{aligned} C_B(T_i) &= \max(U(T_i), H(T_i)), \\ C_B(T_i+) &= H(T_i), \end{aligned}$$

which makes the jump in value evident. Given the existence of these jumps, we may write

$$\begin{aligned} dC_B(t) &= rC_B(t) dt + dM(t) \\ &\quad + \sum_{i=1}^B 1_{\{S(T_i) > S_B(T_i)\}} \delta(T_i - t) (H(T_i) - U(T_i)) dt, \end{aligned}$$

where $H(T_i) = C_B(T_i+)$ is the hold value at time T_i , $U(T_i) = S(T_i) - K$, and $M(t)$ is a martingale,

$$dM(t) = \frac{\partial C_B(t)}{\partial S} (r - q) S(t) dW^\beta(t).$$

Deflating C_B by the money market account and forming expectations, we get, since $c(t) = e^{-r(T-t)} \mathbb{E}_t^Q(C_B(T))$,

$$C_B(t) = c(t) + \sum_{T_i \geq t} e^{-r(T_i-t)} \mathbb{E}_t^Q \left(1_{\{S(T_i) > S_B(T_i)\}} (U(T_i) - H(T_i)) \right).$$

As $H(T_i)$ must be less than the exercise value $U(T_i)$ whenever $S(T_i) > S_B(T_i)$ we can simplify this expression to the following result that we, in Section 18.2.3, call the *marginal exercise value decomposition*.

Proposition 1.10.9. *The Bermudan option price $C_B(t)$ satisfies*

$$C_B(t) = c(t) + E_B(t), \quad t \leq T,$$

where the (Bermudan) early exercise premium $E_B(t)$ is defined as

$$E_B(t) = \sum_{T_i \geq t} e^{-r(T_i-t)} \mathbb{E}_t^Q \left((U(T_i) - H(T_i))^+ \right),$$

with $T_1 < T_2 < \dots < T_B = T$ being the set of exercise dates.

As shown in Section 18.2.3, the result in Proposition 1.10.9 may be extended to more complicated processes and payouts than those considered here.

Finite Difference Methods

In Chapter 1 we described how the pricing of a derivative security typically requires either the solution of a parabolic partial differential equation (PDE) or the evaluation of an expectation of a random variable. In realistic applications, both of these price formulations often do not allow for closed-form solution, in which case we must resort to either analytical approximations or, more generally, numerical techniques. In the next two chapters we will describe a number of numerical algorithms useful in derivatives pricing. Analytical approximations will receive ample treatment later in this book, in the context of specific problems.

Our treatment of numerical methods is broken into two main subjects. In this chapter, we cover finite difference solutions of PDEs; and in Chapter 3 we turn to Monte Carlo evaluation of expectations. Many excellent specialist books exist on both topics, including Mitchell and Griffiths [1980], Tavella and Randall [2000], and Glasserman [2004]; our treatment only surveys the most important concepts, as required for our needs in this book. We do provide, however, a number of schemes rarely described in detail in the finance literature and also supplement our analysis with a number of “tricks of the trade”, particularly in the application of finite difference grids.

The analysis of numerical PDE solutions in this chapter is arranged in two blocks. First, in Sections 2.1–2.8 we study the basic mechanics of the finite difference grid method for one-dimensional PDEs. Subsequently, Sections 2.9–2.12 then apply operator splitting techniques to extend the finite difference method to PDE of dimensions two and higher. The analysis culminates with a presentation of *ADI schemes* for multi-dimensional PDEs with mixed partial derivatives.

2.1 1-Dimensional PDEs: Problem Formulation

Initially, we will consider the numerical solution of the general one-dimensional terminal value PDE problem

$$\frac{\partial V}{\partial t} + \mathcal{L}V = 0, \quad (2.1)$$

where \mathcal{L} is the operator

$$\mathcal{L} = \mu(t, x) \frac{\partial}{\partial x} + \frac{1}{2}\sigma(t, x)^2 \frac{\partial^2}{\partial x^2} - r(t, x),$$

and where $V = V(t, x)$ satisfies a terminal condition $V(T, x) = g(x)$. We recognize the PDE as being an extension of the Black-Scholes PDE (1.47) to general time- and state-dependent drift (μ), volatility (σ), and interest rate (r). Underneath the PDE lies a physical model where a state variable process $x(\cdot)$ follows an SDE of the form

$$dx(t) = \mu(t, x(t)) dt + \sigma(t, x(t)) dW(t) \quad (2.2)$$

where $W(t)$ is a Brownian motion in the risk-neutral probability measure Q . Let the range of values attainable by $x(t)$ on $t \in [0, T]$ be denoted $\mathcal{B} \subseteq \mathbb{R}$, and assume that the functions $\mu, \sigma, r : [0, T] \times \mathcal{B} \rightarrow \mathbb{R}$ are sufficiently regular to make (2.1) and (2.2) meaningful (see Chapter 1).

The terminal value problem above is, as discussed earlier, a *Cauchy problem* to be solved for $V(t, x)$ on $(t, x) \in [0, T] \times \mathcal{B}$. In many cases of practical interest, further boundary conditions are applied in the spatial (x) domain. If such boundary conditions are expressed directly in terms of V (rather than its derivatives) we have a *Dirichlet boundary problem*. For instance, a so-called *up-and-out barrier option* will pay out $g(x(T))$ at time T if and only if $x(t)$ stays strictly below a contractually specified barrier level H at all times $t \leq T$. If, on the other hand, $x(t)$ touches H at any time during the life of the contract, it will expire worthless (or “knock out”). In this case, the PDE is only to be solved on $(t, x) \in [0, T] \times (\mathcal{B} \cap (-\infty, H))$ and is subject to the Dirichlet boundary condition

$$V(t, H) = 0, \quad t \in [0, T],$$

which expresses that the option has no value for $x \geq H$. We note that it is not uncommon to encounter options where the spatial domain boundaries are functions of time, a situation we shall deal with in Section 2.7.1. Also, as we shall see shortly, sometimes boundary conditions are conveniently expressed in terms of derivatives of V .

For numerical solution of the PDE (2.1), we often need to assume that the domain of the state variable x is finite, even in situations where (2.1) is supposed to hold for an infinite domain. Suitable truncation of the domain can often be done probabilistically, based on a confidence interval for $x(T)$. To illustrate the procedure, consider the Black-Scholes PDE (1.47) applied to a call option with strike K . A common first step is to use the transformation $x = \ln S$, such that the PDE has constant coefficients,

$$\frac{\partial V}{\partial t} + \left(r - \frac{1}{2}\sigma^2\right) \frac{\partial V}{\partial x} + \frac{1}{2}\sigma^2 \frac{\partial^2 V}{\partial x^2} - rV = 0, \quad (2.3)$$

with terminal value (for a call option) $V(T, x) = (e^x - K)^+$. The domain of x is here the entire real line, $\mathcal{B} = \mathbb{R}$. We know (from (1.39)) that

$$x(T) = x(0) + \left(r - \frac{1}{2}\sigma^2\right)T + \sigma(W(T) - W(0)), \quad (2.4)$$

which is a Gaussian random variable with mean $\bar{x} = x(0) + (r - \frac{1}{2}\sigma^2)T$ and variance $\sigma^2 T$. Consider now replacing the domain $(-\infty, \infty)$ with the finite interval $[\bar{x} - \alpha\sigma\sqrt{T}, \bar{x} + \alpha\sigma\sqrt{T}]$ for some positive constant α . The likelihood of $x(T)$ falling outside of this interval is easily seen to be $2\Phi(-\alpha)$ (where, as always, $\Phi(z)$ is the standard Gaussian cumulative distribution function). If, say, we set α to 4, $2\Phi(-4) = 6.3 \times 10^{-5}$, which is an insignificant probability for most applications. Larger (smaller) values of α will make the truncation error smaller (larger) and will ultimately require more (less) effort in a numerical scheme. We recommend values of α somewhere between 3 and 5 for most applications. For the Black-Scholes case, a rigorous estimate of the error imposed by domain truncation is given in Kangro and Nicolaides [2000].

In many cases of practical interest, it is not possible to write down an exact confidence interval for $x(T)$. In such cases, one instead may use an approximate confidence interval, found by, for instance, using “average” values for $\mu(t, x)$ and $\sigma(t, x)$. High precision in these estimates is typically not needed.

2.2 Finite Difference Discretization

In order to solve the PDE (2.1) numerically, we now wish to discretize it on the rectangular domain $(t, x) \in [0, T] \times [\underline{M}, \overline{M}]$, where \overline{M} and \underline{M} are finite constants, possibly found by a truncation procedure such as the one outlined above. We first introduce two equidistant¹ grids $\{t_i\}_{i=0}^n$ and $\{x_j\}_{j=0}^{m+1}$ where $t_i = iT/n \triangleq i\Delta_t$, $i = 0, 1, \dots, n$, and $x_j = \underline{M} + j(\overline{M} - \underline{M})/(m+1) \triangleq \underline{M} + j\Delta_x$, $j = 0, 1, \dots, m+1$. The terminal value $V(T, x) = g(x)$ is imposed at $t_n = T$, and spatial boundary conditions are imposed at x_0 and x_{m+1} .

2.2.1 Discretization in x -Direction. Dirichlet Boundary Conditions

We first focus on the spatial operator \mathcal{L} and restrict x to take values in the interior of the spatial grid $x \in \{x_j\}_{j=1}^m$. Consider replacing the first- and second-order partial derivatives with first- and second-order difference operators:

¹Non-equidistant grids are often required in practice and will be covered in Section 2.4.

$$\delta_x V(t, x_j) \triangleq \frac{V(t, x_{j+1}) - V(t, x_{j-1})}{2\Delta_x}, \quad (2.5)$$

$$\delta_{xx} V(t, x_j) \triangleq \frac{V(t, x_{j+1}) + V(t, x_{j-1}) - 2V(t, x_j)}{\Delta_x^2}. \quad (2.6)$$

These operators are accurate to second order. Formally²,

Lemma 2.2.1.

$$\delta_x V(t, x_j) = \frac{\partial V(t, x_j)}{\partial x} + O(\Delta_x^2),$$

$$\delta_{xx} V(t, x_j) = \frac{\partial^2 V(t, x_j)}{\partial x^2} + O(\Delta_x^2).$$

Proof. A Taylor expansion of $V(t, x)$ around the point $x = x_j$ gives

$$\begin{aligned} V(t, x_{j+1}) &= V(t, x_j) + \Delta_x \frac{\partial V(t, x_j)}{\partial x} \\ &\quad + \frac{1}{2} \Delta_x^2 \frac{\partial^2 V(t, x_j)}{\partial x^2} + \frac{1}{6} \Delta_x^3 \frac{\partial^3 V(t, x_j)}{\partial x^3} + O(\Delta_x^4), \end{aligned}$$

and

$$\begin{aligned} V(t, x_{j-1}) &= V(t, x_j) - \Delta_x \frac{\partial V(t, x_j)}{\partial x} \\ &\quad + \frac{1}{2} \Delta_x^2 \frac{\partial^2 V(t, x_j)}{\partial x^2} - \frac{1}{6} \Delta_x^3 \frac{\partial^3 V(t, x_j)}{\partial x^3} + O(\Delta_x^4). \end{aligned}$$

Insertion of these expressions into (2.5) and (2.6) gives the desired result.

□

In other words, if we introduce the discrete operator

$$\widehat{\mathcal{L}} = \mu(t, x)\delta_x + \frac{1}{2}\sigma(t, x)^2\delta_{xx} - r(t, x),$$

we have, for $x \in \{x_j\}_{j=1}^m$,

$$\mathcal{L}V(t, x) = \widehat{\mathcal{L}}V(t, x) + O(\Delta_x^2).$$

With attention restricted to values on the grid $\{x_j\}_{j=1}^m$, we can view $\widehat{\mathcal{L}}$ as a matrix, once we specify the side boundary conditions at x_0 and x_{m+1} . For the Dirichlet case, assume for instance that

$$V(x_0, t) = \underline{f}(t, x_0), \quad V(x_{m+1}, t) = \overline{f}(t, x_{m+1}),$$

²Recall that a function $f(h)$ is of order $O(e(h))$ if $|f(h)|/|e(h)|$ is bounded from above by a positive constant in the limit $h \rightarrow 0$.

for given functions $\underline{f}, \bar{f} : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$. With³ $\mathbf{V}(t) \triangleq (V(t, x_1), \dots, V(t, x_m))^\top$ and, for $j = 1, \dots, m$,

$$c_j(t) \triangleq -\sigma(t, x_j)^2 \Delta_x^{-2} - r(t, x_j), \quad (2.7)$$

$$u_j(t) \triangleq \frac{1}{2}\mu(t, x_j)\Delta_x^{-1} + \frac{1}{2}\sigma(t, x_j)^2\Delta_x^{-2}, \quad (2.8)$$

$$l_j(t) \triangleq -\frac{1}{2}\mu(t, x_j)\Delta_x^{-1} + \frac{1}{2}\sigma(t, x_j)^2\Delta_x^{-2}, \quad (2.9)$$

we can write

$$\hat{\mathcal{L}}\mathbf{V}(t) = \mathbf{A}(t)\mathbf{V}(t) + \boldsymbol{\Omega}(t), \quad (2.10)$$

where \mathbf{A} is a *tri-diagonal matrix*

$$\mathbf{A}(t) = \begin{pmatrix} c_1(t) & u_1(t) & 0 & 0 & 0 & \dots & 0 \\ l_2(t) & c_2(t) & u_2(t) & 0 & 0 & \dots & 0 \\ 0 & l_3(t) & c_3(t) & u_3(t) & 0 & \dots & 0 \\ 0 & 0 & l_4(t) & c_4(t) & u_4(t) & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & 0 & 0 & 0 & l_{m-1}(t) & c_{m-1}(t) & u_{m-1}(t) \\ 0 & 0 & 0 & 0 & 0 & l_m(t) & c_m(t) \end{pmatrix} \quad (2.11)$$

and $\boldsymbol{\Omega}(t)$ is a vector containing boundary values

$$\boldsymbol{\Omega}(t) = \begin{pmatrix} l_1(t)\underline{f}(t, x_0) \\ 0 \\ \vdots \\ 0 \\ u_m(t)\bar{f}(t, x_{m+1}) \end{pmatrix}.$$

As discussed earlier, sometimes one or both of the functions \bar{f} and \underline{f} are explicitly imposed as part of the option specification (as is the case for a knock-out options). In other cases, asymptotics may be necessary to establish these functions. For instance, for the case of a simple call option on a stock paying no dividends, we can set

$$\begin{aligned} \bar{f}(t, x) &= e^x - Ke^{-r(T-t)}, \\ \underline{f}(t, x) &= 0, \end{aligned}$$

where we, as before, have set $x = \ln S$ (S being the stock price) and assumed that the strike K is positive. The result for \underline{f} is obvious; the result for \bar{f} follows from the fact that a deep in-the-money call option will almost certainly pay at maturity the stock (the present value of which is just $S = e^x$) minus the strike (the present value of which is $Ke^{-r(T-t)}$).

³For clarity, this chapter uses boldface type for all vectors and matrices.

2.2.2 Other Boundary Conditions

Deriving asymptotic Dirichlet conditions can be quite involved for complicated option payouts and is often inconvenient in implementations. Rather than having to perform an asymptotic analysis for each and every type of option payout, it would be preferable to have a general-purpose mechanism for specifying the boundary condition. One common idea involves making assumptions on the form of the functional dependency between V and x at the grid boundaries, often from specification of relationships between spatial derivatives. For instance, if we impose the condition that the second derivative of V is zero at the upper boundary (x_{m+1}) — that is, V is a linear function of x — we can write (effectively using a downward discretization of the second derivative)

$$\frac{V(t, x_{m+1}) + V(t, x_{m-1}) - 2V(t, x_m)}{\Delta_x^2} = 0 \\ \Rightarrow V(t, x_{m+1}) = 2V(t, x_m) - V(t, x_{m-1}).$$

A similar assumption at the lower spatial boundary yields

$$V(t, x_0) = 2V(t, x_1) - V(t, x_2).$$

For PDEs discretized in the logarithm of some asset, it may be more natural to assume that $V(t, x) \propto e^x$ at the boundaries; equivalently, we can assume that $\partial V / \partial x = \partial^2 V / \partial x^2$ at the boundary. When discretized in downward fashion at the upper boundary (x_{m+1}), this implies that

$$\frac{V(t, x_{m+1}) - V(t, x_m)}{\Delta_x} = \frac{V(t, x_{m+1}) + V(t, x_{m-1}) - 2V(t, x_m)}{\Delta_x^2}$$

or (assuming that $\Delta_x \neq 1$)

$$V(t, x_{m+1}) = V(t, x_{m-1}) \frac{1}{\Delta_x - 1} + V(t, x_m) \frac{\Delta_x - 2}{\Delta_x - 1}.$$

Similarly,

$$V(t, x_0) = V(t, x_1) \frac{2 + \Delta_x}{1 + \Delta_x} - V(t, x_2) \frac{1}{\Delta_x + 1}.$$

Common for both methods above — and for the Dirichlet specification discussed earlier — is that they give rise to boundary specifications through simple linear systems of the general form

$$V(t, x_{m+1}) = k_m(t)V(t, x_m) + k_{m-1}(t)V(t, x_{m-1}) + \bar{f}(t, x_{m+1}), \quad (2.12)$$

$$V(t, x_0) = k_1(t)V(t, x_1) + k_2(t)V(t, x_2) + \underline{f}(t, x_0). \quad (2.13)$$

This boundary specification can be captured in the matrix system (2.10) by simply rewriting a few components of $\mathbf{A}(t)$; specifically, we must set

$$\begin{aligned}
c_m(t) &= -\sigma(t, x_m)^2 \Delta_x^{-2} - r(t, x_m) + k_m(t) u_m(t), \\
l_m(t) &= -\frac{1}{2} \mu(t, x_m) \Delta_x^{-1} + \frac{1}{2} \sigma(t, x_m)^2 \Delta_x^{-2} + k_{m-1}(t) u_m(t), \\
c_1(t) &= -\sigma(t, x_j)^2 \Delta_x^{-2} - r(t, x_j) + k_1(t) l_1(t), \\
u_1(t) &= \frac{1}{2} \mu(t, x_1) \Delta_x^{-1} + \frac{1}{2} \sigma(t, x_1)^2 \Delta_x^{-2} + k_2(t) l_1(t).
\end{aligned}$$

All other components of \mathbf{A} remain as in (2.11); note that \mathbf{A} remains tri-diagonal.

An alternative approach to specification of boundary conditions in the x -domain involves using the PDE itself to determine the boundary conditions, through replacement of all central difference operators with one-sided differences at the boundaries. Section 10.1.5.2 contains a detailed example of this idea; ultimately, this approach leads to boundary conditions that can also be written in the form (2.12)–(2.13).

2.2.3 Time-Discretization

To simplify notation, assume for now that $\Omega(t) = 0$ for all t , as will be the case if, say, we use the linear or linear-exponential boundary conditions outlined earlier. On the spatial grid, our original PDE can be written

$$\frac{\partial \mathbf{V}(t)}{\partial t} = -\mathbf{A}(t)\mathbf{V}(t) + O(\Delta_x^2)$$

which, ignoring the error term⁴, defines a system of coupled ordinary differential equations (ODEs).

A number of methods are available for the numerical solution of coupled ODEs; see, e.g., Press et al. [1992]. We here only consider basic two-level time-stepping schemes, where grid computations at time t_i involve only PDE values at times t_i and t_{i+1} . Focusing the attention on a particular bucket $[t_i, t_{i+1}]$, the choice for the finite difference approximation of $\partial V / \partial t$ is obvious:

$$\frac{\partial \mathbf{V}}{\partial t} \approx \frac{\mathbf{V}(t_{i+1}) - \mathbf{V}(t_i)}{\Delta_t}.$$

Not so obvious, however, is to which time in the interval $[t_i, t_{i+1}]$ we should associate this derivative. To be general, consider picking a time $t_i^{i+1}(\theta) \in [t_i, t_{i+1}]$, given by

$$t_i^{i+1}(\theta) = (1 - \theta)t_{i+1} + \theta t_i, \quad (2.14)$$

where $\theta \in [0, 1]$ is a parameter. We then write

$$\frac{\partial \mathbf{V}(t_i^{i+1}(\theta))}{\partial t} \approx \frac{\mathbf{V}(t_{i+1}) - \mathbf{V}(t_i)}{\Delta_t}.$$

⁴Note that the error term $O(\Delta_x^2)$ is here to be interpreted as an m -dimensional vector. We will use such short-hand notation throughout this chapter.

By a Taylor expansion, it is easy to see that this expression is first-order accurate in the time step when $\theta \neq \frac{1}{2}$, and second-order accurate when $\theta = \frac{1}{2}$. Written compactly,

$$\frac{\partial \mathbf{V}(t_i^{i+1}(\theta))}{\partial t} = \frac{\mathbf{V}(t_{i+1}) - \mathbf{V}(t_i)}{\Delta_t} + 1_{\{\theta \neq \frac{1}{2}\}} O(\Delta_t) + O(\Delta_t^2). \quad (2.15)$$

This result on the convergence order is intuitive since only in the case $\theta = \frac{1}{2}$ is the difference coefficient precisely central; for all other cases, the difference coefficient is either predominantly backward in time or predominantly forward in time.

The time-discretization technique introduced above is known as a *theta scheme*. The special cases of $\theta = 1$, $\theta = 0$, and $\theta = \frac{1}{2}$ are known as the *fully implicit scheme*, the *fully explicit scheme*, and the *Crank-Nicolson scheme*, respectively. In light of the convergence result (2.15), one may wonder why anything other than the Crank-Nicolson scheme is ever used. The CN method is, indeed, often the method of choice, but there are situations where a straight application of the Crank-Nicolson scheme can lead to oscillations in the numerical solution or its spatial derivatives. Judicial application of the fully implicit method can often alleviate these problems, as we shall discuss later. The fully explicit method should never be used due to poor convergence and stability properties (see Section 2.3), but has nevertheless managed to survive in a surprisingly large number of finance texts and papers.

2.2.4 Finite Difference Scheme

We now proceed to combine the discretizations (2.10) and (2.15) into a complete finite difference scheme. First, we expand

$$\begin{aligned} \mathbf{A}(t_i^{i+1}(\theta)) \mathbf{V}(t_i^{i+1}(\theta)) &= \theta \mathbf{A}(t_i^{i+1}(\theta)) \mathbf{V}(t_i) \\ &\quad + (1 - \theta) \mathbf{A}(t_i^{i+1}(\theta)) \mathbf{V}(t_{i+1}) + 1_{\{\theta \neq \frac{1}{2}\}} O(\Delta_t) + O(\Delta_t^2), \end{aligned}$$

such that our PDE can be represented as

$$\begin{aligned} \frac{\mathbf{V}(t_{i+1}) - \mathbf{V}(t_i)}{\Delta_t} + 1_{\{\theta \neq \frac{1}{2}\}} O(\Delta_t) + O(\Delta_t^2) \\ &= -\mathbf{A}(t_i^{i+1}(\theta)) \mathbf{V}(t_i^{i+1}(\theta)) + O(\Delta_x^2) \\ &= -\theta \mathbf{A}(t_i^{i+1}(\theta)) \mathbf{V}(t_i) - (1 - \theta) \mathbf{A}(t_i^{i+1}(\theta)) \mathbf{V}(t_{i+1}) \\ &\quad + 1_{\{\theta \neq \frac{1}{2}\}} O(\Delta_t) + O(\Delta_t^2) + O(\Delta_x^2). \end{aligned}$$

Multiplying through with Δ_t gives rise to the complete finite difference representation of the PDE solution at times t_i and t_{i+1} :

Proposition 2.2.2. *On the grid $\{x_j\}_{j=1}^m$, the solution to (2.1) at times t_i and t_{i+1} is characterized by*

$$(\mathbf{I} - \theta \Delta_t \mathbf{A}(t_i^{i+1}(\theta))) \mathbf{V}(t_i) = (\mathbf{I} + (1 - \theta) \Delta_t \mathbf{A}(t_i^{i+1}(\theta))) \mathbf{V}(t_{i+1}) + e_i^{i+1}, \quad (2.16)$$

where \mathbf{I} is the $m \times m$ identity matrix, and e_i^{i+1} is an error term

$$e_i^{i+1} = \Delta_t O(\Delta_x^2) + 1_{\{\theta \neq \frac{1}{2}\}} O(\Delta_t^2) + O(\Delta_t^3). \quad (2.17)$$

Let $\widehat{\mathbf{V}}(t_i, x_j)$ denote the approximation to the true solution $V(t_i, x_j)$ obtained by using (2.16) without the error term. Defining

$$\widehat{\mathbf{V}}(t) = \left(\widehat{V}(t, x_1), \dots, \widehat{V}(t, x_m) \right)^\top,$$

we have

$$(\mathbf{I} - \theta \Delta_t \mathbf{A}(t_i^{i+1}(\theta))) \widehat{\mathbf{V}}(t_i) = (\mathbf{I} + (1 - \theta) \Delta_t \mathbf{A}(t_i^{i+1}(\theta))) \widehat{\mathbf{V}}(t_{i+1}). \quad (2.18)$$

For a known value of $\widehat{\mathbf{V}}(t_{i+1})$, (2.18) defines a simple linear system of equations that can be solved for $\widehat{\mathbf{V}}(t_i)$ by standard methods. Simplifying matters is the fact that the matrix $(\mathbf{I} - \theta \Delta_t \mathbf{A}(t_i^{i+1}(\theta)))$ is tri-diagonal, allowing us to solve (2.18) in only $O(m)$ operations; see Press et al. [1992] for an algorithm⁵.

Starting from the prescribed terminal condition $V(t_n, x_j) = g(x_j)$, $j = 1, \dots, m$, we can now use (2.18) to iteratively step backward in time until we ultimately recover $\widehat{\mathbf{V}}(0)$. This procedure is known as *backward induction*.

Proposition 2.2.3. *The theta scheme (2.18) recovers $\widehat{\mathbf{V}}(0)$ in $O(mn)$ operations. If the scheme converges, the error on $\widehat{\mathbf{V}}(0)$ compared to the exact solution $\mathbf{V}(0)$ is of order*

$$O(\Delta_x^2) + 1_{\{\theta \neq \frac{1}{2}\}} O(\Delta_t) + O(\Delta_t^2).$$

Proof. The backward induction algorithm requires the solution of n tri-diagonal systems, one per time step, for a total computational cost of $O(mn)$. The local truncation error on $\widehat{\mathbf{V}}(t_i)$ is e_i^{i+1} , making the global truncation error after n time steps of order ne_i^{i+1} . Combining (2.17) with the fact that $n = T/\Delta_t = O(\Delta_t^{-1})$ gives the order result listed in the proposition. \square

⁵The special case of an explicit scheme ($\theta = 0$) provides us with a direct expression for $V(t_i, x_j)$ in terms of $V(t_{i+1}, x_{j-1})$, $V(t_{i+1}, x_j)$, and $V(t_{i+1}, x_{j+1})$, a scheme that is easily visualized as a “trinomial tree”. The intuitive nature of the explicit scheme coupled with the fact that no matrix equation must be solved may explain the popularity of this scheme in the finance literature, despite its poor numerical qualities (see Section 2.3). We stress that the workload of the explicit scheme is still $O(m)$ per time step, as is the case for all theta schemes.

It follows from Proposition 2.2.3 that the Crank-Nicolson scheme is second-order convergent in the time step, and all other theta schemes are first-order convergent in the time step. All theta-schemes are second-order convergent in the spatial step Δ_x .

In deriving (2.18), we assumed earlier that the boundary vector was zero, $\Omega(t) = 0$. Including a non-zero boundary vector into the scheme is, however, straightforward and results in a time-stepping scheme of the form

$$\begin{aligned} (\mathbf{I} - \theta \Delta_t \mathbf{A}(t_i^{i+1}(\theta))) \hat{\mathbf{V}}(t_i) &= (\mathbf{I} + (1 - \theta) \Delta_t \mathbf{A}(t_i^{i+1}(\theta))) \hat{\mathbf{V}}(t_{i+1}) \\ &\quad + (1 - \theta) \Omega(t_{i+1}) + \theta \Omega(t_i). \end{aligned} \quad (2.19)$$

Again, this system is easily solved for $\hat{\mathbf{V}}(t_i)$ by a standard tri-diagonal equation solver.

As a final point, we stress that the finite difference scheme above ultimately yields a full vector of values $\hat{\mathbf{V}}(0)$ at time 0, with one element per value of x_j , $j = 1, \dots, m$. In general, we are mainly interested in $V(0, x(0))$, where $x(0)$ is the known value of x at time 0. There is no need to include $x(0)$ in the grid, as we can simply employ an interpolator (e.g., a cubic spline) on this vector $\hat{\mathbf{V}}(0)$ to compute $V(0, x(0))$. Clearly, such an interpolator should be at least second-order accurate to avoid interfering with the overall $O(\Delta_x^2)$ convergence of the finite difference scheme. Assuming the interpolator is sufficiently smooth, we can also use it to compute various partial derivatives with respect to x that we may be interested in. Alternatively, these can be computed by the same type of finite difference coefficients discussed in Section 2.2.1. The derivative $\partial V(0, x(0))/\partial t$ — the *time decay* — can be picked up from the grid in the same fashion.

Remark 2.2.4. The scheme (2.18) may, without affecting convergence order, be replaced with

$$(\mathbf{I} - \theta \Delta_t \mathbf{A}(t_i)) \hat{\mathbf{V}}(t_i) = (\mathbf{I} + (1 - \theta) \Delta_t \mathbf{A}(t_{i+1})) \hat{\mathbf{V}}(t_{i+1}).$$

2.3 Stability

2.3.1 Matrix Methods

Ignoring the contributions from boundary conditions, the finite difference scheme developed in the previous section can be rewritten

$$\hat{\mathbf{V}}(t_i) = \mathbf{B}_i^{i+1} \hat{\mathbf{V}}(t_{i+1}), \quad (2.20)$$

where

$$\mathbf{B}_i^{i+1} \triangleq (\mathbf{I} - \theta \Delta_t \mathbf{A}(t_i^{i+1}(\theta)))^{-1} (\mathbf{I} + (1 - \theta) \Delta_t \mathbf{A}(t_i^{i+1}(\theta))).$$

That is, for any $0 \leq k < n$,

$$\widehat{\mathbf{V}}(t_k) = \mathbf{B}_k^n \widehat{\mathbf{V}}(t_n), \quad \mathbf{B}_k^n \triangleq \mathbf{B}_k^{k+1} \mathbf{B}_{k+1}^{k+2} \dots \mathbf{B}_{n-1}^n.$$

We say that the scheme is *stable* if $|\widehat{\mathbf{V}}(t_k)|$ is bounded for all $0 \leq k < n$. Assuming $|\widehat{\mathbf{V}}(T)| < \infty$, a necessary and sufficient condition for stability is that there exists a constant K such that for all $0 \leq k < n$

$$|\mathbf{B}_k^n| \leq K, \quad (2.21)$$

where $|\cdot|$ is any matrix norm, e.g. the spectral norm or the infinity norm⁶. See Mitchell and Griffiths [1980] for further details.

2.3.2 Von Neumann Analysis

For simple problems with time- and space-independent coefficients, it may be possible to establish the spectral norm of \mathbf{B}_k^n by direct methods (see e.g. Mitchell and Griffiths [1980], Kraaijevanger et al. [1987], Lenferink and Spijker [1991], Spijker and Straetemans [1997]), but generally the stability criterion (2.21) is difficult to evaluate. While certain somewhat simpler matrix-based methods exist to establish necessary conditions for stability (again, see Mitchell and Griffiths [1980]), we shall here only consider a “local” method, known as the *von Neumann method*. In principle, the von Neumann method only holds for finite difference schemes where the underlying PDE has constant coefficients, but there is much numerical evidence to support wider application⁷. The von Neumann method does not directly consider the effect of boundary conditions on stability, but (for constant coefficient problems) provides a necessary condition for stability irrespective of the type of boundary condition.

The basis for the von Neumann analysis is the observation that a real function sampled on a finite number of points is uniquely defined by a complex Fourier series. For our PDE solution sampled on the spatial grid, the precise result is

$$V(t_k, x_j) = \sum_l H_l(t_k) e^{-i\omega_l j \Delta_x},$$

where $H_l(t_k)$ and ω_l are the amplification factor (discrete Fourier transform) and wave number for the l -th mode, respectively. Notice that i here denotes

⁶The spectral norm of a matrix \mathbf{C} is defined as the largest absolute eigenvalue of $(\mathbf{C}^\top \mathbf{C})^{1/2}$. The infinity norm is defined as $\max_i \sum_j |C_{i,j}|$.

⁷In the application to PDEs with non-constant coefficients, it may help to think of the von Neumann analysis as being applied to the PDE locally with “frozen” coefficients, followed by an examination of the worst case among all frozen coefficients.

the imaginary unit, $i^2 = -1$, with k (momentarily) having taken the role of the time index in the finite difference grid. For the constant coefficient case, a key fact for our PDE problem is that

$$H_l(t_k) = H_l(t_{k+1})\xi_l^{-1},$$

where ξ_l is a mode-specific *amplification factor* independent of time. To determine how a solution is propagated back through the finite difference grid, it thus suffices to consider a test function of the form

$$v(t_k, x_j) = \xi(\omega)^{n-k} e^{i\omega j \Delta_x}. \quad (2.22)$$

According to the Von Neumann criterion, stability of (2.20) requires that the *modulus of the amplification factor* $\xi(\omega)$ is less or equal to one, independent of the wave number:

$$\forall \omega : |\xi(\omega)| \leq 1. \quad (2.23)$$

This criterion is natural and merely expresses that all eigenmodes should be damped, and not exponentially amplified, by the finite difference scheme.

Turning to our system (2.20), assume for simplicity that $r(t, x) = 0$. A positive interest rate (we will nearly always have $r(t, x) > 0$) introduces some extra dampening through discounting effects and will, if anything, lead to better stability properties than the case of zero interest rates. Writing $v(t_k, x_j) = v_{k,j}$, $\sigma(t_k^{k+1}(\theta), x_j) = \sigma_{k,j}$, and $\mu(t_k^{k+1}(\theta), x_j) = \mu_{k,j}$, the von Neumann analysis gives the following result:

Proposition 2.3.1. Define $\alpha = \Delta_t / (\Delta_x)^2$. For (2.20) with $r(t, x) = 0$, the von Neumann stability criterion is

$$1 \geq \theta \geq \frac{1}{2} - \frac{1}{\alpha} \left(\frac{\sigma_{k,j}^2}{\sigma_{k,j}^4 + \mu_{k,j}^2 \Delta_x^2 + |\mu_{k,j}^2 \Delta_x^2 - \sigma_{k,j}^4|} \right), \quad (2.24)$$

to hold for all $k = 0, 1, \dots, n-1$, $j = 1, 2, \dots, m$.

Proof. Define $\varsigma_{k,j}^\pm = \sigma_{k,j}^2 \pm \Delta_x \mu_{k,j}$. A local application of (2.20) gives

$$\begin{aligned} v_{k,j-1} \left(-\frac{\alpha\theta}{2} \varsigma_{k,j}^- \right) + v_{k,j} (1 + \alpha\theta\sigma_{k,j}^2) + v_{k,j+1} \left(-\frac{\alpha\theta}{2} \varsigma_{k,j}^+ \right) = \\ v_{k+1,j} \left(\frac{\alpha(1-\theta)}{2} \varsigma_{k,j}^- \right) + v_{k+1,j} (1 - \alpha(1-\theta)\sigma_{k,j}^2) + v_{k+1,j+1} \left(\frac{\alpha(1-\theta)}{2} \varsigma_{k,j}^+ \right) \end{aligned}$$

with α defined above. Inserting (2.22) and rearranging (using Euler's formulas for sin and cos) yields

$$\xi(\omega) = \frac{1 - (1-\theta)\alpha\sigma_{k,j}^2(1 - \cos \omega \Delta_x) + i(1-\theta)\alpha\Delta_x \mu_{k,j} \sin \omega \Delta_x}{1 + \theta\alpha\sigma_{k,j}^2(1 - \cos \omega \Delta_x) - i\theta\alpha\Delta_x \mu_{k,j} \sin \omega \Delta_x}.$$

Note that ξ is a function of k and j , due to the non-constant PDE parameters. As discussed earlier (see also Mitchell and Griffiths [1980]), we expect the system to be stable if the criterion (2.23) holds for all k and j in the grid. Computing the modulus of ξ and requiring that it does not exceed one leads, after straightforward manipulations, to the stability criterion

$$\forall \omega : 2\alpha\sigma_{k,j}^2 + (2\theta - 1)\alpha^2 [\sigma_{k,j}^4 + \mu_{k,j}^2 \Delta_x^2 + \cos \omega \Delta_x (\mu_{k,j}^2 \Delta_x^2 - \sigma_{k,j}^4)] \geq 0.$$

As $\cos \omega \Delta_x \in [-1, 1]$, this expression can be simplified to (2.24). \square

From (2.24) we can immediately conclude that the finite difference scheme is always stable if $\frac{1}{2} \leq \theta \leq 1$, irrespective of the magnitudes of Δ_x and Δ_t . For $\frac{1}{2} \leq \theta \leq 1$, we therefore say that the theta scheme is *absolutely stable*, or simply *A-stable*. Both the fully implicit ($\theta = 1$) and the Crank-Nicolson ($\theta = \frac{1}{2}$) finite difference schemes are thus *A*-stable. For the explicit scheme ($\theta = 0$), however, stability is *conditional*, requiring

$$\frac{2}{\alpha}\sigma_{k,j}^2 \geq \sigma_{k,j}^4 + \mu_{k,j}^2 \Delta_x^2 + |\mu_{k,j}^2 \Delta_x^2 - \sigma_{k,j}^4|.$$

For small drifts, this expression amounts to the restriction $\sigma_{k,j}^2 \leq \Delta_x^2/\Delta_t$ which can be quite onerous, often requiring the (laborious) use of thousands of time steps in the finite difference grid. We shall not consider fully explicit methods any further in this book.

Returning to the case $\frac{1}{2} \leq \theta \leq 1$, let us introduce a stronger definition of stability. A time-stepping method is said to be *strongly A-stable* if the modulus of the amplification factor ξ is strictly below 1 for any value of the time step, including the limit⁸ $\Delta_t \rightarrow \infty$. From (2.24), we see that if $\Delta_t \rightarrow \infty$ (which implies $\alpha \rightarrow \infty$), then the modulus of the amplification factor could reach 1 in the special case of $\theta = 1/2$. In other words, the Crank-Nicolson scheme is *not* strongly *A*-stable. For large time steps, harmonics in the Crank-Nicolson finite difference solution will effectively not be damped from one time step to the next, opening up the possibility that unwanted high-frequency oscillations can creep into the numerical solution. In practice, this is primarily a problem if high-frequency eigenmodes have high amplification factors, as can happen if there is an outright discontinuity in the terminal value function g . The problem is especially noticeable if the discontinuity in the value function is “close” in both time and space to $t = 0$ and $x = x(0)$ (as would be the case for a short-dated option with a discontinuity close to the starting value of x). Oscillations can be prevented by setting the time step smaller than twice the maximum stable explicit time step (see Tavella and Randall [2000]), but this can often be computationally expensive. We shall deal with other methods to suppress oscillations in Section 2.5.

We conclude this section by noting a deep connection between the stability of a finite difference scheme and its convergence to the true solution

⁸If further $|\xi|$ approaches zero for $\Delta_t \rightarrow 0$, the scheme is said to be *L-stable*.

of the PDE as $\Delta_t \rightarrow 0$ and $\Delta_x \rightarrow 0$. First, we define a finite difference scheme to be *consistent* if local (Taylor) truncation errors approach zero for $\Delta_t \rightarrow 0$ and $\Delta_x \rightarrow 0$. All the schemes we have encountered so far are consistent. Further, define a finite difference scheme to be *convergent* if the difference between the numerical solution and the exact PDE solution at a fixed point in the domain converges to zero uniformly as $\Delta_t \rightarrow 0$ and $\Delta_x \rightarrow 0$ (not necessarily independently of each other). We then have

Theorem 2.3.2 (Lax Equivalence Theorem). *For a well-posed⁹ linear terminal value PDE, a consistent 2-level finite difference scheme is convergent if and only if it is stable.*

A more precise statement of the above result, as well as a proof, can be found in Mitchell and Griffiths [1980].

2.4 Non-Equidistant Discretization

In practice, we often wish to align the finite difference grid to particular dates (e.g., those on which a coupon or a dividend is paid) and particular values of x (e.g., those on which strikes and barriers are positioned). Also, for numerical reasons we may want to make certain important parts of the finite difference grid more densely spaced to concentrate computational effort on domains of particular importance to the solution of the PDE. To do so, we will now relax our earlier assumption of equidistant discretization in time and space. Doing so for the time domain is actually trivial and merely requires us to replace Δ_t in (2.18) with $\Delta_{t,i} \triangleq t_{i+1} - t_i$, where the spacing of the time grid $\{t_i\}_{i=0}^n$ is now no longer constant. The backward induction algorithm can proceed as before. We note that the ability to freely select the time grid will allow us to line up perfectly with dates that carry high significance for the product in question (e.g. dates on which cash flows take place, see Section 2.7.3) or to, say, use coarser time steps for the part of the finite difference grid that is far in the future. For an adaptive algorithm to automatically select the time-step, see d'Halluin et al. [2001].

For the spatial step, we have a number of options to induce non-equidistant spacing. One method involves a non-linear change of variables $y = h(x)$ in the PDE, followed by a regular equidistant discretization in the new variable y . This maps into a non-equidistant discretization in x which, provided that $h(\cdot)$ is chosen carefully, will have the desired geometry. Discussion of this method along with guidelines for choosing $h(\cdot)$ can be found in Chapter 5 of Tavella and Randall [2000]. We will here pursue a more direct alternative, where we simply introduce an irregular grid $\{x_j\}_{j=0}^{m+1}$

⁹Well-posed means that the PDE we are solving has a unique solution that depends continuously on the problem data (PDE coefficients, domain, boundary conditions, etc.)

and redefine the finite difference operators (2.5)–(2.6) to achieve maximum precision. For this, define

$$\Delta_{x,j}^+ \triangleq x_{j+1} - x_j, \quad \Delta_{x,j}^- \triangleq x_j - x_{j-1},$$

and set

$$\delta_x^+ V(t, x_j) = \frac{V(t, x_{j+1}) - V(t, x_j)}{\Delta_{x,j}^+}, \quad \delta_x^- V(t, x_j) = \frac{V(t, x_j) - V(t, x_{j-1})}{\Delta_{x,j}^-}.$$

By a Taylor expansion, we get

$$\begin{aligned} \delta_x^+ V(t, x_j) &= \frac{\partial V(t, x_j)}{\partial x} + \frac{1}{2} \frac{\partial^2 V(t, x_j)}{\partial x^2} \Delta_{x,j}^+ \\ &\quad + \frac{1}{6} \frac{\partial^3 V(t, x_j)}{\partial x^3} (\Delta_{x,j}^+)^2 + O((\Delta_{x,j}^+)^3), \end{aligned} \quad (2.25)$$

$$\begin{aligned} \delta_x^- V(t, x_j) &= \frac{\partial V(t, x_j)}{\partial x} - \frac{1}{2} \frac{\partial^2 V(t, x_j)}{\partial x^2} \Delta_{x,j}^- \\ &\quad + \frac{1}{6} \frac{\partial^3 V(t, x_j)}{\partial x^3} (\Delta_{x,j}^-)^2 + O((\Delta_{x,j}^-)^3). \end{aligned} \quad (2.26)$$

Maximum accuracy on the first-order derivative approximation is achieved by selecting a weighted combination of (2.25)–(2.26) such that the terms of order $O(\Delta_{x,j}^+)$ and $O(\Delta_{x,j}^-)$ cancel. That is, we set

$$\begin{aligned} \delta_x V(t, x_j) &= \frac{\Delta_{x,j}^-}{\Delta_{x,j}^+ + \Delta_{x,j}^-} \cdot \delta_x^+ V(t, x_j) + \frac{\Delta_{x,j}^+}{\Delta_{x,j}^+ + \Delta_{x,j}^-} \cdot \delta_x^- V(t, x_j) \quad (2.27) \\ &= \frac{\partial V(t, x_j)}{\partial x} + O\left(\frac{(\Delta_{x,j}^+)^2 \Delta_{x,j}^- + (\Delta_{x,j}^-)^2 \Delta_{x,j}^+}{\Delta_{x,j}^+ + \Delta_{x,j}^-}\right) \end{aligned}$$

which is second-order accurate, in the sense that reducing both $\Delta_{x,j}^+$ and $\Delta_{x,j}^-$ by a factor of k will reduce the error by a factor of k^2 . To estimate the derivative $\partial^2 V(t, x_j)/\partial x^2$ we set

$$\begin{aligned} \delta_{xx} V(t, x_j) &= \frac{\delta_x^+ V(t, x_j) - \delta_x^- V(t, x_j)}{\frac{1}{2} (\Delta_{x,j}^+ + \Delta_{x,j}^-)} \quad (2.28) \\ &= \frac{\partial^2 V(t, x_j)}{\partial x^2} + O\left(\frac{(\Delta_{x,j}^+)^2 - (\Delta_{x,j}^-)^2}{\Delta_{x,j}^+ + \Delta_{x,j}^-} + \frac{(\Delta_{x,j}^+)^3 + (\Delta_{x,j}^-)^3}{\Delta_{x,j}^+ + \Delta_{x,j}^-}\right) \end{aligned}$$

which is only first-order accurate, unless $\Delta_{x,j}^+ = \Delta_{x,j}^-$. Despite this, the global discretization error will typically remain second-order in the spatial step, even for a non-equidistant grid. A proof of this perhaps somewhat

surprising result can be found in the monograph Axelsson and Barker [1991] on finite element methods.

Development of a theta scheme around the definitions (2.27) and (2.28) proceeds in the same way as in Section 2.2. The resulting time-stepping scheme is identical to (2.18), after a modification of the matrix \mathbf{A} . Specifically, we must simply redefine the c -, u -, and l -arrays in (2.7)–(2.9) as follows:

$$c_j(t) \triangleq \frac{\Delta_{x,j}^+ - \Delta_{x,j}^-}{\Delta_{x,j}^+ \Delta_{x,j}^-} - \frac{1}{\Delta_{x,j}^- \Delta_{x,j}^+} \sigma(t, x_j)^2 - r(t, x_j), \quad (2.29)$$

$$u_j(t) \triangleq \frac{\Delta_{x,j}^-}{(\Delta_{x,j}^+ + \Delta_{x,j}^-) \Delta_{x,j}^+} \mu(t, x_j) + \frac{1}{(\Delta_{x,j}^+ + \Delta_{x,j}^-) \Delta_{x,j}^+} \sigma(t, x_j)^2, \quad (2.30)$$

$$l_j(t) \triangleq -\frac{\Delta_{x,j}^+}{(\Delta_{x,j}^+ + \Delta_{x,j}^-) \Delta_{x,j}^-} \mu(t, x_j) + \frac{1}{(\Delta_{x,j}^+ + \Delta_{x,j}^-) \Delta_{x,j}^-} \sigma(t, x_j)^2. \quad (2.31)$$

For an example where having a non-equidistant grid is essential to the numerical performance of the scheme, see Section 9.4.3.

2.5 Smoothing and Continuity Correction

2.5.1 Crank-Nicolson Oscillation Remedies

As discussed earlier, for discontinuous terminal conditions, the Crank-Nicolson scheme may exhibit localized oscillations if the time step is too coarse relative to the spatial step. Depending on the timing and spatial position of the discontinuities, these spurious oscillations may negatively affect the computed option value or, more likely, its first (“delta”) or second (“gamma”) x -derivatives. Further, in the presence of discontinuous terminal conditions, the expected $O(\Delta_t^2)$ convergence order of the Crank-Nicolson scheme may not be realized. While $O(\Delta_t^2)$ convergence is possible without spurious oscillations in some multi-level time-stepping schemes, there is evidence that these schemes are less robust than the Crank-Nicolson scheme for many financially relevant problems, see, e.g., Windcliff et al. [2001]. Fortunately, it is relatively easy to remedy the problems in the Crank-Nicolson scheme. Specifically, a theoretical result by Rannacher [1984] shows that second-order convergence can be achieved for the Crank-Nicolson scheme, provided that two simple algorithm modifications are taken:

- The discontinuous terminal payout is least-squares (L^2) projected onto the space of linear Lagrange basis functions¹⁰.

¹⁰Recall that the linear Lagrange basis functions (also called “hat” functions) are simply small triangles given by $l_j(x) = 1_{\{x_{j-1} < x \leq x_j\}} \cdot \frac{x - x_{j-1}}{x_j - x_{j-1}} + 1_{\{x_j < x \leq x_{j+1}\}} \cdot \frac{x_{j+1} - x}{x_{j+1} - x_j}$, $j = 1, \dots, m$. For an algorithm to perform the L^2 -projection, see Pooley et al. [2003].

- Two fully implicit time steps ($\theta = 1$) are taken before we switch to Crank-Nicolson ($\theta = \frac{1}{2}$) time stepping (“Rannacher stepping”).

Both techniques effectively smoothen out the discontinuity before the Crank-Nicolson scheme is applied, dampening the problematic high-frequency modes of the numerical solution. As demonstrated in Pooley et al. [2003] (see also Giles and Carter [2006]), applying either technique in isolation will typically not suffice; both are jointly required to ensure smooth second-order convergence. That said, the application of Lagrange basis function projection may conveniently be substituted with simpler smoothing techniques, with no loss of convergence order. The usefulness of such payoff smoothing extends beyond the case of discontinuous boundary conditions, so we proceed to discuss a few common techniques next.

2.5.2 Continuity Correction

By the Shannon sampling theorem, (see Shannon [1949]) if the spectrum of $g(x)$ contains frequencies higher than $1/(2\Delta_x)$ (the *Nyquist frequency*), information is lost when we sample $g(x)$ on our mesh $\{x_j\}_{j=0}^{m+1}$. In other words, whenever $g(x)$ or its derivatives are non-smooth, we will incur a *quantization* error where important features of the payout (e.g., the discontinuity of the slope of a call option at the strike) will be lost between grid points. As the grid geometry is modified, and the location of critical points (strikes, barriers, etc.) relative to x -grid changes, the computed finite difference solution will jump back and forth in erratic fashion. This so-called *odd-even effect* will result in poor convergence and an undesirably strong dependence of the solution on the grid geometry.

One straightforward way to reduce the odd-even effect (and to smooth out the high-frequency components of the payoff) is to apply a common technique from probability theory known as a *continuity correction*. Here, we simply imagine that the value of g at a grid point x_j represents the average value of the function over the interval $[x_j - (x_j - x_{j-1})/2, x_j + (x_{j+1} - x_j)/2]$. In setting the terminal boundary value $V(T, x_j)$ we thus write

$$V(T, x_j) = \frac{1}{(x_{j+1} - x_{j-1})/2} \int_{x_j - (x_j - x_{j-1})/2}^{x_j + (x_{j+1} - x_j)/2} g(x) dx. \quad (2.32)$$

We note that this implies that $V(T, x_j) \neq g(x_j)$, unless g is linear in x . The application of continuity correction to parabolic PDE solvers was first proposed in Kreiss et al. [1970].

2.5.3 Grid Shifting

Consider the effect of using (2.32) on a *digital call option*, $g(x) = 1_{\{x > H\}}$, where the level H (the digital strike) is located between nodes x_k and

x_{k+1} . For nodes x_j , $j > k + 1$, clearly $V(T, x_j) = 1$; for nodes x_j , $j < k$, $V(T, x_j) = 0$. The smoothing algorithm will have effect only at x_k or x_{k+1} , and will set either $V(T, x_k)$ or $V(T, x_{k+1})$ to a value somewhere between 0 and 1, depending on which of x_k or x_{k+1} is closest to H . If H happens to be exactly midway between x_k or x_{k+1} , the continuity correction is seen to have no effect whatsoever.

The digital option example above gives rise to a method listed in Tavella and Randall [2000] (see also Cheuk and Vorst [1996]). Here, we simply arrange the spatial grid such that the x -values where the payoff (or its derivatives) is discontinuous are exactly midway between grid nodes. If necessary, we can use a scheme with non-equidistant grid spacing to accomplish this (see Section 2.4). Our example above shows that aligning the grid in this way will, in a loose sense, make the payoff smooth.

For digital options, the grid shifting technique can be very efficient, and such “locking” of the location of strikes and barriers relative to the spatial grid can often reduce odd-even effects even better than the continuity correction discussed earlier. To demonstrate, consider the concrete task of using a finite difference grid to price a digital call option on a stock S in the Black-Scholes model. In this case, we conveniently have a theoretical option price to compare against, since it is easily shown that the time 0 value $V(0)$ must be

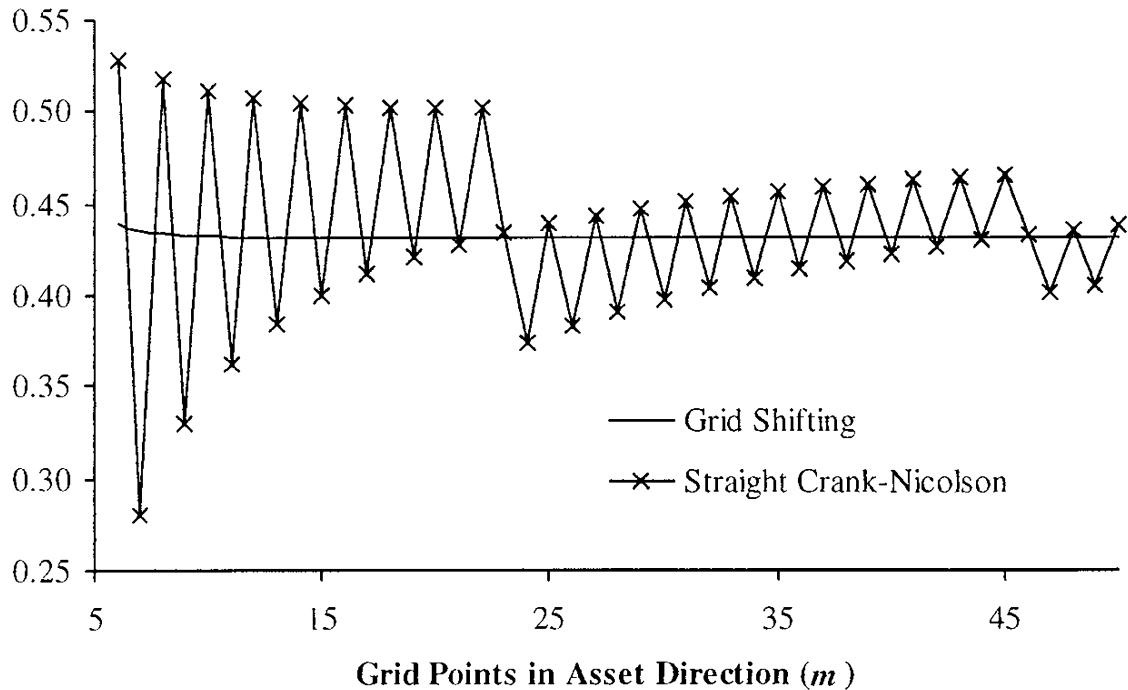
$$V(0) = e^{-rT} Q(S(T) > H) = e^{-rT} \Phi \left(\frac{\ln(S(0)/H) + (r - \sigma^2/2)T}{\sigma\sqrt{T}} \right). \quad (2.33)$$

For our numerical work, we discretize the asset equidistantly in log-space (i.e., we work with the PDE (2.3)) and determine the spatial grid boundaries by probabilistic means using a multiplier of $\alpha = 4.5$, see Section 2.1. Spatial boundary conditions are $\partial V / \partial x = \partial^2 V / \partial x^2$, implemented as described in Section 2.2.2. In one experiment, we apply a straight Crank-Nicolson approach, with no attempt to regularize the payoff condition. In a second experiment, we combine Crank-Nicolson with Rannacher stepping and also nudge the entire spatial grid upwards until the log-barrier $\ln(H)$ is located exactly half-way between two spatial grid points. Numerical results are shown in Figure 2.1.

As Figure 2.1 shows, a naive Crank-Nicolson implementation is plagued by severe odd-even effects and very slow convergence — 100’s of spatial steps appear to be necessary before acceptable levels of the option price are reached. On the other hand, grid shifting combined with Rannacher stepping results in a perfectly smooth¹¹ convergence profile, and 5-digit price precision is here reached in less than 30 steps.

¹¹It can be verified that the convergence order in m is, as expected, close to 2 in this experiment.

Fig. 2.1. 3 Year Digital Option Price



Notes: Finite difference estimates for the Black-Scholes price of a 3 year digital option with a strike of $H = 100$. The initial asset price is $S(0) = 100$, the interest rate is $r = 0$, and the volatility is $\sigma = 20\%$. Time stepping is performed with an equidistant grid containing $n = 50$ points. Spatial discretization in log-space is equidistant, as described in the main text; the number of grid points (m) is as listed on the x -axis of the figure. The “Straight Crank-Nicolson” graph shows the convergence profile for a pure Crank-Nicolson finite difference grid. The “Grid Shifting” graph shows the convergence profile for a Crank-Nicolson finite difference grid with Rannacher stepping and a shift of the spatial grid to center $\ln(H)$ midway between two grid points. From (2.33), the theoretical value of the option is 0.4312451.

2.6 Convection-Dominated PDEs

Recall from Section 2.3 that stability of the explicit finite difference scheme requires that (omitting grid subscripts on μ and σ)

$$\frac{2\Delta_x^2}{\Delta_t} \sigma^2 \geq \sigma^4 + \mu^2 \Delta_x^2 + |\mu^2 \Delta_x^2 - \sigma^4|.$$

As discussed, this condition can be violated if Δ_t is too large relative to Δ_x . However, for fixed Δ_t and Δ_x we notice that instability can also be triggered if the absolute value of the drift μ is raised to be sufficiently large relative to the diffusion coefficient σ .

While theta schemes with $\theta \geq 1/2$ are always stable, large drifts in the PDE can nevertheless cause spurious oscillations and an overall deterioration in numerical performance of these schemes. PDEs for which this effect

occurs are said to be *convection-dominated*. To quantify matters, assume for simplicity that the finite difference grid is equidistant in the x -direction, and consider the matrix \mathbf{A} in (2.11) with tri-diagonal coefficients c , u , and l given by (2.7)–(2.9). As discussed in e.g. d’Halluin et al. [2005], spurious oscillations can occur when, for some t and some j , either $u_j(t) < 0$ or $l_j(t) < 0$. From (2.8) and (2.9), to avoid spurious oscillations we would thus need

$$\sigma(t, x_j)^2 \geq |\mu(t, x_j)|\Delta_x. \quad (2.34)$$

Intuitively, in convection-dominated systems, the central difference coefficient δ_x and δ_{xx} used to discretize the PDE can no longer fully contain the large expected up- or downward trend of the underlying process for x ; as a result, spurious oscillations can occur.

2.6.1 Upwinding

There are a number of well-established techniques to deal with convection-dominated PDEs. First, we can obviously attempt to lower Δ_x such that (2.34) is satisfied. This, however, may not be practical from a computational standpoint (and may require that Δ_t is lowered as well to avoid spurious oscillations originating from the time-stepping scheme). An alternative is to modify the first-order discrete operator δ_x such that it points in the direction of the large absolute drift. For instance, we can simply elect to use a suitably oriented one-sided difference, rather than a central difference, whenever (2.34) is violated. This procedure is known as *upstream differencing* or *upwinding*. To formalize the idea, introduce a new first-order difference operator δ_x^* given as

$$\delta_x^* V(t, x_j) = \begin{cases} \frac{1}{2} (V(t, x_{j+1}) - V(t, x_{j-1})) \Delta_x^{-1}, & |\mu(t, x_j)|\Delta_x \leq \sigma(t, x_j)^2, \\ (V(t, x_j) - V(t, x_{j-1})) \Delta_x^{-1}, & \mu(t, x_j)\Delta_x < -\sigma(t, x_j)^2, \\ (V(t, x_{j+1}) - V(t, x_j)) \Delta_x^{-1}, & \mu(t, x_j)\Delta_x > \sigma(t, x_j)^2. \end{cases}$$

Using δ_x^* instead of δ_x modifies the matrix \mathbf{A} in (2.11). Specifically, if $\mu(t, x_j)\Delta_x < -\sigma(t, x_j)^2$ we replace (2.7)–(2.9) with:

$$c_j(t) = \mu(t, x_j)\Delta_x^{-1} - \sigma(t, x_j)^2\Delta_x^{-2} - r(t, x_j), \quad (2.35)$$

$$u_j(t) = \frac{1}{2}\sigma(t, x_j)^2\Delta_x^{-2}, \quad (2.36)$$

$$l_j(t) = -\mu(t, x_j)\Delta_x^{-1} + \frac{1}{2}\sigma(t, x_j)^2\Delta_x^{-2}. \quad (2.37)$$

And when $\mu(t, x_j)\Delta_x > \sigma(t, x_j)^2$, we use

$$c_j(t) = -\mu(t, x_j)\Delta_x^{-1} - \sigma(t, x_j)^2\Delta_x^{-2} - r(t, x_j), \quad (2.38)$$

$$u_j(t) = \mu(t, x_j)\Delta_x^{-1} + \frac{1}{2}\sigma(t, x_j)^2\Delta_x^{-2}, \quad (2.39)$$

$$l_j(t) = \frac{1}{2}\sigma(t, x_j)^2\Delta_x^{-2}. \quad (2.40)$$

For non-equidistant grids, a similar modification to (2.29)–(2.31) is required. We omit the straightforward details.

Let us try to gain some further understanding of the upwind algorithm. Comparison of (2.35)–(2.40) with (2.7)–(2.9), shows that upwinding amounts to using a regular central difference operator δ_x on a PDE with a diffusion coefficient modified to be $\sigma(t, x) + \sqrt{|\mu(t, x)|\Delta_x}$. The numerical scheme in effect introduces enough artificial diffusion into the PDE to satisfy (2.34). Doing so, however, comes at a cost: the convergence order of the scheme will be reduced to $O(\Delta_x)$ if one-sided differencing ends up being activated in a significant part of the grid. We note that higher-order upwinding schemes are possible if the finite difference operator δ_x^* is allowed to act on more than three neighboring points. For such schemes, the matrix \mathbf{A} will no longer be tri-diagonal.

2.6.2 Other Techniques

As discussed earlier, upwinding amounts to adding numerical diffusion at nodes where the scheme is convection dominated. Alternatively, we can increase $\sigma(t, x)$ directly, to $\sigma(t, x) + \varepsilon$ where ε is chosen to be large enough for the scheme to satisfy (2.34). By solving the resulting PDE for different values of ε , it may be possible to determine how the error associated with ε scales in ε . This, in turn, will allow us to extrapolate to the limit $\varepsilon = 0$. See p. 135 of Tavella and Randall [2000] for an example.

The upwinding scheme presented in Section 2.6.1 switches abruptly from central differencing to one-sided differencing when the condition (2.34) is violated. In some schemes, the switch from central to one-sided differencing is made smooth by using a weighted average of a one-sided and a central difference operator. The weight on the central difference is close to one when $\sigma(t, x)^2 \gg |\mu(t, x)|\Delta_x$, but decreases smoothly to zero as $\sigma(t, x)^2/|\mu(t, x)|$ tends to zero. While it is unclear whether a smooth transition to upwinding is truly important (the convergence order is typically not improved over straight upwinding), Duffy [2000] suggests that the class of exponentially fitted schemes (see Duffy [2000] and Stoyan [1979]) may be quite robust in derivatives pricing applications.

In some finance applications, multi-dimensional PDEs might arise where $\sigma(t, x) = 0$ for one of the underlying variables; see for instance Section 2.7.5. While upwinding techniques still apply here, we note that specialized methods exist with better ($O(\Delta_x^2)$) convergence, should they become necessary. See, for instance, d'Halluin et al. [2005] for details on the so-called *semi-Lagrangian* methods.

2.7 Option Examples

In our discussion so far, we have assumed that options are characterized by a single terminal payoff function $g(x)$ and a set of spatial boundary

conditions determining the option price at the boundaries of the x -domain. In reality, many options are more complicated than this and may involve early exercise decisions, pre-maturity cash flows, path dependency, and more. In this section, we provide some relatively straightforward examples of such complications and how to modify the basic finite difference algorithm to deal with them. More examples will be provided later, in the context of specific fixed income securities.

2.7.1 Continuous Barrier Options

We have already touched upon the concept of an up-and-out knock-out option, an option that expires worthless if the x -process ever rises above a critical level H . As we described, we here must simply solve the PDE (2.1) on a domain $[\underline{M}, H]$, where \underline{M} represents the lowest attainable value of the process $x(t)$ on $[0, T]$. The boundary condition at the upper boundary is then dictated to be $V(t, H) = 0$, i.e. of the Dirichlet type. We can generalize this to allow both “up” and “down” type barriers, and to perhaps give a non-zero payout (a “rebate”) at the time the barrier(s) are hit (provided this happens before the option maturity). Specifically, if we have a lower barrier at \underline{H} , an upper barrier of \bar{H} , a time-dependent lower rebate function of $\underline{f}(t)$, and a time-dependent upper rebate function of $\bar{f}(t)$, we must dimension our spatial grid $\{x_j\}_{j=0}^{m+1}$ to have $x_0 = \underline{H}$, $x_{m+1} = \bar{H}$, and we then simply impose the Dirichlet boundary conditions $V(t, x_0) = \underline{f}(t)$ and $V(t, x_{m+1}) = \bar{f}(t)$. See (2.10) and the definition of Ω for the algorithm required to incorporate such Dirichlet boundary conditions into the finite difference scheme.

In practice, barrier options sometimes involve time-dependent barriers, possibly with discontinuities. For instance, *step-up* and *step-down* barrier options will have piecewise flat barriers that increase (step-up) or decrease (step-down) at discrete points in time. Extension of the finite difference algorithm to cover step-up and step-down options is relatively straightforward. As an illustration, consider a zero-rebate up-and-out single-barrier option where the (upper) barrier is flat, except for a discontinuous change at time $T^* < T$, at which point the barrier moves from a value of H^* to a value of H , with $H > H^*$. We set the x -domain of our finite difference grid to $x \in [\underline{M}, H]$, with \underline{M} a probabilistic lower limit, as defined above; accordingly, our spatial grid would be $\{x_j\}_{j=0}^{m+1}$, where $x_0 = \underline{M}$ and $x_{m+1} = H$. In preparation for the shift in barrier levels at time T^* , we make sure that one level in the spatial grid — say x_{k+1} , $k < m$, — is set exactly at the level H^* . Similarly, we make sure that one level in the time grid is set exactly to T^* . Starting at time T , we then iterate backwards in time by repeated solution of m -dimensional tri-diagonal systems of equations, at each step integrating a prescribed rebate function by supplying the Dirichlet boundary condition $V(t, x_{m+1}) = 0$. The moment we hit T^* , the PDE now only applies to the smaller region $[\underline{M}, H^*]$, covered by the reduced spatial grid $\{x_j\}_{j=0}^{k+1}$ with

$x_{k+1} = H^*$. From T^* back to time 0, the backward induction algorithm then involves only k -dimensional tri-diagonal systems of equations, with the Dirichlet boundary condition $V(t, x_{k+1}) = 0$. Spatial nodes above x_{k+1} correspond to zero option value and can be ignored¹². Modification of the algorithm outlined above to handle more than two barrier discontinuities is straightforward.

We can extend our definition of barrier options even further by making the topology of “alive” and “dead” regions more complicated. At time t , assume for instance that the PDE applies in an “alive” region of $x \in L(t)$ and a rebate function $R(t, x)$ that applies in the “dead” region $D(t) = \mathcal{B} \setminus L(t)$. Assume that we discretize the problem on a single rectangular finite difference grid spanning the spatial domain $[\underline{M}, \bar{M}]$, where \underline{M} and \bar{M} are set such that the alive regions are covered, up to probabilistic limits (if necessary). Given option values at time t_{i+1} , we then only need to run the basic matrix equation (2.18) for values in our grid $\{x_j\}$ that lie inside $L(t_i)$. This requires scaling down the dimension of the matrix \mathbf{A} as needed, and providing the relevant boundary conditions (given through $R(t_i, x_j)$) at the boundary (or boundaries) of $L(t_i)$. The parts of the spatial grid that lie outside of $L(t_i)$ can be directly filled in with values provided by the rebate function R . Notice that, if possible, the spatial grid should be set such that the boundaries of $L(t_i)$ are contained in the mesh; this will likely require us to use the techniques outlined in Section 2.4.

If the alive region has the simple form $L(t) = [\alpha(t), \beta(t)]$ for smooth deterministic functions α and β , an alternative to the scheme above is to introduce a time-dependent transformation that straightens out the barriers, allowing us to return to the standard finite difference setup where the PDE applies to a single rectangular (t, x) domain. One possible transformation involves using a spatial variable of

$$y = y(t, x) = \frac{x - \alpha(t)}{\beta(t) - \alpha(t)}, \quad (2.41)$$

which transforms the curved x -barriers $\alpha(t)$ and $\beta(t)$ into flat y -barriers at $y = 0$ and $y = 1$, respectively. The linearity of the transformation (2.41) makes it easy to work with; see Tavella and Randall [2000] for details and a discussion of extensions to multi-dimensional PDEs and to barriers with discontinuities.

¹²An obvious twist to the algorithm involves using different spatial grids over $[0, T^*]$ and $[T^*, T]$, allowing for more flexibility in node placement. In this case, values computed by backward induction must, at time T^* , be interpolated from one x -grid to another. The interpolation rule should be at least third-order accurate; see the discussion in Section 2.7.3.

2.7.2 Discrete Barrier Options

The barrier options considered in Section 2.7.1 are continuously monitored, in the sense that the barrier condition is observed for all times in a given interval. In practice, monitoring the barrier condition continuously can be impractical, and it may instead only be imposed on a discrete set of dates $T_1 < T_2 < \dots < T_K$, with $T_K \leq T$ and $T_1 > 0$. For the sake of concreteness, let us consider a discretely monitored up-and-out option with a constant barrier H . For a continuously monitored up-and-out barrier option it would suffice to solve the PDE on a domain $x \in [\underline{M}, H]$, where \underline{M} is a probabilistic lower limit. This is, however, no longer the case for a discretely monitored option where we need to allow the value function to “diffuse” above the barrier levels between dates in the monitoring set $\{T_k\}_{k=1}^K$. To allow for this, we discretize the PDE on a larger domain $x \in [\underline{M}, \bar{M}]$, $\bar{M} > H$. We can determine \bar{M} probabilistically by determining a confidence interval for how far above the barrier $x(t)$ can rise between monitoring dates. For instance, for the Black-Scholes PDE (2.3), assume that $\max_{k=2,\dots,K} (T_k - T_{k-1}) = \Delta_T$. Conditioned on $x(t) = H$, the probability that $x(t + \Delta_T)$ exceeds

$$x_\alpha = H + \left(r - \frac{1}{2}\sigma^2 \right) \Delta_T + \alpha\sigma\sqrt{\Delta_T}$$

is $\Phi(-\alpha)$. As in Section 2.1, we recommend setting $\bar{M} = x_\alpha$ for values of α somewhere between 3 and 5. To properly capture diffusion between barrier observation dates, we should also dimension the time grid of the finite difference scheme such that multiple time steps (at least two or three, say) are taken between observation dates. All observation dates $\{T_k\}_{k=1}^K$ should obviously be contained in the time grid.

Between barrier observation dates, we solve our PDE by the standard finite difference algorithm outlined in Section 2.2.4, as always imposing either an asymptotic Dirichlet condition at $x = \bar{M}$ or a condition on the x -derivatives of the value function. At each barrier observation time T_k , we must impose a *barrier jump condition*

$$V(T_k-, x) = V(T_k+, x)1_{\{x < H\}}, \quad k = 1, \dots, K, \quad (2.42)$$

where the notation $T_k \pm$ was introduced in Section 1.10.1 to denote the limit $T_k \pm \varepsilon$ for $\varepsilon \downarrow 0$. This merely states that all values $V(T_k, x)$ are zero for $x \geq H$, consistent with the definition of an up-and-out option. In our finite difference scheme, we incorporate this jump condition by simply interpreting the vector $\widehat{\mathbf{V}}(T_k)$ as found by regular backward induction as $\widehat{\mathbf{V}}(T_k+)$ and then replacing

$$\widehat{\mathbf{V}}(T_k+) = \left(\widehat{V}_1(T_k+), \dots, \widehat{V}_m(T_k+) \right)^\top$$

with

$$\widehat{\mathbf{V}}(T_k-) = \left(\widehat{V}_1(T_k+) \mathbf{1}_{\{x_1 < H\}}, \dots, \widehat{V}_m(T_k+) \mathbf{1}_{\{x_m < H\}} \right)^\top$$

before continuing the algorithm backwards from T_k .

The jump condition (2.42) will generally produce a discontinuity in V as a function of x , around the barrier level H . If we use Crank-Nicolson time-stepping, it will then be prudent to employ a fully implicit scheme for the first few backwards time steps (Rannacher stepping) past each barrier observation date T_k . As discussed in Section 2.5, ideally this should be combined with a smoothing algorithm acting on $\widehat{\mathbf{V}}(T_k-)$ or, perhaps more conveniently, a shift of the spatial grid such that H lies exactly mid-way between two spatial nodes in the grid.

We round off by noting that the discussion above for an up-and-out option easily extends to more complicated discrete barrier options, including those with time-varying barrier levels and rebates. For instance, assume that an option involves upper and lower time-varying barriers of $\overline{H}(t)$ and $\underline{H}(t)$, respectively, as well as a time- and state-dependent rebate of $R(t, x)$. In this case, we simply replace the jump condition (2.42) with

$$\begin{aligned} V(T_k-, x) &= V(T_k+, x) \mathbf{1}_{\{\underline{H}(T_k) < x < \overline{H}(T_k)\}} \\ &\quad + R(T_k, x) \left(\mathbf{1}_{\{x \geq \overline{H}(T_k)\}} + \mathbf{1}_{\{x \leq \underline{H}(T_k)\}} \right), \end{aligned}$$

and otherwise proceed as above. We note that time-dependent barriers will typically require flexibility in setting the spatial grid, as there are now multiple critical x -levels to consider. The discretization in Section 2.4 can obviously assist with this.

2.7.3 Coupon-Paying Securities and Dividends

Many fixed-income securities are coupon-bearing and involve periodic transfer of a cash amount between the buyer and the seller. This can easily be incorporated into a finite difference grid, through a jump condition. Specifically, consider a security that pays its owner a single cash amount of $p(T^*, x)$ at time $T^* < T$, where p is a deterministic function $p : [0, T] \times \mathcal{B} \rightarrow \mathbb{R}$. We dimension our time grid such that T^* is contained in the grid, and then apply at time T^* the condition

$$V(T^*-, x) = V(T^*+, x) + p(T^*, x). \quad (2.43)$$

This simply expresses that V will decrease by an amount p immediately after p is paid (and thereby no longer contained in V). In a finite difference algorithm, (2.43) is incorporated by replacing $\widehat{\mathbf{V}}(T^*+)$, as found by regular backward induction, with

$$\widehat{\mathbf{V}}(T^*-) = \left(\widehat{V}_1(T^*+) + p(T^*, x_1), \dots, \widehat{V}_m(T^*+) + p(T^*, x_m) \right)^\top$$

before continuing the algorithm backwards from T^* . Extensions to multiple coupons are trivial.

In some cases a derivative security does not itself pay coupons, but is written on a security that does. This involves no particular complications, except for the case where payments may affect the state variable underlying the PDE. For instance, consider the classical case of a stock paying a dividend: at the time of the dividend payment, the stock jumps down by an amount equal to the dividend payment. For a model that uses the stock price (or a transformation of the stock price) as the state variable x , a dividend payment at time T^* would thus be associated with a discontinuity in the state variable, $x(T^*+) = x(T^*-) - d(T^*, x(T^* -))$, where d is the magnitude of the jump¹³. As long as the dividend-payment does not come as a surprise (i.e., at a random time), it must already be contained into the option price at $T^* -$, and will have no price effect as we move forward from $T^* -$ to $T^* +$. We can express this continuity restriction through yet another jump condition

$$V(T^* -, x) = V(T^* +, x - d(T^*, x)). \quad (2.44)$$

See Wilmott et al. [1993] for more discussion. Implementation of (2.44) in a finite difference grid proceeds as follows. First, we use regular backward induction to establish

$$\begin{aligned} \widehat{\mathbf{V}}(T^* +) &= \left(\widehat{V}_1(T^* +), \dots, \widehat{V}_m(T^* +) \right)^\top \\ &= \left(\widehat{V}(T^* +, x_1), \dots, \widehat{V}(T^* +, x_m) \right)^\top. \end{aligned}$$

Then we write

$$\widehat{\mathbf{V}}(T^* -) = \left(\widehat{V}(T^* +, x_1 - d(T^*, x_1)), \dots, \widehat{V}(T^* +, x_m - d(T^*, x_m)) \right)^\top.$$

The values $\widehat{V}_j(T^* +, x_j - d)$ here can be found by interpolation in the x -direction on the $\widehat{\mathbf{V}}(T^* +)$ -array. As shown in Tavella and Randall [2000], the order of the interpolator should be strictly higher than two, to avoid inducing spurious numerical diffusion into our θ -style finite difference schemes. We note that this rules out the piecewise linear interpolation rule proposed in Wilmott et al. [1993]. A common choice is to use cubic spline interpolation; see Chapter 6 for much information on cubic splines.

2.7.4 Securities with Early Exercise

In Section 1.10 we introduced the concept of Bermudan and American securities with early exercise features. Under the assumption that exercise

¹³To prevent negative stock prices, it may be necessary to truncate the size of d locally in the finite difference grid. For simplicity, we ignore this complication here.

values are determined by a deterministic function¹⁴ $h(t, x)$, $h : [0, T] \times \mathcal{B} \rightarrow \mathbb{R}$, finite difference grids are ideal for pricing of such securities. Let us first consider a Bermudan option with exercise opportunities restricted to the finite set $\{T_k\}_{k=1}^K$. The Bellman principle (1.67) in Section 1.10 can, as shown there, be expressed as a simple jump condition

$$V(T_k-, x) = \max(V(T_k+, x), h(T_k, x)), \quad k = 1, \dots, K, \quad (2.45)$$

which can be incorporated into a finite difference solver precisely the same way as in previous sections. The condition (2.45) will result in a kink in the value function around the level of x at which we shift from the hold region into the exercise region. If Crank-Nicolson time-stepping is used, one should ideally apply smoothing on the finite difference value vector $\hat{\mathbf{V}}(T^* -)$, particularly around the kink.

If exercise can take place continuously (that is, American-style) on a given time interval, a crude way to incorporate this into a finite difference grid is by simply applying (2.45) to every point in the time grid of the finite difference scheme. By not specifically imposing the partial differential inequalities (see Section 1.10.1), this algorithm, however, will generally only be accurate to first order in the time step, even if a Crank-Nicolson scheme is used; see Carverhill and Clewlow [1990] for a proof. As American-style exercise is rarely used in fixed income markets, we shall not pursue this issue further but just point out that a number of schemes exist to restore second-order time convergence to finite difference pricing of American options, see, e.g., Forsyth and Vetzal [2002].

2.7.5 Path-Dependent Options

Finite difference methods are normally limited to Markovian problems where dynamics are characterized by SDEs and where payouts are simple deterministic functions of the underlying state variables. A number of options, however, have terminal time T payouts that depend not only on the state of x at time T , but on the entire path $\{x(t), t \in [0, T]\}$. In general, such options must be priced by Monte Carlo methods (see Chapter 3), but exceptions exist. Indeed, barrier and American options can be considered path-dependent options, yet, as we have seen, can still be priced in a finite difference grid. Even stronger path-dependence can sometimes be handled, through the introduction of new state variables to the PDE.

To give an example, consider a path-dependent contract where the terminal payout at time T can be written as

¹⁴If h represents the value of a derivative security that has no closed-form pricing formula, it may be necessary to estimate this function by backward induction in the finite difference grid itself. Such a “preprocessing” step is typically straightforward to execute.

$$V(T) = g(x(T), I(T)), \quad (2.46)$$

where I is a path integral of the type

$$I(t) = \int_0^t h(x(s)) ds, \quad (2.47)$$

for some deterministic function h . For instance, if $h(x) = x$, we say that the option is a continuously sampled *Asian option*.

For the payout (2.46) we have $V(t) = V(t, x(t), I(t))$ where $x(t)$ satisfies the SDE (2.2) and

$$dI(t) = h(x(t)) dt, \quad I(0) = 0.$$

From the backward Kolmogorov equation, it follows that $V(t, x, I)$ solves

$$\frac{\partial V}{\partial t} + \mu(t, x) \frac{\partial V}{\partial x} + \frac{1}{2} \sigma(t, x)^2 \frac{\partial^2 V}{\partial x^2} + h(x) \frac{\partial V}{\partial I} = r(t, x)V, \quad (2.48)$$

subject to the terminal condition $V(T, x, I) = g(x, I)$. There are several complications with this PDE. First, it involves *two* spatial variables, x and I , requiring the use of a two-dimensional PDE solver. Second, the PDE contains no second-order derivative in the variable I , i.e. it is convection dominated in the I -direction. We have discussed methods to handle the latter issue in Section 2.6.1 and will turn to address the former in Section 2.9. Another complication is the fact that the term $h(x)$ multiplying $\partial V / \partial I$ may be of a different order of magnitude than the other coefficients in (2.48), increasing the difficulty of solving the equation numerically. We refer to Zvan et al. [1998] for a more detailed discussion of PDEs of the type (2.48).

In practice, it is rare that a continuous-time integral such as (2.47) is used in an option payout. Instead, one normally samples the function $h(x(t))$ only on a discrete set of dates, i.e. we replace $I(T)$ with

$$I(T) = \sum_{i=1}^n h(x(T_i)) (T_i - T_{i-1}),$$

where $T_0 < T_1 < \dots < T_n$ is a discrete schedule, with $T_0 = 0$ and $T_n = T$. Informally, we now have

$$dI(t) = \delta(T_i - t) \cdot h(x(T_i)) (T_i - T_{i-1}), \quad I(0) = 0, \quad (2.49)$$

where $\delta(\cdot)$ is the Dirac delta function. In a PDE setting, we incorporate a process such as (2.49) through appropriate jump conditions, writing

$$V(T_i-, x, I) = V(T_i+, x, I + h(x)(T_i - T_{i-1})). \quad (2.50)$$

In the same fashion as for discrete dividends (Section 2.7.3), the jump condition enforces continuity of the option price across the dates where I

gets updated. The condition is applied at each date in the discrete schedule, $i = 1, \dots, n$; in between schedule dates (where now $dI(t) = 0$), we solve the PDE

$$\frac{\partial V}{\partial t} + \mu(t, x) \frac{\partial V}{\partial x} + \frac{1}{2} \sigma(t, x)^2 \frac{\partial^2 V}{\partial x^2} = r(t, x)V,$$

which has no term involving I . When the I -direction is discretized in, say, m_I different values, the solution scheme thus involves solving m_I different *one-dimensional* PDEs backward in time; the solutions of these m_I PDEs exchange information with each other at each date in the schedule, in accordance with (2.50). As was the case for cash dividends, implementation of (2.50) will normally require support from an interpolation scheme, to align the (x -dependent) jumps in I with the knots of the discretized I -grid used in the finite difference scheme. See, e.g., Zvan et al. [1999] or Wilmott et al. [1993] for further details. An application of this idea in the context of interest rate derivatives is given in Section 18.4.5.

On rare occasions — basically when the homogeneity condition $V(\eta x, \lambda I, t) = \lambda^\eta V(x, I, t)$, $\lambda, \eta > 0$, holds — it is possible to make a change of variables or a change of probability measure that will reduce (2.48) or its discrete-time version to a one-dimensional PDE; see e.g. Rogers and Shi [1995] or Andreasen [1998] for the case of various Asian options. Section 18.4.5 demonstrates one such method, sometimes called the method of *similarity reduction*, for pricing of “weakly path-dependent” securities, including certain callable interest rate derivatives where the notional accretes at a stochastic coupon rate (see Section 5.14.5 for definitions).

2.7.6 Multiple Exercise Rights

Certain financial products with early exercise rights allow the holder to exercise more than once. Such “multi-exercise” options are relatively rare, but the so-called *chooser cap* (also known as a *flexi-cap*) is occasionally traded and constitutes a good example for describing how to handle multi-exercise options in a PDE setting. Let there be given a set of L possible exercise dates, $T_1 < T_2 < \dots < T_L$, and assume that we have the right to exercise no more than l times, with $l < L$. Provided that we exercise at time T_i , in a chooser cap we are paid¹⁵ $(S(T_i) - K)^+$, where $S(\cdot)$ is some interest rate index and K is the strike. Clearly, we would never exercise at time T_i unless $S(T_i) > K$, but how much larger than K the rate $S(T_i)$ needs to be to trigger optimal exercise is not obvious, and must at least depend on i) how many of our l exercise opportunities we have already used up at time T_i ; and ii) how much value is lost by using (rather than postponing) one of the remaining exercise opportunities.

¹⁵We have ignored a day count scaling constant in the payout. Also, in most cases payment takes place at time T_{i+1} , rather than at T_i ; such a payment delay can be handled by a discount operation.

While the question of how to exercise optimally on a chooser cap may appear quite complex, it is surprisingly easy to implement in a finite difference setting by combining techniques from Sections 2.7.4 and 2.7.5 above. The key to the method is to introduce an additional state variable I to keep track of how many exercise opportunities are left. Assume that all interest rates are functions of a Markov state variable $x(\cdot)$, and let therefore $V(t, x, I)$ denote the value of the chooser cap at time t , given $x(t) = x$ and given that there are still I exercise opportunities left. Notice that the variable I can only take $l + 1$ distinct values: $0, 1, \dots, l$; notice also that $V(t, x, 0) = 0$ for all t and x , since $I = 0$ corresponds to the situation where there are no exercise opportunities left. Additionally, at the terminal time T_L we clearly have

$$V(T_L, x, I) = (S(T_L, x) - K)^+, \quad I = 1, 2, \dots, l, \quad (2.51)$$

where we have written $S(T_L) = S(T_L, x)$ to emphasize the deterministic dependence of S on the state variable x .

For given dynamics of $x(t)$, starting with the terminal conditions in (2.51), we may roll the l different value functions $V(\cdot, x, I)$, $I = 1, 2, \dots, l$, back through time in standard finite difference manner. At each time T_i , $i = 1, \dots, L - 1$, jump conditions similar to (2.45) must be applied, for all $I = 1, 2, \dots, l$:

$$V(T_{i-}, x, I) = \max(V(T_i+, x, I), V(T_i+, x, I - 1) + (S(T_i, x) - K)^+).$$

Notice that these conditions simply express that exercise is optimal only if the exercise value (the cap payout plus the value of a chooser cap with one less exercise opportunity) exceeds the hold value (the non-exercised chooser cap). Once we have rolled all the way back to $t = 0$, the chooser cap value at time $t = 0$ may be identified as $V(0) = V(0, x(0), l)$.

We should note that the “chooser” or “flexi” feature can be added to securities other than caps (and floors). For instance, in Section 19.5 we study the so-called *flexi-swap*, another security with multiple embedded exercise rights.

2.8 Special Issues

In this section, we briefly show a few techniques that may come in handy for certain applications.

2.8.1 Mesh Refinements for Multiple Events

As discussed in Section 2.1, the domain of the state variable x is often determined as an exact or approximate confidence interval for the random variable $x(T)$, where T is the final time of interest for a particular valuation problem we want to solve. Given the number of desired spatial steps in the

scheme, the discretization step in x -direction is then obtained by dividing the size of the confidence interval by the number of steps. Similarly, the discretization step in t -direction is typically obtained by dividing T by the number of desired time steps. This is a standard procedure for building a simple rectangular mesh, and it works well if the derivative we wish to value does not have any “interesting” features between the valuation time 0 and the final time T (e.g., for a simple European option). However, as should be evident from the examples in Section 2.7, many real-life derivative securities are characterized by a multitude of events during their lifetimes, all of which must be adequately captured in the PDE scheme. It is not hard to see that a grid dimensioning scheme based solely on the last event date may yield inappropriate mesh resolution at earlier dates.

To make the discussion above concrete, let us consider the example of a Bermudan option (see Section 2.7.4) with two exercise dates, T_1 and T_2 . Assume that $0 < T_1 \ll T_2$, i.e. that the first exercise date is much closer to the valuation date than the second (and last) one. Also assume that there is a decent chance that the option actually will be exercised at time T_1 , making it important to capture to good precision the value of the option expiring at T_1 . Now, if we build our mesh based only on the distribution of the state variable $x(T_2)$ at time T_2 , there would typically be too few t -points in the interval $[0, T_1]$. Also, the x -direction discretization step would be too large compared to the range of possible values of the state variable $x(T_1)$ at time T_1 , i.e. the x -grid would be too coarse for the process $x(\cdot)$ on the time interval $[0, T_1]$. Both issues would typically lead to a large discretization error in the finite difference stepping of the option over the time period $[0, T_1]$, leading to problems with accuracy in values and risk sensitivities.

The issue of the sparsity of the time grid is fairly easy to deal with, as we are free to add extra points to the time grid before time T_1 . This by itself, however, will not solve precision problems, as the space step remains large. Any proper solution should, of course, come in the form of refining both the t - and x -grids at the same time.

One possible way of refining the x -discretization is to abandon the usage of a single rectangular (t, x) -domain, and instead link together different equidistant rectangular meshes for different periods in the life of the derivative. These mesh “blocks” would generally increase in spatial width with time and would connect to each other via an interpolation scheme. To be more specific, let us assume, as in Section 2.1, that the state variable $x(\cdot)$ is the logarithm of the stock in the Black-Scholes model and is given by (2.4), with the PDE to solve given by (2.3). We extend our simple two-period example above to a derivative with K times of interest, $0 < T_1 < \dots < T_K$; these times could be specified as an additional input into valuation, or derived from the trade description (e.g. they could represent the exercise dates for a Bermudan option, or the knock-out dates for the discretely-monitored barrier option of Section 2.7.2). Suppose we are given values of m and n , and now wish to construct the mesh for the time period $[T_{k-1}, T_k]$, by using

the same time and space steps Δ_t^k , Δ_x^k as would be used in the standard scheme of Section 2.1 for a derivative security with the terminal payoff at T_k . That is, having fixed the cutoff α we would set

$$\Delta_t^k = T_k/n, \quad \Delta_x^k = 2\alpha\sigma\sqrt{T_k}/(m+1). \quad (2.52)$$

Then the rectangular, equidistant mesh for the time period $[T_{k-1}, T_k]$ is given by

$$\{t_i^k\}_{i=0}^{\lfloor(1-T_{k-1}/T_k)n\rfloor} \times \{x_j^k\}_{j=0}^{m+1}, \quad t_i^k = T_{k-1} + i\Delta_t^k, \quad x_j^k = x_{\min}^k + j\Delta_x^k, \quad (2.53)$$

where $\lfloor \cdot \rfloor$ denotes the integer part of a real number and (see (2.4))

$$x_{\min}^k = x(0) + \left(r - \frac{1}{2}\sigma^2\right)T_k - \alpha\sigma\sqrt{T_k}. \quad (2.54)$$

Note that in reality we would want to make sure that the point T_k is also in the mesh for the time period $[T_{k-1}, T_k]$, even though for simplicity of notations we did not reflect it in (2.53). It is also useful to note that the total number of time points is not going to be n , but is actually equal to

$$\sum_{k=1}^K \lfloor(1 - T_{k-1}/T_k)n\rfloor,$$

which scales linearly with n . Clearly, if exactly n points were required, a simple adjustment to the definition of the time step in (2.52) could be applied.

With a mesh as defined above, when arriving at time T_k in a backward induction scheme the solution $V(T_k, \cdot)$ would be discretized on the x -grid $\{x_j^{k+1}\}_{j=0}^{m+1}$. To solve the PDE backwards over the time period $[T_{k-1}, T_k]$, we would need to resample it on the different x -grid $\{x_j^k\}_{j=0}^{m+1}$. As with interpolation across dividends (Section 2.7.3), simple cubic interpolation would be a good choice here. Specifically, one would fit a cubic spline to the values $V(T_k, x_j^{k+1})$, $j = 0, \dots, m+1$, and then calculate $V(T_k, x_j^k)$, $j = 0, \dots, m+1$, by valuing the spline at the required grid points.

The “interpolated mesh” scheme above is rather intuitive and straightforward, but it does suffer from the need to do interpolation work that could slow down the PDE (especially in dimensions higher than 1 and/or for a large number of interface points K). Also, it is not entirely clear how interpolation will affect stability and convergence properties of the PDE. Finally, linking the interface mesh geometry to the trade specifics (such as exercise dates) may not be ideal from the point of view of designing an efficient valuation flow in a risk management system. These considerations lead us to an alternative approach that relies on *non-equidistant* discretization as developed in Section 2.4. The idea of this method is to use non-uniform

discretization to concentrate more points, both in time and space, around the initial point $t = 0$, $x = x(0)$. Clearly many ways of achieving this are possible — below we present a simple scheme we have used with good results.

We define K , the user input, to be the number of spatial refinement levels (with $K = 2$ or 3 typically used), and τ , another user input, to be a time scaling constant (typically $\tau = 4$). If T is the final horizon for valuation, we then introduce times

$$0 = T_0 < T_1 < \dots < T_K = T$$

by

$$T_k = \frac{T}{\tau^{K-k}}, \quad k = 1, \dots, K.$$

Then, the time grid for the time period $[T_{k-1}, T_k]$ is given by uniformly distributing $\tilde{n} \triangleq \lfloor n/K \rfloor$ points¹⁶ over $[T_{k-1}, T_k]$, i.e. is given by $\{t_i^k\}_{i=0}^{\tilde{n}}$ with

$$t_i^k = T_{k-1} + i \frac{T_k - T_{k-1}}{\tilde{n}} = T_{k-1} + i \frac{T_k - T_{k-1}}{\lfloor n/K \rfloor}.$$

(Note that we can use this specification with the interpolated mesh as well, instead of the time grid definition in (2.53)). The fact that the width of the intervals $[T_{k-1}, T_k]$ grow with k means that the time grid is more finely spaced in the beginning of the interval $[0, T]$ than at the end.

The x -grid we are going to define will be universal — i.e. the same for all time steps on the whole time interval $[0, T]$ — and non-uniform. To construct it, we first define a set of nested x -subdomains $[x_{\min}^k, x_{\max}^k]$, with x_{\min}^k defined by (2.54) and x_{\max}^k defined accordingly, i.e.

$$x_{\max}^k = x(0) + \left(r - \frac{1}{2}\sigma^2 \right) T_k + \alpha\sigma\sqrt{T_k},$$

for $k = 0, \dots, K$. Then we define step sizes by

$$\Delta_x^k = \frac{x_{\max}^k - x_{\min}^k}{\tilde{m} + 1}, \quad \tilde{m} = \left\lfloor \frac{m}{K} \right\rfloor.$$

The x -grid is then constructed by distributing grid points uniformly in subintervals $[x_{\min}^k, x_{\min}^{k-1}]$ and $[x_{\max}^{k-1}, x_{\max}^k]$ with the space step Δ_x^k , and is given by

$$\left(\bigcup_{k=1}^K \left\{ x_j^{\min,k} \right\}_{j=0}^{m^k+1} \right) \cup \left(\bigcup_{k=1}^K \left\{ x_j^{\max,k} \right\}_{j=0}^{m^k+1} \right),$$

where

$$x_j^{\min,k} = x_{\min}^k + j\Delta_x^k, \quad x_j^{\max,k} = x_{\max}^{k-1} + j\Delta_x^k,$$

¹⁶And, as advised earlier, adding trade event dates that fall into this period — although we do not reflect this in our notations for simplicity.

and

$$m^k = \left\lfloor \frac{x_{\min}^{k-1} - x_{\min}^k}{\Delta_x^k} \right\rfloor - 1 = \left\lfloor \frac{x_{\max}^k - x_{\max}^{k-1}}{\Delta_x^k} \right\rfloor - 1.$$

This distribution of space points results in an x -grid that is more dense around the point $x = x(0)$ than at the edges. It is worth noting that with only one refinement level $K = 1$, the standard rectangular uniform mesh sized by the terminal distribution of the state variable is recovered.

2.8.2 Analytics at the Last Time Step

In cases where the dynamics of underlying PDE variables are tractable, one naturally wonders whether finite difference methods could somehow be improved by incorporating analytical results into the scheme. Here, and in the next section, we discuss two simple ideas.

Suppose that we are faced with the problem of pricing a contingent claim with terminal boundary condition $g(x(T))$, where $x(t)$ is a Markovian process with known Arrow-Debreu state prices:

$$G(t, x; s, y) = \mathbb{E}^Q \left(\delta(x(s) - y) e^{-\int_t^s r(u, x(u)) du} | x(t) = x \right), \quad s > t.$$

Assume also that the claim in question involves a jump condition at time $0 < T^* < T$ (but no jump conditions between T^* and T). If our finite difference grid is $\{x_j\}_{j=0}^{m+1}$, we can now use a series of $m + 2$ outright convolutions to compute

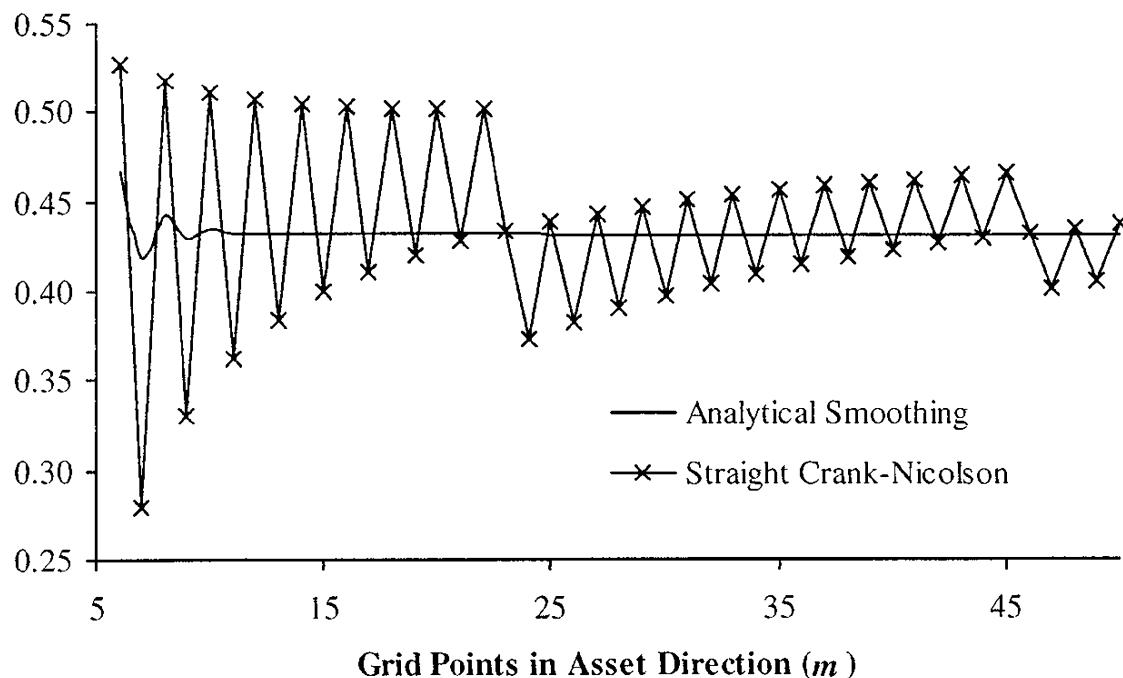
$$V(T^*, x_j) = \int_{\mathbb{R}} G(T^*, x_j; T, y) g(y) dy, \quad j = 0, 1, \dots, m + 1. \quad (2.55)$$

If we are lucky (i.e., if both g and G are sufficiently simple), then the integral on the right-hand side may be known in closed form for all values of x_j . If not, we can always perform a series of numerical integrations, the total cost of which is typically¹⁷ $O(m^2)$, i.e. more expensive than the typical $O(m)$ cost of a single time step in a finite difference method. There are several reasons why we may want to perform the numerical integrations nevertheless. First, the convolution expression (2.55) is exact, as it is based on the true transition density. Second, if the gap between T^* and T is large, an ordinary finite difference grid would need to roll back from T to T^* using multiple time steps n^* , at a total cost of $O(n^* m)$; if n^* is of the same magnitude as m , the computational effort of the convolution scheme would be comparable to that of a finite difference grid. Third, for discontinuous payouts, the integration in (2.55) will have a naturally smoothing effect, similar to (but often better than) the continuity correction method of Section 2.5.2. The smoothing

¹⁷There are exceptions. For instance, if fast Fourier transform (FFT) methods are applicable, the cost may be reduced to $O(m \ln(m))$. See Section 8.4 for details.

effect is discussed in more detail in Section 23.2.4 and is also demonstrated below, in Figure 2.2, where we have continued our investigation of the 3 year digital option considered earlier in Section 2.5.3. Since the model used in Figure 2.2 is ordinary Black-Scholes and $g(x) = 1_{\{x>H\}}$, the integrals in (2.55) can here be computed in closed form from (2.33).

Fig. 2.2. 3 Year Digital Option Price



Notes: Finite difference estimates for the Black-Scholes price of a 3 year digital option. All contract and model parameters are as in Figure 2.1. Time stepping is performed with an equidistant grid containing $n = 50$ points. Spatial discretization in log-space is equidistant, as described in the main text; the number of grid points (m) is as listed on the x -axis of the figure. The “Straight Crank-Nicolson” graph shows the convergence profile for a pure Crank-Nicolson finite difference grid. The “Analytical Smoothing” graph shows the convergence profile for a Crank-Nicolson finite difference grid starting at $T^* = 2.5$ years, with the terminal boundary condition set equal to a 0.5 year digital option price (as in (2.55)). The theoretical value of the option is 0.4312451.

In principle, we could continue rolling back from T^* (through, possibly, jump conditions at earlier times) by performing convolutions, rather than solving finite difference grids. In practice, this rarely leads to improvements over a finite difference grid, unless the densities and payoffs are quite simple¹⁸. Moreover, in many cases we may not have *exact* Arrow-Debreu

¹⁸For simple densities (especially Gaussian), special-purpose methods exist to compute convolutions rapidly, typically involving payoff approximations through piecewise polynomials or other simple functions. We do not cover these methods in

prices, only approximate ones based on, say, a small-time expansion (see, e.g., Section 13.1.9.1). In this case, a one-time convolution may be safe — especially if $T - T^*$ is small — whereas repeated convolutions may lead to unacceptable biases.

2.8.3 Analytics at the First Time Step

The idea in Section 2.8.2 of replacing the finite difference stepping with analytical integration is even easier to apply over the *first*, rather than the *last*, time step. Suppose T^* is the first “interesting” time for a given derivative security, i.e. there might be a jump condition at time T^* but none over the time interval $[0, T^*]$. Then, rather than stepping the finite difference scheme from T^* to 0, we can perform a *single* integration to calculate the value $V(0, x(0))$ of the derivative at time zero from the discretized values $\{V(T^*, x_j)\}_{j=0}^{m+1}$ of the derivative at time T^* (using the same notations as in Section 2.8.2),

$$V(0, x(0)) = \int_{\mathbb{R}} G(0, x(0); T^*, y) \tilde{V}(T^*, y) dy,$$

where $\tilde{V}(T^*, y)$ is interpolated (using cubic splines, say) from the values $\{V(T^*, x_j)\}_{j=0}^{m+1}$ on the grid. If the integral is computed numerically — as is most often the case — the numerical cost is often comparable with that of the finite difference stepping because only one value $V(0, x(0))$ is required at time 0, not the whole slice.

While there are typically no numerical cost savings that arise from using integration over the first time step, there are accuracy and stability considerations that favor this approach. We have already seen in Section 2.8.1 that the standard discretization of a PDE often leads to insufficient fidelity in resolving any features of the payoff that are close to today, and numerical integration can be of considerable help in this regard. Moreover, as we discuss in much detail later in Chapter 23, an integration scheme typically allows us to treat discontinuities in the value $V(T^*, x)$ arising from the jump condition at time T^* explicitly. If the discontinuity is introduced at the value of the state variable x^* , then the integration scheme can (and should) explicitly take this information into account. For example we would write

$$\begin{aligned} V(0, x(0)) &= \int_{-\infty}^{x^*} G(0, x(0); T^*, y) \tilde{V}^-(T^*, y) dy \\ &\quad + \int_{x^*}^{\infty} G(0, x(0); T^*, y) \tilde{V}^+(T^*, y) dy \end{aligned}$$

this book except for a brief mention in Section 11.A. For a representative example see Hu et al. [2006].

and calculate $\tilde{V}^-(T^*, y)$ by interpolating the grid values in the time interval $(-\infty, x^*)$, and $\tilde{V}^+(T^*, y)$ by interpolating the grid values in (x^*, ∞) , separately¹⁹.

The usefulness of the method is only limited by the availability of the closed-form expression for the time 0 Arrow-Debreu prices $G(0, x(0); T^*, \cdot)$. For some models this is not an issue; for most others, sufficiently close approximations could be obtained in a small-time limit (see e.g. Section 13.1.9.1 for a typical approach) that can be useful for times T^* that are not too large. By a change of measure, we see that

$$\begin{aligned} V(0, x(0)) &= \mathbb{E} \left(e^{-\int_0^{T^*} r(s) ds} V(T^*, x(T^*)) \right) \\ &= P(0, T^*) \mathbb{E}^{T^*} (V(T^*, x(T^*))), \end{aligned}$$

where \mathbb{E}^{T^*} is the expected value operator under the T^* -forward measure Q^{T^*} ; so we really only need the expression for the *density* (rather than Arrow-Debreu security prices) of $x(T^*)$ under Q^{T^*} , either exact or approximate.

Finally, we note that while the integration over the first time step can be seen to offer similar advantages to those of the methods in Section 2.8.1, the two approaches are not substitutes for each other, but are complementary. We typically recommend using direct integration over the time step $[0, T^*]$, where T^* is the smaller of the time of the first jump condition or the limit of applicability of the approximation to the density of $x(T^*)$, and then (if needed) use the methods in Section 2.8.1 over the time interval $[T^*, T]$, with T being the final maturity of the option in question.

2.9 Multi-Dimensional PDEs: Problem Formulation

We now turn our attention to the numerical solution of multi-dimensional terminal value problems. Let the spatial variable x be p -dimensional, $x = (x_1, \dots, x_p)^\top$, and consider the PDE

$$\frac{\partial V}{\partial t} + \sum_{h=1}^p \mu_h(t, x) \frac{\partial V}{\partial x_h} + \frac{1}{2} \sum_{h=1}^p \sum_{l=1}^p s_{h,l}(t, x) \frac{\partial^2 V}{\partial x_h \partial x_l} - r(t, x) = 0, \quad (2.56)$$

where $s_{h,h}(t, x) \geq 0$ and $s_{h,l}(t, x) = s_{l,h}(t, x)$ for $h, l = 1, \dots, p$. The PDE is assumed subject to the terminal value condition $V(T, x) = g(x)$, $g : \mathbb{R}^p \rightarrow \mathbb{R}$.

From the results in Chapter 1, we recognize that the PDE provides the solution to the expectation

¹⁹One of the functions $\tilde{V}^-(T^*, y)$, $\tilde{V}^+(T^*, y)$ is often known analytically and for all values of y (rather than sampled on the grid); this is for instance the case for the Bermudan options of Section 2.7.4. The integration algorithm should obviously take advantage of this.

$$V(t, x) = \mathbb{E}_t \left(e^{-\int_0^T r(u, x) du} g(x(T)) \mid x(t) = x \right),$$

where the components of $x(t)$ satisfy risk-neutral SDEs of the type

$$dx_h(t) = \mu_h(t, x(t)) dt + \sigma_h(t, x(t)) dW(t), \quad h = 1, \dots, p. \quad (2.57)$$

Here $W(t)$ is a d -dimensional Brownian motion, $\mu_h : [0, T] \times \mathbb{R}^p \rightarrow \mathbb{R}$, $h = 1, \dots, p$, are (scalar) drifts, and $\sigma_h : [0, T] \times \mathbb{R}^p \rightarrow \mathbb{R}^{1 \times d}$, $h = 1, \dots, p$, are d -dimensional (row vector) diffusion coefficients. The PDE coefficients $s_{h,l}$ in (2.56) represent the instantaneous covariance matrix for the components of $x(\cdot)$, i.e., $s_{h,l}(t, x) = \sigma_h(t, x)\sigma_l(t, x)^\top$. We assume enough regularity on μ_h , σ_h , r , and g to ensure that (2.56) has a unique solution.

For the purpose of solving (2.56) numerically, we assume that the PDE is to be solved on a (finite) spatial domain in x , $x \in [\underline{M}_1, \bar{M}_1] \times \dots \times [\underline{M}_p, \bar{M}_p]$, where $\underline{M}_h, \bar{M}_h$, $h = 1, \dots, p$, are constants either dictated by the contract at hand (barrier options) or found by a suitable probabilistic truncation (see Section 2.1).

2.10 Two-Dimensional PDE with No Mixed Derivatives

To illustrate the construction of finite difference discretization of (2.56), we start out with the simple case where $p = d = 2$ and there are no mixed partial derivatives in the PDE: $s_{1,2}(t, x) = s_{2,1}(t, x) = 0$ for all t and x . Probabilistically, the absence of mixed derivatives corresponds to the case where the stochastic process increments $dx_1(t)$ and $dx_2(t)$ are independent. Defining $\gamma_h(t, x)^2 = s_{h,h}(t, x)$, $h = 1, 2$, the PDE to be solved now becomes

$$\frac{\partial V}{\partial t} + (\mathcal{L}_1 + \mathcal{L}_2) V = 0, \quad (2.58)$$

where

$$\mathcal{L}_h \triangleq \mu_h(t, x) \frac{\partial}{\partial x_h} + \frac{1}{2} \gamma_h(t, x)^2 \frac{\partial^2}{\partial x_h^2} - \frac{1}{2} r(t, x), \quad h = 1, 2.$$

Notice that we have divided the term $r(t, x)$ into equal pieces in \mathcal{L}_1 and \mathcal{L}_2 .

To discretize (2.58) in x , introduce grids $x_1 \in \{x_1^{j_1}\}_{j_1=0}^{m_1+1}$ and $x_2 \in \{x_2^{j_2}\}_{j_2=0}^{m_2+1}$. To simplify notation, assume these grids are equidistant such that $x_1^{j_1} = \underline{M}_1 + j_1 \Delta_1$ and $x_2^{j_2} = \underline{M}_2 + j_2 \Delta_2$. Let $V_{j_1, j_2}(t) \triangleq V(t, x_1^{j_1}, x_2^{j_2})$. We define discrete central difference operators as before

$$\begin{aligned} \delta_{x_1} V_{j_1, j_2}(t) &= \frac{V_{j_1+1, j_2}(t) - V_{j_1-1, j_2}(t)}{2\Delta_1}, \\ \delta_{x_2} V_{j_1, j_2}(t) &= \frac{V_{j_1, j_2+1}(t) - V_{j_1, j_2-1}(t)}{2\Delta_2}, \end{aligned}$$

and

$$\delta_{x_1 x_1} V_{j_1, j_2}(t) = \frac{V_{j_1+1, j_2}(t) - 2V_{j_1, j_2}(t) + V_{j_1-1, j_2}(t)}{\Delta_1^2},$$

$$\delta_{x_2 x_2} V_{j_1, j_2}(t) = \frac{V_{j_1, j_2+1}(t) - 2V_{j_1, j_2}(t) + V_{j_1, j_2-1}(t)}{\Delta_2^2}.$$

These operators, in turn, give rise to the discrete operators

$$\widehat{\mathcal{L}}_h \triangleq \mu_h(t, x)\delta_{x_h} + \frac{1}{2}\gamma_h(t, x)^2\delta_{x_h x_h} - \frac{1}{2}r(t, x), \quad h = 1, 2,$$

where x is constrained to take values in the spatial grid. A Taylor expansion shows that this operator is second-order accurate (compare to Lemma 2.2.1),

$$(\mathcal{L}_1 + \mathcal{L}_2)V(t, x) = (\widehat{\mathcal{L}}_1 + \widehat{\mathcal{L}}_2)V(t, x) + O(\Delta_1^2 + \Delta_2^2).$$

2.10.1 Theta Method

Turning to a theta-style time discretization, consider first proceeding exactly as in Section 2.2.3. Assuming equidistant time spacing Δ_t , we get for the period $[t_i, t_{i+1}]$,

$$\begin{aligned} & \left(1 - \theta\Delta_t (\widehat{\mathcal{L}}_1 + \widehat{\mathcal{L}}_2)\right) V_{j_1, j_2}(t_i) \\ &= \left(1 + (1 - \theta)\Delta_t (\widehat{\mathcal{L}}_1 + \widehat{\mathcal{L}}_2)\right) V_{j_1, j_2}(t_{i+1}) + e_i^{i+1}, \end{aligned}$$

where

$$e_i^{i+1} = O\left(\Delta_t \left(\Delta_1^2 + \Delta_2^2 + 1_{\{\theta \neq \frac{1}{2}\}}\Delta_t + \Delta_t^2\right)\right),$$

and where it is understood that $\widehat{\mathcal{L}}_1$ and $\widehat{\mathcal{L}}_2$ are to be evaluated at $(t, x) = (t_i^{i+1}(\theta), x_1^{j_1}, x_2^{j_2})$ with $t_i^{i+1}(\theta)$ defined as in (2.14). If $\widehat{V}_{j_1, j_2}(t) \triangleq \widehat{V}(t, x_1^{j_1}, x_2^{j_2})$ is a finite difference approximation to $V_{j_1, j_2}(t)$, we thus get the scheme

$$\left(1 - \theta\Delta_t (\widehat{\mathcal{L}}_1 + \widehat{\mathcal{L}}_2)\right) \widehat{V}_{j_1, j_2}(t_i) = \left(1 + (1 - \theta)\Delta_t (\widehat{\mathcal{L}}_1 + \widehat{\mathcal{L}}_2)\right) \widehat{V}_{j_1, j_2}(t_{i+1}), \quad (2.59)$$

to be solved for the $m_1 m_2$ interior points $\widehat{V}_{j_1, j_2}(t_i)$, $j_1 = 1, \dots, m_1$, $j_2 = 1, \dots, m_2$, given the values of $\widehat{V}_{j_1, j_2}(t_{i+1})$, and given appropriate boundary conditions at $j_1 = 0$, $j_1 = m_1 + 1$, $j_2 = 0$, and $j_2 = m_2 + 1$.

The scheme (2.59) represents a system of linear equations in $m_1 m_2$ unknowns $\{\widehat{V}_{j_1, j_2}(t_i)\}$. When written out as a matrix equation (which requires us to arrange the various $\widehat{V}_{j_1, j_2}(t_i)$ in some order in a $(m_1 m_2)$ -dimensional vector), the matrix to be inverted is sparse but, unfortunately, no longer tri-diagonal. Solution of the system of equations by standard methods (e.g., Gauss-Jordan elimination or LU decomposition) is out of the question due

to the size of the matrix²⁰. We can proceed in two ways: either we use a specialized sparse-matrix solver; or we attempt to redo the discretization (2.59) to make it computationally efficient. We personally prefer the second approach and shall outline one method in the next section. As for the first approach, we simply note that a good iterative sparse solver should be able to solve (2.59) in order $O((m_1 m_2)^{5/4})$ operations. See Saad [2003] for concrete algorithms.

2.10.2 The Alternating Direction Implicit (ADI) Method

The ADI method is an example of a so-called *operator splitting* method, where the simultaneous application of two operators (here $\widehat{\mathcal{L}}_1$ and $\widehat{\mathcal{L}}_2$) is split into two *sequential* operator applications. To illustrate the idea, set $\theta = \frac{1}{2}$ (Crank-Nicolson scheme) in (2.59) and approximate

$$\left(1 - \frac{1}{2}\Delta_t (\widehat{\mathcal{L}}_1 + \widehat{\mathcal{L}}_2)\right) \approx \left(1 - \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_1\right) \left(1 - \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_2\right), \quad (2.60)$$

$$\left(1 + \frac{1}{2}\Delta_t (\widehat{\mathcal{L}}_1 + \widehat{\mathcal{L}}_2)\right) \approx \left(1 + \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_1\right) \left(1 + \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_2\right). \quad (2.61)$$

It is easy to see²¹ (and to verify, by a Taylor expansion) that the operators on the right-hand sides of these approximations have the same order truncation error as do the left-hand sides, namely $O(\Delta_t(\Delta_1^2 + \Delta_2^2 + \Delta_t^2))$. To the order of our original scheme, no accuracy is gained or lost in using the right-hand sides of (2.60)–(2.61). What is gained, however, is a considerable improvement in computational efficiency, originating in the fact that the resulting scheme

$$\begin{aligned} & \left(1 - \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_1\right) \left(1 - \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_2\right) \widehat{V}_{j_1,j_2}(t_i) \\ &= \left(1 + \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_1\right) \left(1 + \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_2\right) \widehat{V}_{j_1,j_2}(t_{i+1}) \end{aligned} \quad (2.62)$$

can be *split* into the system

$$\left(1 - \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_1\right) U_{j_1,j_2} = \left(1 + \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_2\right) \widehat{V}_{j_1,j_2}(t_{i+1}), \quad (2.63)$$

$$\left(1 - \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_2\right) \widehat{V}_{j_1,j_2}(t_i) = \left(1 + \frac{1}{2}\Delta_t \widehat{\mathcal{L}}_1\right) U_{j_1,j_2}, \quad (2.64)$$

²⁰Recall that the solution of a general linear system with $m_1 m_2$ unknowns is an $O(m_1^2 m_2^2)$ operation. For, say, m_1 and m_2 in the order of 100, this would involve around 1,000,000 times more work than what is required for a one-dimensional (tri-diagonal) scheme ($O(m)$).

²¹To those versed in operator notation, we notice that the right- and left-hand sides both approximate, to identical order, $\exp(\pm 0.5\Delta_t(\widehat{\mathcal{L}}_1 + \widehat{\mathcal{L}}_2))$.

where we have introduced an *intermediate value* U_{j_1,j_2} . The advantage of this decomposition is the fact that in each of (2.63) and (2.64), there is only one operator on the left-hand side, leading to simple tri-diagonal equation systems. To formalize this, first define

$$\mathbf{U}_1^{j_2} = (U_{1,j_2}, U_{2,j_2}, \dots, U_{m_1,j_2})^\top.$$

Then, for a fixed value of j_2 we can write for the first step

$$\left(\mathbf{I} - \frac{1}{2} \Delta_t \mathbf{A}_1^{j_2} \left(\frac{t_{i+1} + t_i}{2} \right) \right) \mathbf{U}_1^{j_2} = \mathbf{M}_2^{j_2} \left(\frac{t_{i+1} + t_i}{2} \right), \quad (2.65)$$

where $\mathbf{A}_1^{j_2}$ is an $(m_1 \times m_1)$ -dimensional tri-diagonal matrix of the same form as (2.11) (to get $\mathbf{A}_1^{j_2}$, basically freeze $x_2 = x_2^{j_2}$ and substitute μ_1 and γ_1 for μ and σ in the definition of the one-dimensional matrix \mathbf{A}). The m_1 -dimensional vector $\mathbf{M}_2^{j_2}$ has components $M_{2,j_1}^{j_2}$, $j_1 = 1, \dots, m_1$, given by

$$\begin{aligned} M_{2,j_1}^{j_2} \left(\frac{t_{i+1} + t_i}{2} \right) &= \left(1 + \frac{1}{2} \Delta_t \hat{\mathcal{L}}_2 \right) \hat{V}_{j_1,j_2}(t_{i+1}) \\ &= \frac{1}{2} \varsigma_{j_1,j_2}^- \hat{V}_{j_1,j_2-1}(t_{i+1}) + \frac{1}{2} \varsigma_{j_1,j_2}^+ \hat{V}_{j_1,j_2+1}(t_{i+1}) \\ &\quad + \left(1 - \frac{1}{2} \varsigma_{j_1,j_2} \right) \hat{V}_{j_1,j_2}(t_{i+1}), \end{aligned} \quad (2.66)$$

where we have defined

$$\begin{aligned} \varsigma_{j_1,j_2}^\pm &\triangleq \frac{\Delta_t}{2\Delta_2^2} \left(\gamma_2 \left(\frac{t_{i+1} + t_i}{2}, x_1^{j_1}, x_2^{j_2} \right)^2 \pm \Delta_2 \mu_2 \left(\frac{t_{i+1} + t_i}{2}, x_1^{j_1}, x_2^{j_2} \right) \right), \\ \varsigma_{j_1,j_2} &\triangleq \frac{\Delta_t}{\Delta_2^2} \left(\gamma_2 \left(\frac{t_{i+1} + t_i}{2}, x_1^{j_1}, x_2^{j_2} \right)^2 + \frac{1}{2} \Delta_2^2 r \left(\frac{t_{i+1} + t_i}{2}, x_1^{j_1}, x_2^{j_2} \right) \right). \end{aligned}$$

For known values of $\hat{V}(t_{i+1})$, (2.65) defines a simple tri-diagonal equation system which can be solved for $\mathbf{U}_1^{j_2}$ in $O(m_1)$ operations. Repeating the procedure above for $j_2 = 1, \dots, m_2$ allows us to find U_{j_1,j_2} for all $j_1 = 1, \dots, m_1$, $j_2 = 1, \dots, m_2$, at a total computational cost of $O(m_1 m_2)$.

Turning to the second step of (2.63)–(2.64), we first fix j_1 and define

$$\hat{\mathbf{V}}_2^{j_1}(t) = (\hat{V}_{j_1,1}(t), \hat{V}_{j_1,2}(t), \dots, \hat{V}_{j_1,m_2}(t))^\top.$$

In the same fashion as earlier, we can then write

$$\left(\mathbf{I} - \frac{1}{2} \Delta_t \mathbf{A}_2^{j_1} \left(\frac{t_{i+1} + t_i}{2} \right) \right) \hat{\mathbf{V}}_2^{j_1}(t_i) = \mathbf{M}_1^{j_1} \left(\frac{t_{i+1} + t_i}{2} \right), \quad (2.67)$$

where $\mathbf{A}_2^{j_1}$ is an $(m_2 \times m_2)$ -dimensional tri-diagonal matrix and where the right-hand side vector now has components

$$M_{1,j_2}^{j_1} \left(\frac{t_{i+1} + t_i}{2} \right) = \left(1 + \frac{1}{2} \Delta_t \widehat{\mathcal{L}}_1 \right) U_{j_1,j_2}, \quad j_2 = 1, \dots, m_2.$$

For brevity we omit writing out the $M_{1,j_2}^{j_1}$ (which will be similar to (2.66)), but just notice that the right-hand side of (2.67) is known after the first step of the ADI algorithm (above) is complete. For a given value of j_1 , we can solve the tri-diagonal system (2.67) for $\widehat{\mathbf{V}}_2^{j_1}(t_i)$ in $O(m_2)$ operations. Looping over all m_1 different values of j_1 , the full matrix of time t_i values $\widehat{\mathbf{V}}_{j_1,j_2}(t_i)$, $j_1 = 1, \dots, m_1$, $j_2 = 1, \dots, m_2$, can then be found at a total computational cost of $O(m_1 m_2)$.

The scheme outlined above is known as the *Peaceman-Rachford* scheme. As is the case for all ADI schemes, the scheme works by alternating the directions that are treated fully implicitly in the finite difference grid: in the first step, the x_1 -direction is fully implicit and the x_2 -direction is fully explicit, and in the second step the order is reversed. In effect, both spatial variables end up being discretized “semi-implicitly”, i.e. similar to a Crank-Nicolson scheme, resulting in convergence order is $O(\Delta_1^2 + \Delta_2^2 + \Delta_t^2)$. We emphasize, however, that whereas a direct application of the Crank-Nicolson scheme will involve (if an efficient sparse-matrix solver is used) a computational cost of $O((m_1 m_2)^{5/4})$ per time step, the computational cost of the Peaceman-Rachford ADI scheme is only $O(m_1 m_2)$. A (tedious) von Neumann analysis reveals that the scheme is A -stable, but, like the Crank-Nicolson scheme, not strongly A -stable.

While the Peaceman-Rachford scheme is a classical example of an ADI scheme, there are many others. For instance, consider a theta-version of the *Douglas-Rachford* scheme:

$$(1 - \theta \Delta_t \widehat{\mathcal{L}}_1) U_{j_1,j_2} = (1 + (1 - \theta) \Delta_t \widehat{\mathcal{L}}_1 + \Delta_t \widehat{\mathcal{L}}_2) \widehat{\mathbf{V}}_{j_1,j_2}(t_{i+1}), \quad (2.68)$$

$$(1 - \theta \Delta_t \widehat{\mathcal{L}}_2) \widehat{\mathbf{V}}_{j_1,j_2}(t_i) = U_{j_1,j_2} - \theta \Delta_t \widehat{\mathcal{L}}_2 \widehat{\mathbf{V}}_{j_1,j_2}(t_{i+1}), \quad (2.69)$$

where we understand that in $\widehat{\mathcal{L}}_1$ and $\widehat{\mathcal{L}}_2$ the PDE coefficients are to be evaluated at time $t_i^{i+1}(\theta)$. Again, notice how the scheme consists of two steps, each involving the solution of tri-diagonal sets of equations along only one of the x_1 - or x_2 -directions. The computational cost thus remains at $O(m_1 m_2)$. It can be shown that the convergence order of this scheme is $O(\Delta_1^2 + \Delta_2^2 + 1_{\{\theta \neq \frac{1}{2}\}} \Delta_t + \Delta_t^2)$ and it is A -stable for $\theta \geq \frac{1}{2}$, and strongly A -stable for $\theta > \frac{1}{2}$. By elimination of U_{j_1,j_2} we note that the unsplit version of the Douglas-Rachford scheme is

$$\begin{aligned} & (1 - \theta \Delta_t \widehat{\mathcal{L}}_1) (1 - \theta \Delta_t \widehat{\mathcal{L}}_2) \widehat{\mathbf{V}}_{j_1,j_2}(t_i) \\ &= ((1 - \theta \Delta_t \widehat{\mathcal{L}}_1) (1 - \theta \Delta_t \widehat{\mathcal{L}}_2) + \Delta_t \widehat{\mathcal{L}}_1 + \Delta_t \widehat{\mathcal{L}}_2) \widehat{\mathbf{V}}_{j_1,j_2}(t_{i+1}). \end{aligned}$$

It is not difficult to see that this approximates (2.59) to second order.

2.10.3 Boundary Conditions and Other Issues

The fact that ADI schemes reduce to solving sequences of matrix systems identical to the ones arising in the one-dimensional case is convenient, in the sense that many of the issues we have encountered for one-dimensional finite difference grids (oscillations, stability, convection dominance, etc.) and their remedies (smoothing, non-equidistant discretization, upwinding, etc.) carry over to the ADI setting with only minor modifications. Consider for instance the issue of applying spatial boundary conditions along the edges of the (x_1, x_2) domain, which we have so far not discussed. As for the one-dimensional PDEs, the most convenient way to express such boundary conditions is typically by imposing conditions on derivatives, like $\partial^2 V(t, x_1^0, x_2^{j_2}) / \partial x_1^2 = \partial V(t, x_1^0, x_2^{j_2}) / \partial x_1$ and so forth. For the Peaceman-Rachford scheme, say, such conditions can be incorporated directly into (2.65) and (2.67) by altering the matrices $\mathbf{A}_1^{j_2}$ and $\mathbf{A}_2^{j_1}$, as well as the boundary elements of $\mathbf{M}_1^{j_1}$ and $\mathbf{M}_2^{j_2}$, in the manner outlined in Section 2.2.1. If instead we wish to impose Dirichlet boundary conditions, we need to add corrective terms to the tri-diagonal systems, as in (2.19). To complete the first part of the split scheme, this then requires us to establish what boundary terms are needed for the intermediate quantity U_{j_1, j_2} , i.e. we must define $U_{j_1, 0}$ and U_{j_1, m_2+1} for $j_1 = 1, \dots, m_1$, as well as U_{0, j_2} and U_{m_1+1, j_2} for $j_2 = 1, \dots, m_2$. While U_{j_1, j_2} is a purely mathematic construct, sometimes it is adequate to think of U_{j_1, j_2} as a proxy for V_{j_1, j_2} evaluated at $t_i^{i+1}(\theta)$, which obviously makes determination of boundary conditions straightforward. For maximum precision, however, we should use the ADI equations themselves to express the boundary conditions of U directly in terms of boundary conditions for $V(t_i)$ and $V(t_{i+1})$. Here, the Douglas-Rachford scheme is particularly easy to deal with, as a rearrangement of (2.69) directly relates U_{j_1, j_2} to $\widehat{V}_{j_1, j_2}(t_i)$ and $\widehat{V}_{j_1, j_2}(t_{i+1})$,

$$U_{j_1, j_2} = \left(1 - \theta \Delta_t \widehat{\mathcal{L}}_2\right) \widehat{V}_{j_1, j_2}(t_i) + \theta \Delta_t \widehat{\mathcal{L}}_2 \widehat{V}_{j_1, j_2}(t_{i+1}).$$

The Peaceman-Rachford scheme requires some further manipulations to express U in terms of $V(t_i)$ and $V(t_{i+1})$; see Mitchell and Griffiths [1980] for the details.

2.11 Two-Dimensional PDE with Mixed Derivatives

Consider now the case where the 2-dimensional PDE (2.58) has a mixed partial derivative,

$$\frac{\partial V}{\partial t} + (\mathcal{L}_1 + \mathcal{L}_2 + \mathcal{L}_{1,2}) V = 0, \quad (2.70)$$

where \mathcal{L}_1 and \mathcal{L}_2 are as in (2.58), and where

$$\mathcal{L}_{1,2} = s_{1,2}(t, x) \frac{\partial^2}{\partial x_1 \partial x_2} \triangleq \rho(t, x) \gamma_1(t, x) \gamma_2(t, x) \frac{\partial^2}{\partial x_1 \partial x_2}. \quad (2.71)$$

The quantity $\rho(t, x)$ is the instantaneous correlation between the processes $x_1(t)$ and $x_2(t)$ in (2.57), i.e. $\rho(t, x) \in [-1, 1]$.

The presence of $\mathcal{L}_{1,2}$ prevents a direct application of the ADI methods in Section 2.10.2, since the mixed operator $\mathcal{L}_{1,2}$ is not amenable to operator splitting. We shall demonstrate two ways to overcome this problem: a) orthogonalization of the PDE; and b) predictor-corrector schemes.

2.11.1 Orthogonalization of the PDE

The idea here is to introduce new variables $y_1(t, x_1, x_2)$ and $y_2(t, x_1, x_2)$ such that the PDE loses its mixed derivative term when stated in terms of these variables. To demonstrate this idea, assume first that $\rho(t, x)$, $\gamma_1(t, x)$, and $\gamma_2(t, x)$ are all functions of time only and independent of x . Then define, say,

$$y_1(t, x_1, x_2) = x_1, \quad (2.72)$$

$$y_2(t, x_1, x_2) = -\rho(t) \frac{\gamma_2(t)}{\gamma_1(t)} x_1 + x_2 \triangleq a(t)x_1 + x_2, \quad (2.73)$$

where we must assume that $\gamma_1(t) \neq 0$ for all t .

Lemma 2.11.1. *Consider the PDE (2.70) subject to the terminal value condition $V(T, x) = g(x)$. Define $y = (y_1, y_2)^\top$ and $v(t, y) = V(t, x)$. With the variable change defined in (2.72)–(2.73), v satisfies*

$$\begin{aligned} \frac{\partial v}{\partial t} + \mu_1^y(t, y) \frac{\partial v}{\partial y_1} + \mu_2^y(t, y) \frac{\partial v}{\partial y_2} + \frac{1}{2} \gamma_1(t)^2 \frac{\partial^2 v}{\partial y_1^2} \\ + \frac{1}{2} (1 - \rho(t)^2) \gamma_2(t)^2 \frac{\partial^2 v}{\partial y_2^2} - r(t, y_1, y_2 - a(t)y_1) = 0, \end{aligned} \quad (2.74)$$

where

$$\mu_1^y(t, y) \triangleq \mu_1(t, x_1, x_2) = \mu_1(t, y_1, y_2 - a(t)y_1), \quad (2.75)$$

$$\begin{aligned} \mu_2^y(t, y) &\triangleq \frac{da(t)}{dt} x_1 + a(t) \mu_1(t, x_1, x_2) + \mu_2(t, x_1, x_2) \\ &= \frac{da(t)}{dt} y_1 + a(t) \mu_1^y(t, y) + \mu_2(t, y_1, y_2 - a(t)y_1). \end{aligned} \quad (2.76)$$

The equation (2.74) is subject to the terminal value condition $v(T, y_1, y_2) = g(x_1, x_2) = g(y_1, y_2 - a(T)y_1)$.

Proof. While the result can be established by the usual mechanics of ordinary calculus, we will take the opportunity to show how stochastic calculus can

also conveniently prove results of this type. Going back to the processes underlying the PDE (see (2.57)), we write

$$dx_1(t) = \mu_1(t, x) dt + \gamma_1(t) dW_1(t), \quad (2.77)$$

$$dx_2(t) = \mu_2(t, x) dt + \gamma_2(t) \left(\rho(t) dW_1(t) + \sqrt{1 - \rho(t)^2} dW_2(t) \right), \quad (2.78)$$

for independent scalar Brownian motions $W_1(t)$ and $W_2(t)$; this is easily seen to generate the correct correlation $\rho(t)$ between x_1 and x_2 . An application of Ito's lemma then shows that the processes for y_1 and y_2 are

$$\begin{aligned} dy_1(t) &= dx_1(t) = \mu_1(t, x) dt + \gamma_1(t) dW_1(t), \\ dy_2(t) &= \frac{da(t)}{dt} x_1(t) dt + a(t) \mu_1(t, x) dt + a(t) \gamma_1(t) dW_1(t) \\ &\quad + \mu_2(t, x) dt + \gamma_2(t) \left(\rho(t) dW_1(t) + \sqrt{1 - \rho(t)^2} dW_2(t) \right) \\ &= \left(\frac{da(t)}{dt} x_1(t) + a(t) \mu_1(t, x) + \mu_2(t, x) \right) dt \\ &\quad + \gamma_2(t) \sqrt{1 - \rho(t)^2} dW_2(t). \end{aligned}$$

With the definitions (2.75)–(2.76), this becomes simply

$$dy_1(t) = \mu_1^y(t, y(t)) dt + \gamma_1(t) dW_1(t), \quad (2.79)$$

$$dy_2(t) = \mu_2^y(t, y(t)) dt + \gamma_2(t) \sqrt{1 - \rho(t)^2} dW_2(t). \quad (2.80)$$

Equations (2.79)–(2.80) define a Markov SDE in $y_1(t)$ and $y_2(t)$ where, importantly, the Brownian motions on $y_1(t)$ and $y_2(t)$ are now *independent*. Writing $V(t, x) = v(t, y)$, it then follows immediately from the backward Kolmogorov equation (see Section 1.8) that v satisfies the PDE (2.74). \square

Through the chosen transformation (2.72)–(2.73), our original PDE has now been put into a form where we can immediately apply the ADI schemes outlined in Section 2.10.2.

In performing the orthogonalization of the PDE in Lemma 2.11.1 we relied on $\rho(t, x)$, $\gamma_1(t, x)$, and $\gamma_2(t, x)$ all being independent of x . This can often be relaxed. Consider for instance the case where $\rho(t, x) = \rho(t)$, $\gamma_1(t, x) = \gamma_1(t, x_1)$, and $\gamma_2(t, x) = \gamma_2(t, x_2)$; here the correlation ρ is still assumed deterministic, but we now allow for some (though not full) x -dependence in γ_1 and γ_2 . Assuming that $\gamma_1(t, x_1) > 0$ and $\gamma_2(t, x_2) > 0$ we can introduce new variables

$$z_1(t, x_1) = \int \frac{1}{\gamma_1(t, x_1)} dx_1, \quad (2.81)$$

$$z_2(t, x_2) = \int \frac{1}{\gamma_2(t, x_2)} dx_2. \quad (2.82)$$

Applying Ito's lemma to (2.77)–(2.78) we see that

$$dz_1(t, x_1) = \left(- \int \frac{\partial \gamma_1(t, x_1)}{\partial t} \frac{1}{\gamma_1(t, x_1)^2} dx_1 + \frac{\mu_1(t, x)}{\gamma_1(t, x_1)} - \frac{1}{2} \frac{\partial \gamma_1(t, x_1)}{\partial x_1} \right) dt + dW_1(t) \quad (2.83)$$

and

$$dz_2(t, x_2) = \left(- \int \frac{\partial \gamma_2(t, x_2)}{\partial t} \frac{1}{\gamma_2(t, x_2)^2} dx_2 + \frac{\mu_2(t, x)}{\gamma_2(t, x_2)} - \frac{1}{2} \frac{\partial \gamma_2(t, x_2)}{\partial x_2} \right) dt + \rho(t) dW_1(t) + \sqrt{1 - \rho(t)^2} dW_2(t). \quad (2.84)$$

As we assumed that $\gamma_1(t, x_1) > 0$ and $\gamma_2(t, x_2) > 0$, the functions z_1 and z_2 are increasing in x_1 and x_2 , respectively, and are thereby invertible. As such, we can rewrite (2.83)–(2.84) in the more appealing form

$$\begin{aligned} dz_1(t, x_1) &= \mu_1^z(t, z_1, z_2) dt + dW_1(t), \\ dz_2(t, x_1) &= \mu_2^z(t, z_1, z_2) dt + \rho(t) dW_1(t) + \sqrt{1 - \rho(t)^2} dW_2(t). \end{aligned}$$

Through the transformation (2.81)–(2.82), we have reduced our original system to one where the coefficients on $W_1(t)$ and $W_2(t)$ are no longer state-dependent, similar to the case that lead to Lemma 2.11.1. We can now proceed with another variable transformation, as in (2.72)–(2.73), to orthogonalize the system and prepare it for an application of the ADI method.

While the orthogonalization method outlined here can be very effective on a range of practical problems, it suffers from a few drawbacks. Most obviously, the method is not completely general and requires a certain structure on the parameters of the PDE. Another drawback is that the introduction of a time-dependent transformation on one or more variables (Lemma 2.11.1) often makes the alignment of the finite difference grid along (time-independent) critical level points in x -space impossible. Also, the introduction of terms like $y_1 da(t)/dt$ in the drift of y_2 (see (2.76)) can be problematic, particularly if the functions $\gamma_1(t)$ and $\gamma_2(t)$ are not smooth. For instance, it is not unlikely that $y_1 da(t)/dt$ will locally be of such magnitude that upwinding will be necessary to prevent oscillations; see Section 2.6.1. Further, we note that inversion of the transformations (2.81)–(2.82) will not always be possible to perform analytically and may require numerical (root-search) work, complicating the scheme and potentially slowing it down. Finally, as we shall highlight in future chapters, maintaining the “continuity” of a numerical scheme with respect to input parameters is of critical importance for the smoothness of risk sensitivities. Such continuity is difficult to ensure if complicated transformations are applied to model variables. So, in the end, we recommend formulating the PDEs in terms of financially meaningful variables, avoiding excessive transformations, and relying on methods such as developed in the next section when dealing with mixed derivatives and other numerical complications.

2.11.2 Predictor-Corrector Scheme

In this section we shall consider a completely general method for handling mixed derivatives in two-dimensional PDEs. While a bit slower than the method outlined in Section 2.11.1, it does not involve any variable transformations and, by extension, does not suffer from the drawbacks associated with such transformations. As a first step, consider the discretization of the mixed derivative $\partial^2 V / \partial x_1 \partial x_2$. There are a few possibilities (see Mitchell and Griffiths [1980]), but we shall just use

$$\begin{aligned} \delta_{x_1 x_2} V_{j_1, j_2}(t) &= \delta_{x_1} \delta_{x_2} V_{j_1, j_2}(t) \\ &= \frac{V_{j_1+1, j_2+1}(t) - V_{j_1+1, j_2-1}(t) - V_{j_1-1, j_2+1}(t) + V_{j_1-1, j_2-1}(t)}{4\Delta_1 \Delta_2}. \end{aligned} \quad (2.85)$$

Extensions to non-equidistant grids follow directly from (2.27) and the relation $\delta_{x_1 x_2} V_{j_1, j_2}(t) = \delta_{x_1} \delta_{x_2} V_{j_1, j_2}(t)$. As we have not encountered mixed difference operators before, for completeness we show the following lemma.

Lemma 2.11.2. *For the discrete operator (2.85) we have*

$$\delta_{x_1 x_2} V_{j_1, j_2}(t) = \frac{\partial^2 V(t, x_1^{j_1}, x_2^{j_2})}{\partial x_1 \partial x_2} + O(\Delta_1^2 + \Delta_2^2).$$

Proof. A Taylor expansion of $V(t, x)$ around the point $x = (x_1^{j_1}, x_2^{j_2})^\top$ gives

$$\begin{aligned} V_{j_1+1, j_2 \pm 1}(t) &= V_{j_1, j_2}(t) + \Delta_1 \frac{\partial V}{\partial x_1} \pm \Delta_2 \frac{\partial V}{\partial x_2} + \frac{1}{2} \Delta_1^2 \frac{\partial^2 V}{\partial x_1^2} + \frac{1}{2} \Delta_2^2 \frac{\partial^2 V}{\partial x_2^2} \\ &\quad \pm \Delta_1 \Delta_2 \frac{\partial^2 V}{\partial x_1 \partial x_2} + \frac{1}{6} \Delta_1^3 \frac{\partial^3 V}{\partial x_1^3} \pm \frac{1}{6} \Delta_2^3 \frac{\partial^3 V}{\partial x_2^3} \\ &\quad + \frac{1}{2} \Delta_1 \Delta_2^2 \frac{\partial^3 V}{\partial x_1 \partial x_2^2} \pm \frac{1}{2} \Delta_1^2 \Delta_2 \frac{\partial^3 V}{\partial x_1^2 \partial x_2} + \dots, \\ V_{j_1-1, j_2 \pm 1}(t) &= V_{j_1, j_2}(t) - \Delta_1 \frac{\partial V}{\partial x_1} \pm \Delta_2 \frac{\partial V}{\partial x_2} + \frac{1}{2} \Delta_1^2 \frac{\partial^2 V}{\partial x_1^2} + \frac{1}{2} \Delta_2^2 \frac{\partial^2 V}{\partial x_2^2} \\ &\quad \mp \Delta_1 \Delta_2 \frac{\partial^2 V}{\partial x_1 \partial x_2} - \frac{1}{6} \Delta_1^3 \frac{\partial^3 V}{\partial x_1^3} \pm \frac{1}{6} \Delta_2^3 \frac{\partial^3 V}{\partial x_2^3} \\ &\quad - \frac{1}{2} \Delta_1 \Delta_2^2 \frac{\partial^3 V}{\partial x_1 \partial x_2^2} \pm \frac{1}{2} \Delta_1^2 \Delta_2 \frac{\partial^3 V}{\partial x_1^2 \partial x_2} + \dots. \end{aligned}$$

A little thought then shows that

$$\begin{aligned} V_{j_1+1, j_2+1}(t) - V_{j_1+1, j_2-1}(t) - V_{j_1-1, j_2+1}(t) + V_{j_1-1, j_2-1}(t) \\ = 4\Delta_1 \Delta_2 \frac{\partial^2 V}{\partial x_1 \partial x_2} + O(\Delta_1^3 \Delta_2 + \Delta_1 \Delta_2^3), \end{aligned}$$

as error terms of order Δ_1^4 , Δ_2^4 , and $\Delta_1^2 \Delta_2^2$ will cancel. The result follows.

□

Equipped with (2.85), we can approximate the operator $\mathcal{L}_{1,2}$ in (2.71) as

$$\widehat{\mathcal{L}}_{1,2} V_{j_1,j_2}(t) \triangleq \rho(t, x_1^{j_1}, x_2^{j_2}) \gamma_1(t, x_1^{j_1}, x_2^{j_2}) \gamma_2(t, x_1^{j_1}, x_2^{j_2}) \delta_{x_1 x_2} V_{j_1,j_2}(t),$$

which is accurate to order $O(\Delta_1^2 + \Delta_2^2)$. The first easy way to modify our ADI scheme to incorporate $\widehat{\mathcal{L}}_{1,2}$ is to treat the mixed derivative fully explicitly. In the Douglas-Rachford scheme (2.68)–(2.69), for instance, we thus modify the right-hand side of the first step as follows:

$$(1 - \theta \Delta_t \widehat{\mathcal{L}}_1) U_{j_1,j_2} = (1 + (1 - \theta) \Delta_t \widehat{\mathcal{L}}_1 + \Delta_t \widehat{\mathcal{L}}_2 + \Delta_t \widehat{\mathcal{L}}_{1,2}) \widehat{V}_{j_1,j_2}(t_{i+1}), \quad (2.86)$$

$$(1 - \theta \Delta_t \widehat{\mathcal{L}}_2) \widehat{V}_{j_1,j_2}(t_i) = U_{j_1,j_2} - \theta \Delta_t \widehat{\mathcal{L}}_2 \widehat{V}_{j_1,j_2}(t_{i+1}). \quad (2.87)$$

The addition of $\widehat{\mathcal{L}}_{1,2}$ this way clearly preserves the ADI structure of the scheme which will continue to involve only sequences of tri-diagonal linear equations. However, having, in effect, only a one-sided time-differencing of the mixed derivative term will lower the convergence order of the time step to $O(\Delta_t)$, irrespective of the choice of θ .

To change the time at which the mixed operator $\widehat{\mathcal{L}}_{1,2}$ is evaluated, consider using a *predictor-corrector* scheme, where the results of (2.86)–(2.87) are re-used in a one-time²² iteration. Specifically, we write, for some $\lambda \in [0, 1]$,

Predictor:

$$(1 - \theta \Delta_t \widehat{\mathcal{L}}_1) U_{j_1,j_2}^{(1)} = (1 + (1 - \theta) \Delta_t \widehat{\mathcal{L}}_1 + \Delta_t \widehat{\mathcal{L}}_2 + \Delta_t \widehat{\mathcal{L}}_{1,2}) \widehat{V}_{j_1,j_2}(t_{i+1}), \quad (2.88)$$

$$(1 - \theta \Delta_t \widehat{\mathcal{L}}_2) U_{j_1,j_2}^{(2)} = U_{j_1,j_2}^{(1)} - \theta \Delta_t \widehat{\mathcal{L}}_2 \widehat{V}_{j_1,j_2}(t_{i+1}). \quad (2.89)$$

Corrector:

$$\begin{aligned} (1 - \theta \Delta_t \widehat{\mathcal{L}}_1) Z_{j_1,j_2}^{(1)} &= (1 + (1 - \theta) \Delta_t \widehat{\mathcal{L}}_1 + \Delta_t \widehat{\mathcal{L}}_2 \\ &\quad + (1 - \lambda) \Delta_t \widehat{\mathcal{L}}_{1,2}) \widehat{V}_{j_1,j_2}(t_{i+1}) + \lambda \Delta_t \widehat{\mathcal{L}}_{1,2} U_{j_1,j_2}^{(2)}, \end{aligned} \quad (2.90)$$

$$(1 - \theta \Delta_t \widehat{\mathcal{L}}_2) \widehat{V}_{j_1,j_2}(t_i) = Z_{j_1,j_2}^{(1)} - \theta \Delta_t \widehat{\mathcal{L}}_2 \widehat{V}_{j_1,j_2}(t_{i+1}). \quad (2.91)$$

²²We can run the iteration more than once if desired, but a single iteration will normally suffice.

Notice how the Douglas-Rachford scheme is first run once, in (2.88)–(2.89), to yield a first guess (a “predictor”), $U_{j_1,j_2}^{(2)}$, for the time t_i value $V_{j_1,j_2}(t_i)$. In a second run of the Douglas-Rachford scheme, in (2.90)–(2.91), this guess is used as a “corrector” to affect the time at which $\widehat{\mathcal{L}}_{1,2}$ is evaluated, by applying this operator to $(1 - \lambda)\widehat{V}_{j_1,j_2}(t_{i+1}) + \lambda U_{j_1,j_2}^{(2)}$; when $\lambda = \frac{1}{2}$ we effectively center the time-differencing of the mixed term. The scheme now relies on three intermediate variables, $U_{j_1,j_2}^{(1)}$, $U_{j_1,j_2}^{(2)}$, and $Z_{j_1,j_2}^{(1)}$.

The combined predictor-corrector scheme above (in a slightly less general form, with $\Delta_1 = \Delta_2$) was suggested by Craig and Sneyd [1988]. It can be shown that the scheme has convergence order

$$O\left((\Delta_1 + \Delta_2)^2 + 1_{\{\theta \neq \frac{1}{2}\}}\Delta_t + 1_{\{\lambda \neq \frac{1}{2}\}}\Delta_t + \Delta_t^2\right),$$

so second order convergence in the time domain is still achievable by setting $\theta = \lambda = \frac{1}{2}$. The scheme will be A -stable for $\theta \geq \frac{1}{2}$ and $\frac{1}{2} \leq \lambda \leq \theta$. The computational cost of the predictor-corrector is clearly still $O(m_1 m_2)$ per time step, as both the predictor and corrector schemes have $O(m_1 m_2)$ cost per time-step. Even though the standard Douglas-Rachford scheme is effectively run twice, we should point out that when intelligently implemented, (2.88)–(2.91) is typically only about 30-40% slower than the Douglas-Rachford scheme, as a number of results from the predictor step can be cached and reused in the corrector step.

As for the standard ADI grids, extensions to non-equidistant grids are straightforward using the techniques in Section 2.4. Boundary conditions in the x -domain are imposed along the lines outlined in Section 2.10.3.

2.12 PDEs of Arbitrary Order

We now turn our attention back to the general p -dimensional PDE (2.56). To prepare for a numerical scheme, let us rewrite the PDE as follows:

$$\frac{\partial V}{\partial t} + \sum_{h=1}^p \mathcal{L}_h V + \sum_{h=1}^p \sum_{l=h+1}^p \mathcal{L}_{h,l} V = 0, \quad (2.92)$$

where

$$\begin{aligned} \mathcal{L}_h &= \mu_h(t, x) \frac{\partial}{\partial x_h} + \frac{1}{2} s_{h,h}(t, x) \frac{\partial^2}{\partial x_h^2} - p^{-1} r(t, x), \\ \mathcal{L}_{h,l} &= s_{h,l}(t, x) \frac{\partial^2}{\partial x_h \partial x_l}. \end{aligned}$$

The method we present here for solution of (2.92) is a p -dimensional version of the predictor-corrector scheme outlined above. The extension

is straightforward and we simply list it here without further discussion; see Craig and Sneyd [1988] for additional background. To simplify notation, we have omitted sub-indices everywhere (i.e., $\widehat{V}(t_i)$ is used instead of $\widehat{V}_{j_1, j_2, \dots, j_p}(t_i)$).

Predictor:

$$\begin{aligned} & \left(1 - \theta \Delta_t \widehat{\mathcal{L}}_1\right) U^{(1)} \\ &= \Delta_t \left(\Delta_t^{-1} + (1 - \theta) \widehat{\mathcal{L}}_1 + \sum_{h=2}^p \widehat{\mathcal{L}}_h + \sum_{h=1}^p \sum_{l=h+1}^p \widehat{\mathcal{L}}_{h,l} \right) \widehat{V}(t_{i+1}), \\ & \left(1 - \theta \Delta_t \widehat{\mathcal{L}}_2\right) U^{(2)} = U^{(1)} - \theta \Delta_t \widehat{\mathcal{L}}_2 \widehat{V}(t_{i+1}), \\ & \vdots \\ & \left(1 - \theta \Delta_t \widehat{\mathcal{L}}_p\right) U^{(p)} = U^{(p-1)} - \theta \Delta_t \widehat{\mathcal{L}}_p \widehat{V}(t_{i+1}). \end{aligned}$$

Corrector:

$$\begin{aligned} & \left(1 - \theta \Delta_t \widehat{\mathcal{L}}_1\right) Z^{(1)} \\ &= \Delta_t \left(\Delta_t^{-1} + (1 - \theta) \widehat{\mathcal{L}}_1 + \sum_{h=2}^p \widehat{\mathcal{L}}_h \right. \\ & \quad \left. + (1 - \lambda) \sum_{h=1}^p \sum_{l=h+1}^p \widehat{\mathcal{L}}_{h,l} \right) \widehat{V}(t_{i+1}) + \lambda \Delta_t \sum_{h=1}^p \sum_{l=h+1}^p \widehat{\mathcal{L}}_{h,l} U^{(p)}, \\ & \left(1 - \theta \Delta_t \widehat{\mathcal{L}}_2\right) Z^{(2)} = Z^{(1)} - \theta \Delta_t \widehat{\mathcal{L}}_2 \widehat{V}(t_{i+1}), \\ & \vdots \\ & \left(1 - \theta \Delta_t \widehat{\mathcal{L}}_p\right) \widehat{V}(t_i) = Z^{(p-1)} - \theta \Delta_t \widehat{\mathcal{L}}_p \widehat{V}(t_{i+1}). \end{aligned}$$

With m_h points in the x_h -direction, $h = 1, \dots, p$, the computational cost of the predictor-corrector scheme is $O(\prod_{h=1}^p m_h)$. For $p \leq 3$, sufficient conditions for A -stability are $\theta \geq \frac{1}{2}$ and $\frac{1}{2} \leq \lambda \leq \theta$. For $p \geq 4$, sufficient conditions are $\theta \leq \frac{1}{2}$ and

$$\frac{1}{2} \leq \lambda \leq \frac{p^{p-1}}{(p-1)^p} \theta.$$

See Craig and Sneyd [1988] for a proof. Convergence is similar to the two-dimensional case.

As a final comment, let us note that as dimensionality increases, the computational complexity of an iterative sparse solver will start approaching that of ADI. Specifically, for a p -dimensional problem, the complexity of the former is $O(m_{\text{total}})$ and for the latter $O(m_{\text{total}}^{(2p+1)/2p})$, with $m_{\text{total}} = m_1 \cdot m_2 \cdot \dots \cdot m_p$.

Monte Carlo Methods

While the finite difference method is flexible and powerful, it has a number of limitations. First, its usage is restricted to problems where the state variable dynamics are Markovian. Second, for strongly path-dependent problems, the method often does not apply. And third, it is unsuited for problems where the dimension of the underlying vector of state variables is high. To expand on the last point, recall from Section 2.9 that the (ADI) finite difference method applied to a p -dimensional problem has computational complexity $O(m^p)$ per time step, where m is the average number of spatial points per dimension. The exponential growth in p — the “curse of dimensionality” — is typical of grid-based methods and prevents the practical usage of the method for p larger than about 4 or 5.

In this chapter, we study the *Monte Carlo method*, a numerical technique where the computational effort grows only linearly in the problem dimension p . While convergence of the Monte Carlo method is relatively slow, it is nearly always the method of choice for high-dimensional pricing problems. Compared to finite difference methods, Monte Carlo methods are easy to apply to problems with non-Markovian dynamics as well as strong path-dependency in the payout. On the other hand, as Monte Carlo methods inherently run forward in time, dynamic programming techniques are challenging to implement, making Monte Carlo pricing of American and Bermudan options significantly more involved than for the naturally backward-working finite difference method.

3.1 Fundamentals

Consider a European-style derivative V with time T payout $V(T) = g(T)$, where $g(T)$ is an \mathcal{F}_T -measurable (and integrable) random variable. Where finite difference methods start with a PDE representation of the price of a contingent claim at times $t < T$, the starting point for the Monte Carlo method is the basic martingale relation (see (1.15))

$$V(t) = N(t) \mathbb{E}_t^{Q^N} (g(T)/N(T)), \quad (3.1)$$

where $N(\cdot)$ is a numeraire and Q^N is the measure induced by $N(\cdot)$. To evaluate this expression numerically, we need a numerical technique to compute expectations of a random variable. For this, we turn to the law of large numbers:

Theorem 3.1.1 (Strong Law of Large Numbers). *Let Y_1, Y_2, \dots be a sequence of independent identically distributed (i.i.d.) random variables with expectation $\mu < \infty$. Define the sample mean*

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i. \quad (3.2)$$

Then

$$\lim_{n \rightarrow \infty} \bar{Y}_n = \mu, \quad a.s.$$

This result forms the basis for the *Monte Carlo method*, which computes the expectation in (3.1) by simply i) generating independent realizations of $g(T)/N(T)$ under Q^N ; and ii) forming their average. Specifically, let $g_1/N_1, \dots, g_n/N_n$ denote n independent samples from the distribution of $g(T)/N(T)$, conditional on \mathcal{F}_t . Then our Monte Carlo estimator for $V(t)$ is the sample mean

$$\bar{V}(t) = N(t) \frac{1}{n} \sum_{i=1}^n g_i/N_i. \quad (3.3)$$

We shall delve into how to generate samples from the distribution of $g(T)/N(T)$ shortly, but before doing so let us consider the expected convergence rate of the Monte Carlo method as n is increased. The key result is here the central limit theorem:

Theorem 3.1.2 (Central Limit Theorem). *Let Y_1, Y_2, \dots be a sequence of i.i.d. random variables with expectation μ and standard deviation $\sigma < \infty$. Let the sample mean be defined as in (3.2). Then, for $n \rightarrow \infty$,*

$$\frac{\bar{Y}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1),$$

where $\mathcal{N}(0, 1)$ is a standard Gaussian distribution and \xrightarrow{d} denotes convergence in distribution¹. Further, if we define

$$s_n = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2},$$

¹Recall that a sequence of variables X_n with cumulative distribution functions F_n converge in distribution to a random variable X with distribution F if $\lim_{n \rightarrow \infty} F_n(x) = F(x)$ for all $x \in \mathbb{R}$ at which $F(x)$ is continuous.

then also

$$\frac{\bar{Y}_n - \mu}{s_n/\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1).$$

Define the Gaussian percentile u_γ as $\Phi(u_\gamma) = 1 - \gamma$, where Φ is the Gaussian cumulative distribution function. From Theorem 3.1.2, and from the definition of convergence in distribution (see footnote 1), the probability that the confidence interval

$$[\bar{V}(t) - u_{\gamma/2} \cdot s_n/\sqrt{n}, \bar{V}(t) + u_{\gamma/2} \cdot s_n/\sqrt{n}] \quad (3.4)$$

fails to include the true value $V(t)$ approaches γ for large n . Here

$$s_n \triangleq \sqrt{\frac{1}{n-1} \sum_{i=1}^n \left(\frac{g_i N(t)}{N_i} - \bar{V}(t) \right)^2},$$

with the quantity s_n/\sqrt{n} known as the *standard error*. For given γ , the rate at which the confidence interval for $V(t)$ contracts is $O(n^{-\frac{1}{2}})$. This is relatively slow: to reduce the width of the interval by a factor of 2, n must increase by a factor of 4. On the other hand, we notice that the (asymptotic) convergence rate only depends on n , not on the specifics of the g_i 's. In particular, if $g(T) = g(X(T))$ where X is p -dimensional, the asymptotic convergence rate is independent of p . As we shall see shortly, in most applications the work required to generate samples of $g(X(T))$ is (at most) linear in p .

3.1.1 Generation of Random Samples

At the most basic level, the Monte Carlo method requires the ability to draw independent realizations of a scalar random variable Z with a specified cumulative distribution function $F(z) = P(Z \leq z)$, where P is a probability measure. On a computer, the starting point for this exercise is normally a *pseudo-random number generator*, a software program that will generate a sequence of numbers uniformly distributed on $[0, 1]$ (i.e. from $\mathcal{U}(0, 1)$). Press et al. [1992] list a number of generators producing sequences of uniform numbers u_1, u_2, \dots from iterative relationships of the form

$$\begin{aligned} I_{i+1} &= (aI_i + c) \bmod(m), \\ u_{i+1} &= I_{i+1}/m. \end{aligned}$$

The externally specified starting point I_0 is the *seed* of the random number generator. In this so-called *general linear congruential generator*, the choice of the *multiplier* a , the *modulus* m , and the *increment* c must be done

with great care to ensure that the period length of the generator is large² and that the resulting algorithm is efficient on a computer. The latter, for instance, can be accomplished by setting m to be a power of 2 such that the modulo operation can be done by bit-shifting. For detailed discussion and a number of concrete algorithms (including computer code), we refer to Press et al. [1992]. The algorithms in Press et al. [1992] should suffice for most fixed income applications, but we should note the existence of more sophisticated methods that (theoretically, at least) have better performance than linear congruential generators. For instance, the so-called *Mersenne twister* proposed in Matsumoto and Nishimura [1998] has become popular, especially the specific variant MT19937 which has a period of $2^{19937} - 1$. For an extensive survey of pseudo-random number generators, see L'Ecuyer [1994].

So far we have only discussed techniques to generate $\mathcal{U}(0, 1)$ numbers, but many methods exist to convert uniformly distributed numbers into draws from the distribution F of Z . We cover a few important techniques next.

3.1.1.1 Inverse Transform Method

The idea of the inverse transform method is straightforward. Let U be a random variable uniformly distributed on $[0, 1]$, and consider setting

$$Z = F^{-1}(U), \quad (3.5)$$

where we assume that F^{-1} is well-defined, for all but a finite number of points³. As desired,

$$\mathbb{P}(Z \leq z) = \mathbb{P}(F^{-1}(U) \leq z) = \mathbb{P}(U \leq F(z)) = F(z),$$

where the last equality follows from the property of uniformly distributed random variables. The inverse transform method (3.5) is quite general, but its practical usefulness hinges on being able to compute F^{-1} fast. Many distributions allow for closed-form inversion; this includes the *exponential distribution* where $F(z) = 1 - e^{-z\lambda}$ for some positive constant λ , and the *Cauchy distribution* where $F(z) = 1/2 + (1/\pi) \arctan((z-t)/s)$ for constants t and $s > 0$.

For the important case of the Gaussian distribution, no closed-form expression for the inverse distribution exists. Nevertheless, the inverse transform

²Note that if a number $I_k = I_i$, the sequences starting from I_k and I_i are identical. In practice, we would want the generator to have *full period*, in the sense that the sequence would produce $m - 1$ distinct values before repeating the sequence.

³For discrete random variables, the distribution function is discontinuous around each of the possible (discrete) outcomes of Z . We can handle this by simply defining $F^{-1}(u) = \inf\{z : F(z) \geq u\}$.

method can still be applied as fast and extremely accurate approximations for Φ^{-1} exist. For instance, Beasley and Springer [1977] suggest the rational approximation

$$\Phi^{-1}(x) \approx \frac{\sum_{i=0}^3 a_i (x - \frac{1}{2})^{2i+1}}{1 + \sum_{i=0}^3 b_i (x - \frac{1}{2})^{2i}}, \quad 0.5 \leq x \leq 0.92, \quad (3.6)$$

for constants $a_i, b_i, i = 0, \dots, 3$, listed in Appendix 3.A. For values of x greater than 0.92, Moro [1995] proposes the approximation

$$\Phi^{-1}(x) \approx \sum_{i=1}^8 c_i [\ln(-\ln(1-x))]^i, \quad 0.92 \leq x < 1, \quad (3.7)$$

for constants $c_i, i = 0, \dots, 8$, given in Appendix 3.A. Taken together, (3.6) and (3.7) provide an approximation valid for $0.5 \leq x < 1$; when $0 < x < 0.5$ we can compute $\Phi^{-1}(x)$ by symmetry: $\Phi^{-1}(1-x) = -\Phi^{-1}(x)$. The precision of (3.6)–(3.7) is excellent⁴, with the error less than 3×10^{-9} for x in the range $x \in [\Phi(-7), \Phi(7)]$. For alternative algorithms, see for instance Acklam [2003] and Wichura [1988].

Well-known alternative methods for sampling in the Gaussian distribution include the *Box-Muller method* and the related *Marsaglia polar method* (see Press et al. [1992]).

3.1.1.2 Acceptance-Rejection Method

In cases where F^{-1} is cumbersome to compute, the so-called *acceptance-rejection method* may be preferable. To describe the method, suppose that we want to sample from a density $f(z) = dF(z)/dz$, and further suppose that we have a good method to sample from a density $e(z)$, where

$$e(z)c \geq f(z), \quad z \in \mathbb{R}, \quad (3.8)$$

for some positive constant c . By necessity, $c \geq 1$ as both e and f integrate to 1. In the acceptance-rejection method, we

1. Draw a sample Z from $e(z)$.
2. Draw an independent uniform variable U , $U \sim \mathcal{U}(0, 1)$.
3. Accept the sample Z if $U \leq f(Z)/(ce(Z))$; otherwise discard it.

⁴If even higher precision is required, we can use (3.6)–(3.7) as a guess for the root y in the equation $\Phi(y) = x$. Any number of numerical root search routines (e.g. Newton-Raphson) can then be applied to improve the precision of the solution further. Typically only one or two iterations will be required to get the solution to within machine precision on a PC.

The proof of why this algorithm works is straightforward and we omit it. Note that the third step of the acceptance-rejection method can be wasteful if too many samples need rejection. The key to the numerical efficiency of the acceptance-rejection method is thus evidently the ability to identify densities $e(z)$ that are “close” to $f(z)$, in the sense that c is close to 1 for all x . Indeed, it can easily be shown that the probability of rejecting a sample is $1/c$. Press et al. [1992] list good choices for $e(z)$ for a number of standard densities $f(z)$.

To demonstrate the mechanics of setting up an acceptance-rejection scheme for a particular distribution, let us consider sampling of a variable χ_ν^2 from a *chi-square distribution* with ν degrees of freedom. This distribution arises in a number of interest rate applications and is characterized by the cumulative distribution function

$$P(\chi_\nu^2 \leq z) = \frac{1}{2^{\nu/2} \Gamma(\nu/2)} \int_0^z e^{-y/2} y^{(\nu/2)-1} dy, \quad \nu > 0, z \geq 0,$$

where Γ is the gamma function. For reasonably large degrees of freedom ν , the chi-square density is typically bell-shaped. The chi-square distribution is a special case of the *gamma distribution* with density

$$f(z; a, b) = \frac{a(az)^{b-1} e^{-az}}{\Gamma(b)}, \quad a, b > 0, z \geq 0. \quad (3.9)$$

The chi-square distribution corresponds to $a = \frac{1}{2}$ and $b = \frac{\nu}{2}$. Rather than considering how to simulate a chi-square distribution, we will consider the more general question of how to draw from (3.9). We note that if a variable X has gamma density $f(z; 1, b)$, then aX , $a > 0$, has gamma density $f(z; a, b)$, so, in fact, it suffices to consider a simulation algorithm for the unit-scale density

$$f(z) = \frac{z^{b-1} e^{-z}}{\Gamma(b)},$$

where we assume that $b \geq 1$. One simple choice of “comparison” density for an acceptance-rejection algorithm is the exponential density

$$e(z) = \lambda e^{-\lambda z},$$

which, as mentioned earlier, can easily be simulated by inverse transform techniques. Note that

$$\frac{f(z)}{e(z)} = \frac{1}{\lambda \Gamma(b)} z^{b-1} e^{(\lambda-1)z},$$

which can be checked to have a maximum value of

$$\sup \left(\frac{f(z)}{e(z)} \right) = \frac{1}{\lambda \Gamma(b)} \left(\frac{b-1}{e(1-\lambda)} \right)^{b-1}, \quad (3.10)$$

where we must assume that $\lambda < 1$. To satisfy (3.8) we take $c = \sup(f(z)/e(z))$ and now search for the value of λ that minimizes c , thereby optimizing computational speed. It is easy to see that (3.10) is minimized for $\lambda = 1/b$, corresponding to $c = b^b e^{1-b}/\Gamma(b)$. Note that

$$\frac{f(z)}{ce(z)} = \frac{z^{b-1}}{\lambda} e^{b-1+(\lambda-1)z} b^{-b},$$

with the third step of the acceptance-rejection algorithm best done in logarithms.

The algorithm outlined above was proposed by Fishman [1976] and works best for moderate values of b . For larger values, the Gamma distribution starts looking like a bell-shaped Gaussian distribution and is no longer well-approximated by an exponential distribution. Indeed, we notice that the probability of rejection ($1/c$) is approximately $e\sqrt{b/(2\pi)}$, so of order $O(\sqrt{b})$. Modifications to the basic Fishman algorithm to accelerate sampling can be found in Cheng and Feast [1980]. Another common idea is to set $e(z)$ to the *Cauchy density*

$$e(z) = \frac{1}{s\pi \left(1 + ((z-t)/s)^2 \right)},$$

where $s > 0$ and t are constants. This distribution is bell-shaped and, as discussed earlier, can be simulated by the inverse transform method. Press et al. [1992] list computer code and references for this case. For values $b \in [0, 1]$, the acceptance-rejection technique of Ahrens and Dieter [1974] can also be used.

3.1.1.3 Composition

A third and final method to generate random variables from a given distribution function exploits known functional relationships that map variables sampled from one or more distributions to variables sampled from a target distribution. This technique is known as *composition*. A classical example of composition is the *log-normal distribution* $\mathcal{LN}(\mu, \sigma^2)$ which, as we saw earlier in Chapter 1, is defined through the relation

$$X \sim \mathcal{N}(\mu, \sigma^2) \Rightarrow e^X \sim \mathcal{LN}(\mu, \sigma^2),$$

where \sim denotes “distributed as”, and where $\mathcal{N}(\mu, \sigma^2)$ is the Gaussian distribution with mean μ and variance σ^2 . In other words, a sample Z from $\mathcal{LN}(\mu, \sigma^2)$ can be generated by drawing (by the inverse transformation method, say) a $\mathcal{N}(0, 1)$ variable X , and then setting $Z = e^{\mu+\sigma X}$.

Another classical example of a functional map is the *Student's t-distribution*, where samples can be generated by multiplying independent

samples from a standard Gaussian and a chi-square distribution; see Andersen et al. [2003] for a financial application of this. While we earlier demonstrated that the chi-square and gamma distributions can be generated by acceptance-rejection techniques, in fact we can also use composition for this. For instance, it is known that if X_1, X_2, \dots, X_ν are independent standard Gaussian variables, then

$$Z = \sum_{i=1}^{\nu} X_i^2 \quad (3.11)$$

is distributed chi-square with ν degrees of freedom. Also, if U_1, \dots, U_b are independent uniformly distributed variables, then

$$Z = -a \sum_{i=1}^b \ln U_i \quad (3.12)$$

is gamma distributed with density (3.9). For small integer-valued distribution parameters b or ν , (3.11) or (3.12) often define a faster simulation scheme than acceptance-rejection methods.

For later use, we note that the relationship (3.11) can be generalized to

$$\tilde{\chi}_{\nu}^2(\lambda) = \sum_{i=1}^{\nu} (X_i + a_i)^2$$

for a series of constants a_i , $i = 1, \dots, \nu$. The random variable $\tilde{\chi}_{\nu}^2(\lambda)$ follows a so-called *non-central chi-square distribution* with ν degrees of freedom and *non-centrality parameter* $\lambda = \sum_i a_i^2$. The distribution function is given by

$$P(\tilde{\chi}_{\nu}^2(\lambda) \leq z) = e^{-\lambda/2} \sum_{j=0}^{\infty} \frac{\left(\frac{1}{2}\lambda\right)^j}{j! \Gamma\left(\frac{\nu}{2} + j\right) 2^{(\nu/2)+j}} \int_0^z y^{\nu/2+j-1} e^{-y/2} dy, \quad (3.13)$$

an expression that also holds for non-integer ν . If $\nu > 1$, samples from a non-central chi-squared distribution can be generated by composition, using the relation

$$\tilde{\chi}_{\nu}^2(\lambda) = (Z + \sqrt{\lambda})^2 + \chi_{\nu-1}^2,$$

where Z is a standard Gaussian random variable independent of $\chi_{\nu-1}^2$. To handle the case $\nu \leq 1$, one can observe from the expression (3.13) that a non-central chi-square variable can be expressed as a regular chi-square variable $\chi_{\nu+2N}^2$, where N is an independent *Poisson-distributed* discrete variable with intensity $\lambda/2$,

$$P(N = j) = e^{-\lambda/2} \frac{(\lambda/2)^j}{j!}, \quad j = 0, 1, \dots$$

This suggests a composition rule for arbitrary ν : draw Poisson variables N (by the inverse transformation method, say) and then draw $\chi_{\nu+2N}^2$ using the methods in Section 3.1.1.2.

3.1.2 Correlated Gaussian Samples

The previous section dealt with the generation of scalar random variables. In applications, however, we may face the task of generating *vectors* of random variables, drawn from a joint multi-variate distribution. Of primary importance in financial applications is the multi-variate Gaussian distribution, so we devote this section to issues surrounding the generation of correlated Gaussian samples.

Recall that a p -dimensional Gaussian distribution $\mathcal{N}(\mu, \Sigma)$ is characterized by a p -dimensional vector-valued mean μ and a $p \times p$ symmetric, positive semi-definite⁵ covariance matrix Σ . The joint density is

$$\phi_p(z; \mu, \Sigma) = \frac{1}{(2\pi)^{p/2}(\det \Sigma)^{1/2}} \exp\left(-\frac{1}{2}(z - \mu)^\top \Sigma^{-1}(z - \mu)\right), \quad z \in \mathbb{R}^p.$$

The following result is useful:

Lemma 3.1.3 (Linear Transformation). *Let $Z \sim \mathcal{N}(\mu, \Sigma)$ be p -dimensional. Given a $d \times p$ matrix A and a d -dimensional vector B , then*

$$AZ + B \sim \mathcal{N}(A\mu + B, A\Sigma A^\top).$$

We can use this lemma as follows. Suppose that we generate p independent standard (that is, $\mathcal{N}(0, 1)$) Gaussian samples and collect them in a p -dimensional vector X . This can be accomplished using the techniques in Section 3.1.1. Clearly $X \sim \mathcal{N}(0, I)$, where I is the p -dimensional identity matrix. Define a $(p \times p)$ -dimensional matrix C satisfying

$$CC^\top = \Sigma. \tag{3.14}$$

Then

$$Z = \mu + CX$$

is distributed $\mathcal{N}(\mu, \Sigma)$.

It remains to determine a matrix C that satisfies (3.14). While there is generally an infinite number of such matrices, two particular choices are of primary importance. We discuss these below.

3.1.2.1 Cholesky Decomposition

In the Cholesky decomposition, we impose the constraint that the matrix C be lower triangular (that is, having all zeros above the diagonal), thereby conveniently reducing the number of multiplications required to compute CX to $p(1 + (p - 1)/2)$, rather than p^2 . Assuming that the matrix is positive definite (not only positive semi-definite), the Cholesky decomposition is well-defined, and given by

⁵That is, all eigenvalues of Σ are non-negative.

$$C_{i,i} = \sqrt{\Sigma_{i,i} - \sum_{k=1}^{i-1} C_{i,k}^2}, \quad i = 1, \dots, p,$$

$$C_{i,j} = \frac{1}{C_{j,j}} \left(\Sigma_{i,j} - \sum_{k=1}^{j-1} C_{i,k} C_{j,k} \right), \quad j = 1, \dots, p-1, \quad j < i.$$

For instance, if

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix},$$

where $\rho \in [-1, 1]$ and $\sigma_1, \sigma_2 > 0$, then

$$C = \begin{pmatrix} \sigma_1 & 0 \\ \sigma_2\rho & \sigma_2\sqrt{1-\rho^2} \end{pmatrix},$$

a result that we have already used in Section 2.11. Press et al. [1992], among others, list computer code implementing the relations above.

If the matrix Σ is only positive semi-definite (but not positive definite), the Cholesky decomposition will fail. In this case, linear algebra tells us that the matrix Σ is rank-deficient, with rank $r < p$. As such, we must be able to set $Z = \mu + MY$, where M is a $p \times r$ matrix and $Y \sim \mathcal{N}(0, \Sigma_Y)$ is r -dimensional, with the covariance matrix having full rank r . Using Cholesky composition instead on Σ_Y , we can find a lower diagonal matrix C_Y satisfying $C_Y C_Y^\top = \Sigma_Y$. Thus, in this case

$$Z = \mu + MC_Y X$$

where X is a vector of r (not p) independent standard Gaussian samples. The matrix M can be found by the singular value decomposition (SVD) algorithm, see Press et al. [1992], or the algorithm in the next section.

3.1.2.2 Eigenvalue Decomposition

As an alternative to Cholesky decomposition, we can also consider diagonalizing Σ through an eigenvalue decomposition. Here, we write

$$\Sigma = E\Lambda E^\top, \tag{3.15}$$

where Λ is a diagonal matrix of eigenvalues λ_i , $i = 1, \dots, p$, and the columns of E contain the orthonormal eigenvectors of Σ . Some eigenvalues may be zero, if Σ is rank-deficient (positive semi-definite). Comparison with (3.14) implies that one choice of C is

$$C = E\sqrt{\Lambda} = E \begin{pmatrix} \sqrt{\lambda_1} & 0 & \cdots & 0 \\ 0 & \sqrt{\lambda_2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \sqrt{\lambda_p} \end{pmatrix}. \tag{3.16}$$

The eigenvalue decomposition (3.15) is relatively straightforward, at least as eigenvalue problems go, due to the fact that Σ is symmetric and positive semi-definite; see Press et al. [1992] for an algorithm. While both Cholesky decomposition and eigenvalue decompositions have computational complexity $O(n^3)$, in practice the Cholesky method is often much faster than the eigenvalue method, making the Cholesky method preferable in practice. Nevertheless, decompositions of the type (3.16) have certain appealing theoretical properties that shall be useful later, so the next section explores (3.16) further.

3.1.3 Principal Components Analysis (PCA)

Consider a p -dimensional Gaussian variable Z with a given covariance matrix Σ . Assume, with no loss of generality, that the mean of Z is 0 and that Σ has full rank (positive definite). Consider now writing, as an approximation,

$$Z \approx DX, \quad (3.17)$$

where X is an r -dimensional vector of independent standard Gaussian variables, $r \leq p$, and D is a $(p \times r)$ -dimensional matrix. How should we choose D in an optimal way?

First, we obviously need to define what constitutes an “optimal” approximation in (3.17). We here have in mind L^2 closeness of the covariance matrix DD^\top to Σ (see Lemma 3.1.3), so let us define the optimal D^* as the matrix that minimizes the norm

$$f(D) = \text{tr} \left((\Sigma - DD^\top) (\Sigma - DD^\top)^\top \right).$$

This is just the matrix representation of the usual Frobenius norm on the squared differences between Σ and DD^\top . The value of D that minimizes $f(D)$ can be shown to be

$$D^* = E_r \sqrt{\Lambda_r}, \quad (3.18)$$

where Λ_r is an $r \times r$ diagonal matrix containing the *largest* r eigenvalues of Σ , and E_r is a $p \times r$ matrix of r p -dimensional eigenvectors corresponding to the eigenvalues in Λ_r .

Equipped with the optimal D , we now go back to the approximation (3.17) and write

$$Z \approx \tilde{Z} \triangleq E_r \sqrt{\Lambda_r} X = \sqrt{\lambda_1} e_1 X_1 + \sqrt{\lambda_2} e_2 X_2 + \dots + \sqrt{\lambda_r} e_r X_r, \quad (3.19)$$

where e_i denotes the i -th column of E_r and the λ_i 's are the eigenvalues, sorted in decreasing order of magnitude. The (deterministic) vector e_i is known as the i -th *principal component* of Z , and the (random) variable $\sqrt{\lambda_i} X_i$ as the i -th *principal factor*. With (3.19), we have $\text{tr}(\text{Cov}(Z, Z)) = E(Z^\top Z) = \sum_{i=1}^p \lambda_i$ and $\text{tr}(\text{Cov}(\tilde{Z}, \tilde{Z})) = E(\tilde{Z}^\top \tilde{Z}) = \sum_{i=1}^r \lambda_i$, i.e. the first r terms in the decomposition (3.19) explain a fraction

$$\frac{\sum_{i=1}^r \lambda_i}{\sum_{i=1}^p \lambda_i}$$

of the sum of the diagonal elements of the covariance matrix of Z . Principal components decomposition will thus result in a loss of total variance, unless the covariance matrix is either rank-deficient (i.e. has eigenvalues that are strictly zero), or we use a full set of principal components ($p = r$). In many cases of interest to us here, the loss of variance can be small, even if r is a modest number, e.g. 2 or 3. We notice that the covariance matrix for Z , as approximated by (3.19), will be *rank-deficient*, as the number r of non-zero eigenvalues is less than p .

While we have used a setting with Gaussian variables to motivate our treatment of principal components analysis (PCA), it is, in fact, a generically useful tool for uncovering the structure of large-dimensional random vectors, and replacing them with more manageable, lower-dimensional variables; see, e.g., Theil [1971] for more details and an application to empirical non-Gaussian data. Also, PCA identifies which directions of a multi-dimensional random variable are “important”, potentially allowing us to allocate computational resources in an intelligent manner. One example of this is shown later in this chapter, in Section 3.2.10.

3.2 Generation of Sample Paths

So far, we have assumed that random variables are characterized by a known distribution function. In most of our applications, however, the random variables $g(T)/N(T)$ used in the basic pricing equation (3.1) are specified through an SDE or, more generally, an Ito process. In this section, we shall discuss Monte Carlo simulation of such processes. We start out with a motivating example, set in the Black-Scholes-Merton economy.

3.2.1 Example: Asian Basket Options in Black-Scholes Economy

Consider a dividend-free stock S , with Black-Scholes dynamics

$$dS(t)/S(t) = r dt + \sigma dW(t), \quad (3.20)$$

where $W(t)$ is a Brownian motion in the risk-neutral measure \mathbb{Q} , r is a constant interest rate, and σ is a constant volatility. Let there be given an increasing set of observation times $\{t_1, t_2, \dots, t_m\}$, with $t_m = T$, and define the \mathcal{F}_T -measurable (discretely observed) stock average

$$A(T) = \frac{1}{m} \sum_{i=1}^m S(t_i). \quad (3.21)$$

An *Asian* (or *average rate*) call option with strike K is defined by the terminal payout

$$g(T) = (A(T) - K)^+; \quad (3.22)$$

we wish to price this option by Monte Carlo simulation.

As discussed earlier (see (1.39)), the geometric Brownian motion process (3.20) allows us to express S directly in terms of the Brownian motion,

$$S(t) = S(0)e^{rt - \frac{1}{2}\sigma^2 t + \sigma W(t)}, \quad t > 0,$$

whereby, with $\Delta_i \triangleq t_i - t_{i-1}$ and $t_0 = 0$,

$$S(t_i) = S(t_{i-1}) \exp \left(\left[r - \frac{1}{2}\sigma^2 \right] \Delta_i + \sigma [W(t_i) - W(t_{i-1})] \right),$$

$i = 1, \dots, m$. By the properties of Brownian motion, the increments $W(t_i) - W(t_{i-1})$ are independent Gaussian variables distributed as $\mathcal{N}(0, \Delta_i)$. For the purposes of Monte Carlo simulation, we can therefore write

$$S(t_i) = S(t_{i-1}) \exp \left(\left(r - \frac{1}{2}\sigma^2 \right) \Delta_i \right) \exp \left(\sigma \sqrt{\Delta_i} Z_i \right), \quad i = 1, \dots, m, \quad (3.23)$$

where the Z_i are independent standard $\mathcal{N}(0, 1)$ Gaussian random variables. To produce a single sample draw of $g(T)$, we thus

1. Draw independent standard Gaussian samples $Z_i, i = 1, \dots, m$ (see Section 3.1.1).
2. Starting from $S(0)$, generate $S(t_i), i = 1, \dots, m$, from the iteration (3.23).
3. Compute $g(T)$ from (3.21)–(3.22).

Repeating this procedure n times (with Gaussian samples independent from one path to the next), we can generate n random samples g_1, g_2, \dots, g_n of $g(T)$. Our estimate of the time 0 price of the Asian option is then, from (3.3) with $N(t) = e^{rt}$ and non-random,

$$\bar{V}(0) = e^{-rT} \frac{1}{n} \sum_{j=1}^n g_j.$$

Asymptotic confidence intervals can be computed from (3.4). The pricing algorithm involves drawing mn Gaussian variables, so the computational cost of the pricing algorithm is $O(mn)$.

Increasing the complexity, let us now consider an Asian option on a p -dimensional basket of stocks S_1, S_2, \dots, S_p , each following geometric Brownian motion,

$$dS_k(t)/S_k(t) = r dt + \sigma_k dW_k(t), \quad k = 1, \dots, p.$$

The Brownian motions W_k and W_j are assumed correlated with constant correlation coefficient $\rho_{k,j}$, $j, k = 1, \dots, p$, $j \neq k$. Define a unit-weighted basket price as

$$B(t) = \sum_{k=1}^p S_k(t),$$

and set the terminal Asian option payout to be

$$g(T) = \left(\frac{1}{m} \sum_{i=1}^m B(t_i) - K \right)^+, \quad (3.24)$$

where the time line $\{t_i\}$ is as before. Equivalent to (3.23), we draw sample paths for each asset according to the prescription

$$S_k(t_i) = S_k(t_{i-1}) \exp \left(\left(r - \frac{1}{2} \sigma_k^2 \right) \Delta_i + \sigma_k \sqrt{\Delta_i} Z_{k,i} \right), \quad (3.25)$$

$$i = 1, \dots, m, \quad k = 1, \dots, p,$$

where the $Z_{k,i}$ are Gaussian samples, independently drawn at each time step but correlated across k 's. Let C be the Cholesky decomposition of the correlation matrix $\{\rho_{k,j}\}$ (see Section 3.1.2.1), in which case we can generate the correlated sample vectors $Z_i = (Z_{1,i}, Z_{2,i}, \dots, Z_{p,i})^\top$ as

$$Z_i = CX_i$$

for a p -dimensional vector X_i of independent Gaussian samples. Given joint sample paths of all basket component assets S_k , $k = 1, \dots, p$, pricing of the Asian basket option proceeds as above, substituting (3.24) for (3.22).

Completion of (3.25) requires pm samples to complete a full path of all p assets, making the total computational effort of an n -sample Monte Carlo scheme $O(nmp)$, with the (probabilistic) convergence order $O(n^{-1/2})$ and dependent only on n . As mentioned earlier, the linearity of computational cost on the dimension of the asset vector p compares favorably to the exponential growth in p of finite difference schemes. Notice also the ease with which the Monte Carlo scheme is able to incorporate path-dependence.

3.2.2 Discretization Schemes, Convergence, and Stability

At the heart of the example in Section 3.2.1 was an iterative scheme for the production of a sample path for a vector-valued SDE; see (3.25). For the simple Black-Scholes model, SDE state variables (stock prices) could be expressed analytically in terms of independent increments of a Brownian motion, making path generation straightforward. In practice, however, we are often working with SDEs that do not permit closed-form solution. In such cases, we need to *time-discretize* the SDE, much the same way as we did for the numerical solution of PDEs.

In the next few sections, we shall consider a few important SDE discretization schemes. Before moving on to this, it is useful to discuss the sense in which we consider a discretization scheme to converge to the true SDE solution. For this, consider a vector-valued SDE

$$dX(t) = \mu(t, X(t)) dt + \sigma(t, X(t)) dW(t), \quad (3.26)$$

where $X(t)$ is p -dimensional, W is a d -dimensional vector of independent Brownian motions, and $\mu : [0, T] \times \mathbb{R}^p \rightarrow \mathbb{R}^p$ and $\sigma : [0, T] \times \mathbb{R}^p \rightarrow \mathbb{R}^{p \times d}$ satisfy the usual regularity conditions. Consider an equidistant⁶ time grid $\{0, \Delta, 2\Delta, \dots, m\Delta\}$, the number of references and let \widehat{X} be an approximation to X , based on some kind of time-discretization scheme on the grid $\{i\Delta\}$. For simplicity of notation, set $\widehat{X}_i \triangleq \widehat{X}(i\Delta)$. We say that the underlying approximation is *weakly consistent* if there exists a function $c(\Delta)$ with

$$\lim_{\Delta \downarrow 0} c(\Delta) = 0$$

such that (dropping the measure superscript on the expectation operator)

$$\mathbb{E} \left(\left| \mathbb{E} \left(\Delta^{-1} (\widehat{X}_{i+1} - \widehat{X}_i) \middle| \mathcal{F}_{i\Delta} \right) - \mu(i\Delta, \widehat{X}_i) \right|^2 \right) \leq c(\Delta), \quad (3.27)$$

and

$$\begin{aligned} \mathbb{E} \left(\left| \mathbb{E} \left(\Delta^{-1} (\widehat{X}_{i+1} - \widehat{X}_i) (\widehat{X}_{i+1} - \widehat{X}_i)^\top \middle| \mathcal{F}_{i\Delta} \right) \right. \right. \\ \left. \left. - \sigma(i\Delta, \widehat{X}_i) \sigma(i\Delta, \widehat{X}_i)^\top \right|^2 \right) \leq c(\Delta), \end{aligned} \quad (3.28)$$

for all $i = 0, \dots, m-1$. The notion of weak consistency⁷ thus amounts to requiring that the mean and variance of the increments of the approximating process be close to those of the true SDE solution.

A concept related to consistency is the notion of *weak convergence*. We say that an approximate solution converges weakly to X at time $T = m\Delta$ with respect to a class \mathcal{C} of test functions $g : \mathbb{R}^p \rightarrow \mathbb{R}$ if

$$\lim_{\Delta \downarrow 0} \left| \mathbb{E}(g(X(T))) - \mathbb{E}(g(\widehat{X}(T))) \right| = 0, \quad (3.29)$$

for all $g \in \mathcal{C}$. Notice that the limit necessarily involves $m \rightarrow \infty$.

⁶To keep notation manageable, we use a constant time step Δ in most of this chapter. All results are, however, easily extendable to non-equidistant grids.

⁷*Strong consistency* (which is of little use to us in this book) requires that (3.27) is satisfied, and that the variance of the difference between increments of the true process and the approximation vanish. The second requirement is stronger than (3.28).

The class of test functions used in (3.29) is normally always in the set \mathcal{C}_P^l of functions with polynomially bounded⁸ derivatives of order $0, 1, \dots, l$ with maximum power l . We say that a scheme converges with *weak order* β if, for all $g \in C_P^{2(\beta+1)}$, (3.29) can be strengthened to

$$\left| \mathbb{E}(g(X(T))) - \mathbb{E}\left(g(\hat{X}(T))\right) \right| \leq c\Delta^\beta, \quad (3.30)$$

for all $\Delta \in (0, \Delta_0)$, where Δ_0 and c are constants and c does not depend on Δ (but may depend on g).

One would generally expect that a weakly consistent scheme is weakly convergent. Indeed, this can be established to be the case under certain additional regularity conditions. We will not list the exact result here, but refer to Kloeden and Platen [2000], Theorem 9.7.4.

Finally, a brief word on stability of a time-discretized SDE. A commonly used definition of *A*-stability focuses on the behavior of a discretized test SDE of the type

$$dX(t) = \lambda X(t) dt + dW(t), \quad (3.31)$$

where λ is a complex-valued constant with real part $\text{Re}(\lambda) < 0$. We suppose that a discretization scheme can be represented as

$$\hat{X}_{i+1} = \hat{X}_i G(\lambda\Delta) + Z_i^\Delta, \quad i = 0, 1, \dots, m-1, \quad (3.32)$$

where G is a mapping of the complex plane onto itself and the Z_i^Δ 's are random variables independent of the \hat{X}_i 's. In this case, the *region of stability* for a scheme is the set of $\lambda\Delta$ for which $\text{Re}(\lambda) < 0$ and

$$|G(\lambda\Delta)| < 1. \quad (3.33)$$

Similar to the definition used for finite difference scheme discretizations, we say that an SDE time-discretization scheme is *A-stable*, if the region of stability includes all values of λ with $\text{Re}(\lambda) < 0$ and all $\Delta > 0$.

3.2.3 The Euler Scheme

An obvious first scheme to discretize (3.26) treats both dt and $dW(t)$ fully explicitly, evaluating all SDE coefficients on time step $[i\Delta, i\Delta + \Delta]$ at the left interval point $i\Delta$. In other words, we write, starting from $\hat{X}_0 = X(0)$,

$$\begin{aligned} \hat{X}_{i+1} &= \hat{X}_i + \mu(i\Delta, \hat{X}_i) \Delta + \sigma(i\Delta, \hat{X}_i) (W(i\Delta + \Delta) - W(i\Delta)), \\ i &= 0, 1, \dots, m-1. \end{aligned} \quad (3.34)$$

⁸A function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is polynomially bounded if $|f(x)| \leq k(1 + |x|^q)$, $x \in \mathbb{R}^p$, for constants k and q .

With this scheme, Monte Carlo generation of paths is straightforward and involves, as in Section 3.2.1, replacing the increments $W(i\Delta + \Delta) - W(i\Delta)$ with $Z_i\sqrt{\Delta}$, for a d -dimensional vector of independent standard Gaussian samples Z_i .

The discretization scheme (3.34) is known as the *Euler scheme*, sometimes also called the *Euler-Maruyama* scheme. The Euler scheme is easy to implement and is a true workhorse that we will often use in this book. We note that the scheme is weakly consistent, as

$$\mathbb{E} \left(\left| \mathbb{E} \left(\Delta^{-1} (\widehat{X}_{i+1} - \widehat{X}_i) \mid \mathcal{F}_{i\Delta} \right) - \mu(i\Delta, \widehat{X}_i) \right|^2 \right) = 0,$$

and

$$\begin{aligned} \mathbb{E} \left(\left| \mathbb{E} \left(\Delta^{-1} (\widehat{X}_{i+1} - \widehat{X}_i) (\widehat{X}_{i+1} - \widehat{X}_i)^\top \mid \mathcal{F}_{i\Delta} \right) \right. \right. \\ \left. \left. - \sigma(i\Delta, \widehat{X}_i) \sigma(i\Delta, \widehat{X}_i)^\top \right|^2 \right) = O(\Delta^2). \end{aligned}$$

While one might believe that the explicit discretization of the diffusion term — which is only accurate to order $O(\sqrt{\Delta})$ — would give the scheme weak convergence order⁹ $1/2$, in fact we typically have that the Euler scheme has *weak convergence order* $\beta = 1$. We note that for this result to hold, however, regularity conditions on μ and σ stronger than those of the existence and uniqueness results (Theorem 1.6.1) are needed. For instance, in the case where μ and σ are functions of X alone, Theorem 9.7.6 in Kloeden and Platen [2000] requires that μ and σ be four times continuously differentiable with polynomial growth and uniformly bounded derivatives. See also their Theorem 15.4.2 for a more general result.

Given that the Euler scheme is fully explicit, our experience from finite difference methods suggests that the scheme may have stability problems. To investigate, we follow Section 3.2.2 and consider the test SDE

$$dX(t) = \lambda X(t)dt + dW(t),$$

which is discretized as

$$\widehat{X}_{i+1} = \widehat{X}_i (1 + \lambda\Delta) + \sqrt{\Delta} Z_i, \quad (3.35)$$

where Z_i 's are standard Gaussian. Comparison to (3.32) and (3.33) shows that the region of stability for the Euler scheme is

$$|(1 + \lambda\Delta)| < 1, \quad \text{Re}(\lambda) < 0,$$

which is the unit disc in the complex plane centered at $\lambda\Delta = -1$. For a given λ , there are thus restrictions on how big a time step Δ can be used.

⁹The so-called *strong convergence order* of the Euler scheme is in fact only $1/2$. The concept of strong convergence order is defined in Kloeden and Platen [2000] and is of little importance to applications in this book.

3.2.3.1 Linear-Drift SDEs

The restricted stability region of the Euler scheme can be a practical concern. For instance, SDEs of the important type

$$dX(t) = \kappa(\theta(t) - X(t)) dt + \sigma(t, X(t)) dW(t) \quad (3.36)$$

arise quite frequently in fixed income modeling, and in cases where κ is big (which is often the case for, say, stochastic volatility models such as those covered in Chapters 8, 9 and 13) the Euler scheme can become unstable and return meaningless results. One way to solve the problem is to switch to an implicit scheme (see next section), but in the case (3.36) we can use the fact that the drift term can be removed by a simple change of variable. For instance, for the case where $X(t)$ is scalar we can set

$$Y(t) = e^{\kappa t} X(t) - \kappa \int_0^t e^{\kappa u} \theta(u) du,$$

such that, from Ito's lemma,

$$\begin{aligned} dY(t) &= e^{\kappa t} \sigma(t, X(t)) dW(t) \\ &= e^{\kappa t} \sigma \left(t, e^{-\kappa t} \left(Y(t) + \kappa \int_0^t e^{\kappa u} \theta(u) du \right) \right) dW(t). \end{aligned}$$

Euler simulation of the process for $Y(t)$, rather than for $X(t)$, will center X around its analytically known mean

$$\mathbb{E}(X((i+1)\Delta) | X(i\Delta)) = e^{-\kappa\Delta} X(i\Delta) + \kappa \int_{i\Delta}^{i\Delta+\Delta} e^{-\kappa((i+1)\Delta-u)} \theta(u) du$$

and will often alleviate any stability problems.

3.2.3.2 Log-Euler Scheme

One potential problem with the pure Euler scheme (3.34) is the fact that all increments are locally Gaussian, thereby implying a non-zero probability of \hat{X} crossing zero and becoming negative. Many SDEs, however, are known to produce only non-negative solutions, and the functions μ and σ may not allow for negative arguments. This, for instance, is the case for the square-root process

$$dX(t) = \sqrt{X(t)} dW(t), \quad X(0) > 0,$$

where the Euler scheme cannot be directly applied. Some authors (e.g., Kloeden and Platen [2000]) suggest heuristic modifications of the Euler scheme, such as

$$\widehat{X}_{i+1} = \widehat{X}_i + \sqrt{|\widehat{X}_i|} (W(i\Delta + \Delta) - W(i\Delta)),$$

but ultimately this is not very satisfying and the resulting scheme will often have large errors¹⁰. An alternative is to introduce an invertible transformation $X(t) = f(Y(t))$, with $f : \mathbb{R} \rightarrow \mathbb{R}_+$, and then apply the Euler scheme to Y , at each step recovering X as $f(Y)$. In finance applications, where many processes are based on SDEs that bear some resemblance to geometric Brownian motion, an often-used choice for f is $f(y) = e^y$. The resulting scheme is known as the *log-Euler scheme*.

Consider the SDE (3.26) and assume for simplicity that X is scalar (if X is vector valued, the log-transform can be applied to all, or a few selected, components of X). Set $X(t) = \exp(Y(t))$, such that $Y(t) = \ln(X(t))$. The process for Y then follows from Ito's lemma:

$$dY(t) = \left(\frac{\mu(t, X(t))}{X(t)} - \frac{1}{2} \frac{\sigma(t, X(t))^2}{X(t)^2} \right) dt + \frac{\sigma(t, X(t))}{X(t)} dW(t), \quad X(t) = e^{Y(t)}.$$

Writing out a standard Euler scheme for Y and making the transformation $\widehat{X}_i = \exp(\widehat{Y}_i)$ gives us the (scalar) log-Euler scheme for X :

$$\widehat{X}_{i+1} = \widehat{X}_i \exp \left(\left(\frac{\mu(t, \widehat{X}_i)}{\widehat{X}_i} - \frac{1}{2} \frac{\sigma(t, \widehat{X}_i)^2}{\widehat{X}_i^2} \right) \Delta + \frac{\sigma(t, \widehat{X}_i)}{\widehat{X}_i} Z_i \sqrt{\Delta} \right),$$

where $Z_i \sim \mathcal{N}(0, 1)$. Generalizations of the technique above to situations where the valid range of X is some general set \mathcal{C} are obvious and involve identifying an invertible mapping function $f : \mathbb{R} \rightarrow \mathcal{C}$, preferably one that can be inverted analytically. For instance, if $\mathcal{C} = [a, \infty)$, we could use $f(y) = a + e^y$.

3.2.4 The Implicit Euler Scheme

The implicit Euler scheme for the vector-valued SDE (3.26) takes the form

$$\widehat{X}_{i+1} = \widehat{X}_i + \mu(i\Delta + \Delta, \widehat{X}_{i+1}) \Delta + \sigma(i\Delta, \widehat{X}_i) (W(i\Delta + \Delta) - W(i\Delta)), \quad (3.37)$$

for $i = 0, 1, \dots, m - 1$. We highlight the fact that the drift coefficient μ is now evaluated at time $i\Delta + \Delta$, rather than at time $i\Delta$. It is easy to show that the implicit Euler scheme is consistent. Under regularity conditions, it can also be shown that the weak convergence order is $\beta = 1$, just as was the case for the explicit Euler scheme.

The main advantage of the implicit Euler scheme over the explicit Euler scheme is numerical stability. To examine the region of stability for the implicit Euler scheme, consider again the test SDE

¹⁰For a dedicated treatment of the rather delicate problem of simulating square-root process, see Chapter 9.

$$dX(t) = \lambda X(t) dt + dW(t).$$

It will now be discretized as (compare to (3.35))

$$\widehat{X}_{i+1} = \widehat{X}_i + \widehat{X}_{i+1} \lambda \Delta + \sqrt{\Delta} Z_i,$$

or

$$\widehat{X}_{i+1} (1 - \lambda \Delta) = \widehat{X}_i + \sqrt{\Delta} Z_i.$$

Comparison to (3.32) and (3.33) shows that now

$$G(\lambda \Delta) = \frac{1}{1 - \lambda \Delta}$$

such that the stability criterion $|G(\lambda \Delta)| < 1$ is satisfied for any value of $\lambda \Delta$ where $\text{Re}(\lambda) < 0$. In other words, the implicit scheme is A -stable.

3.2.4.1 Implicit Diffusion Term

The reader may at this point wonder why the implicit scheme (3.37) only discretized the drift term (μ) implicitly, and not the diffusion term (σ). The answer lies in the differences between a regular Riemann integral and the stochastic integral. Recall in particular that the stochastic Ito integral is defined to be non-anticipative, in the sense that the integrand is always evaluated “to the left” on any partitions of the Brownian motion. As a consequence, if $\sigma(i\Delta, \widehat{X}_i)$ were replaced with $\sigma(i\Delta + \Delta, \widehat{X}_{i+1})$ in (3.37), the resulting scheme would not be weakly consistent, in the sense defined earlier. To illustrate this point, just consider the simple scalar process

$$dX(t) = \sigma X(t) dW(t),$$

which we contemplate discretizing as

$$\widehat{X}_{i+1} = \widehat{X}_i + \sigma \widehat{X}_{i+1} (W(i\Delta + \Delta) - W(i\Delta)),$$

or

$$\widehat{X}_{i+1} (1 - \sigma Z_i \sqrt{\Delta}) = \widehat{X}_i, \quad i = 0, \dots, m-1. \quad (3.38)$$

Here, a first difficulty arises: the term $(1 - \sigma Z_i \sqrt{\Delta})$ may become 0 (or very close to zero) if Z_i is an (unbounded) Gaussian variable. For fully implicit discretization schemes, it becomes necessary to use a bounded approximation to the Brownian motion. As discussed in Kloeden and Platen [2000], weak convergence order is preserved if in (3.38) we set the Z_i to be independent binomial variables with

$$P(Z_i = 1) = P(Z_i = -1) = \frac{1}{2}.$$

We assume that $1 - \sigma \sqrt{\Delta} > 0$. Rearranging and Taylor-expanding, we get

$$\begin{aligned}
\frac{\widehat{X}_{i+1} - \widehat{X}_i}{\Delta} &= \frac{\widehat{X}_i}{\Delta} \left(\frac{1}{1 - \sigma Z_i \sqrt{\Delta}} - 1 \right) \\
&= \frac{\widehat{X}_i}{\Delta} \left(1 + \sigma Z_i \sqrt{\Delta} + \sigma^2 Z_i^2 \Delta + O(\sigma^3 Z_i^3 \Delta^{3/2}) - 1 \right) \\
&= \widehat{X}_i \left(\sigma Z_i \Delta^{-1/2} + \sigma^2 Z_i^2 + O(\sigma^3 Z_i^3 \Delta^{1/2}) \right)
\end{aligned}$$

such that

$$\mathbb{E} \left(\frac{\widehat{X}_{i+1} - \widehat{X}_i}{\Delta} \middle| \widehat{X}_i \right) = \widehat{X}_i \left(\sigma^2 + O(\Delta^{1/2}) \right).$$

Clearly, this will cause a violation of the consistency condition (3.27).

In the example above, we notice that consistency can be restored if the drift of the original SDE is changed from 0 to $-\sigma^2 X(t)$ before the “doubly” implicit Euler discretization is employed. More generally, it is not difficult to show that (3.37) can be modified to treat the diffusion term implicitly, provided that the drift of the original vector-valued SDE (3.26) is first changed from μ to

$$\bar{\mu} = \mu - \sum_{j=1}^d \sum_{k=1}^p (\sigma_{X_k})_{\cdot,j} \sigma_{k,j}$$

where the p -dimensional vector $(\sigma_{X_k})_{\cdot,j}$ is the j -th column of the $(p \times d)$ -dimensional matrix $\sigma_{X_k} = \{\partial \sigma_{i,j} / \partial X_k\}$. Inspired by the theta methods of Chapter 2, we can, in fact, introduce a family of discretizations

$$\begin{aligned}
\widehat{X}_{i+1} &= \widehat{X}_i + \left[(1 - \theta) \bar{\mu}_\eta \left(i\Delta, \widehat{X}_i \right) + \theta \bar{\mu}_\eta \left(i\Delta + \Delta, \widehat{X}_{i+1} \right) \right] \Delta \\
&\quad + \left[(1 - \eta) \sigma \left(i\Delta, \widehat{X}_i \right) + \eta \sigma \left(i\Delta + \Delta, \widehat{X}_{i+1} \right) \right] Z_i \sqrt{\Delta}, \quad (3.39)
\end{aligned}$$

where the Z_i are binomially distributed variables, $\theta, \eta \in [0, 1]$ are parameters, and

$$\bar{\mu}_\eta = \mu - \eta \sum_{j=1}^d \sum_{k=1}^p (\sigma_{X_k})_{\cdot,j} \sigma_{k,j}. \quad (3.40)$$

As it turns out, all these schemes theoretically have identical convergence order $\beta = 1$, but in practice some choices of θ, η may turn out to work better than others. We shall discuss methods to raise the theoretical convergence order in Section 3.2.6. The scheme (3.39) can be verified to be A -stable for $\theta \in [1/2, 1]$.

3.2.5 Predictor-Corrector Schemes

A closer examination of the implicit Euler scheme (3.37) demonstrates the need to recover $\widehat{X}(i\Delta + \Delta)$ as the vector-valued root of a possibly non-linear equation. In general, this must be done numerically (using, say, the

Newton-Raphson algorithm), causing a severe deterioration of computational performance. An alternative is to use the explicit Euler scheme as a *predictor* and the implicit scheme as a *corrector*, much the same way we used explicit finite difference approximations as predictors in the Craig-Sneyd algorithm of Section 2.11. Moving straight to the general implicit discretization family (3.39), we write the predictor-corrector as

$$\bar{X}_{i+1} = \hat{X}_i + \mu(i\Delta, \hat{X}_i) \Delta + \sigma(i\Delta, \hat{X}_i) (W(i\Delta + \Delta) - W(i\Delta)), \quad (3.41)$$

$$\begin{aligned} \hat{X}_{i+1} &= \hat{X}_i + [(1-\theta)\bar{\mu}_\eta(i\Delta, \hat{X}_i) + \theta\bar{\mu}_\eta(i\Delta + \Delta, \bar{X}_{i+1})] \Delta \\ &\quad + [(1-\eta)\sigma(i\Delta, \hat{X}_i) + \eta\sigma(i\Delta + \Delta, \bar{X}_{i+1})] (W(i\Delta + \Delta) - W(i\Delta)), \end{aligned} \quad (3.42)$$

where $\theta, \eta \in [0, 1]$, and $\bar{\mu}_\eta$ is as given in (3.40). It is understood that the Brownian motion increments in (3.41) and (3.42) are to be identical.

For sufficiently smooth coefficients, it can be shown that the predictor-corrector scheme (3.41)–(3.42) converges weakly with order $\beta = 1$, independent of the choice of θ and η . As for stability, discretization of (3.31) leads to

$$\begin{aligned} \bar{X}_{i+1} &= \hat{X}_i (1 + \lambda\Delta) + W(i\Delta + \Delta) - W(i\Delta), \\ \hat{X}_{i+1} &= \hat{X}_i + [(1-\theta)\lambda\hat{X}_i + \theta\lambda\bar{X}_{i+1}] \Delta + W(i\Delta + \Delta) - W(i\Delta) \\ &= \hat{X}_i (1 + \lambda\Delta(1 + \theta\lambda\Delta)) + (W(i\Delta + \Delta) - W(i\Delta))(1 + \theta\lambda\Delta). \end{aligned}$$

The region of stability can be verified to be

$$|1 + \lambda\Delta(1 + \theta\lambda\Delta)| < 1, \quad \text{Re}(\lambda) < 0.$$

For $\theta = \frac{1}{2}$, the stability criterion above is identical to that of the classical *Heun scheme* (or *modified trapezoidal scheme*) used for ordinary differential equations. Indeed, the predictor-corrector scheme above can be seen as an adaptation of this scheme for SDEs. We note that SDE adaptations of more sophisticated ODE solvers (such as Runge-Kutta) are also possible, but this goes beyond the scope of this text.

3.2.6 Ito-Taylor Expansions and Higher-Order Schemes

Despite our various efforts at centering derivatives, none of the schemes listed above theoretically attain second-order weak convergence. To develop such schemes, we need to delve further into adapting classical Taylor expansions to the rules of stochastic (Ito) calculus. As we shall ultimately not have much use for higher-order schemes, we keep the treatment informal and limit ourselves to the scalar case where $p = d = 1$ in (3.26).

3.2.6.1 Ordinary Taylor Expansion of ODEs

To gain intuition, start by setting $\sigma = 0$ in (3.26), such that we first deal with an ordinary ODE

$$dX(t) = \mu(t, X(t)) dt. \quad (3.43)$$

For a given value of t , we can use Taylor's theorem to write

$$X(t + \Delta) = X(t) + \frac{dX(t)}{dt} \Delta + \frac{1}{2} \frac{d^2 X(t)}{dt^2} \Delta^2 + O(\Delta^3),$$

where we stop at order $O(\Delta^3)$. We notice that

$$\frac{dX(t)}{dt} = \mu(t, X(t))$$

and

$$\begin{aligned} \frac{d^2 X(t)}{dt^2} &= \frac{\partial}{\partial t} \mu(t, X(t)) + \frac{\partial}{\partial x} \mu(t, X(t)) \cdot \frac{dX(t)}{dt} \\ &= \left(\frac{\partial}{\partial t} + \mu(t, X(t)) \frac{\partial}{\partial x} \right) \mu(t, X(t)). \end{aligned}$$

Setting

$$\mathcal{L} \triangleq \frac{\partial}{\partial t} + \mu \frac{\partial}{\partial x},$$

we thus have

$$X(t + \Delta) = X(t) + \mu(t, X(t)) \Delta + \frac{1}{2} \mathcal{L} \mu(t, X(t)) \Delta^2 + O(\Delta^3). \quad (3.44)$$

Another way to develop (3.44) proceeds by iteration on the integral representation

$$X(t + \Delta) = X(t) + \int_t^{t+\Delta} \mu(u, X(u)) du. \quad (3.45)$$

First we recognize that (as seen above)

$$d\mu(t, X(t)) = \mathcal{L} \mu(t, X(t)) dt$$

such that

$$\mu(u, X(u)) = \mu(t, X(t)) + \int_t^u \mathcal{L} \mu(s, X(s)) ds, \quad u > t. \quad (3.46)$$

Inserting this into (3.45) gives

$$X(t + \Delta) = X(t) + \mu(t, X(t)) \int_t^{t+\Delta} du + \int_t^{t+\Delta} \int_t^u \mathcal{L} \mu(s, X(s)) ds du.$$

Applied to $\mathcal{L}\mu(s, X(s))$ the steps that lead to (3.46) yield

$$\mathcal{L}\mu(s, X(s)) = \mathcal{L}\mu(t, X(t)) + \int_t^s \mathcal{L}^2\mu(v, X(v)) dv, \quad s > t,$$

such that

$$\begin{aligned} X(t + \Delta) &= X(t) + \mu(t, X(t)) \int_t^{t+\Delta} du + \mathcal{L}\mu(t, X(t)) \int_t^{t+\Delta} \int_t^u ds du \\ &\quad + \int_t^{t+\Delta} \int_t^u \int_t^s \mathcal{L}^2\mu(v, X(v)) dv ds du \\ &= X(t) + \mu(t, X(t)) \Delta + \frac{1}{2} \mathcal{L}\mu(t, X(t)) \Delta^2 + O(\Delta^3), \end{aligned} \quad (3.47)$$

which is just (3.44). We can continue the iteration to arbitrary high order.

3.2.6.2 Ito-Taylor Expansions

One may wonder why in the previous section we bothered with the integral representation of Taylor's theorem when the usual (differential) Taylor expansion lead to the correct result. The reason is that the integral approach can be extended to SDEs, leading to stochastic *Ito-Taylor expansions*. To give a flavor of these, reintroduce a diffusion term to (3.43), and start out with the integral representation

$$X(t + \Delta) = X(t) + \int_t^{t+\Delta} \mu(u, X(u)) du + \int_t^{t+\Delta} \sigma(u, X(u)) dW(u). \quad (3.48)$$

Applying Ito's lemma to μ gives (compare to (3.46))

$$\mu(u, X(u)) = \mu(t, X(t)) + \int_t^u \mathcal{L}_0\mu(s, X(s)) ds + \int_t^u \mathcal{L}_1\mu(s, X(s)) dW(s), \quad (3.49)$$

where

$$\mathcal{L}_0 \triangleq \frac{\partial}{\partial t} + \mu \frac{\partial}{\partial x} + \frac{1}{2} \sigma^2 \frac{\partial^2}{\partial x^2}, \quad \mathcal{L}_1 \triangleq \sigma \frac{\partial}{\partial x}.$$

Similarly,

$$\sigma(u, X(u)) = \sigma(t, X(t)) + \int_t^u \mathcal{L}_0\sigma(s, X(s)) ds + \int_t^u \mathcal{L}_1\sigma(s, X(s)) dW(s). \quad (3.50)$$

Plugging (3.49) and (3.50) into (3.48) yields

$$\begin{aligned} X(t + \Delta) &= X(t) + \mu(t, X(t)) \int_t^{t+\Delta} du + \sigma(t, X(t)) \int_t^{t+\Delta} dW(u) + R_1 \\ &= X(t) + \mu(t, X(t)) \Delta + \sigma(t, X(t)) (W(t + \Delta) - W(t)) + R_1, \end{aligned} \quad (3.51)$$

where the remainder R_1 is

$$\begin{aligned} R_1 = & \int_t^{t+\Delta} \int_t^u \mathcal{L}_0 \mu(s, X(s)) ds du \\ & + \int_t^{t+\Delta} \int_t^u \mathcal{L}_1 \mu(s, X(s)) dW(s) du \\ & + \int_t^{t+\Delta} \int_t^u \mathcal{L}_0 \sigma(s, X(s)) ds dW(u) \\ & + \int_t^{t+\Delta} \int_t^u \mathcal{L}_1 \sigma(s, X(s)) dW(s) dW(u). \end{aligned}$$

As for the ODE example above, we can repeat this procedure arbitrarily many times. Going just one step further, we arrive at

$$\begin{aligned} X(t + \Delta) = & X(t) + \mu(t, X(t)) \Delta + \sigma(t, X(t)) (W(t + \Delta) - W(t)) \\ & + \mathcal{L}_0 \mu(t, X(t)) \frac{1}{2} \Delta^2 \\ & + \mathcal{L}_1 \mu(t, X(t)) \int_t^{t+\Delta} \int_t^u dW(s) du \\ & + \mathcal{L}_0 \sigma(t, X(t)) \int_t^{t+\Delta} \int_t^u ds dW(u) \\ & + \mathcal{L}_1 \sigma(t, X(t)) \int_t^{t+\Delta} \int_t^u dW(s) dW(u) + R_2, \end{aligned} \quad (3.52)$$

where R_2 contains triple integrals over t and W .

Stochastic Taylor expansions can be continued to arbitrary order, but we shall not go any further.

3.2.6.3 Milstein Second-Order Discretization Scheme

Discarding the remainder R_1 in the one-step iteration (3.51) is seen to lead to the Euler scheme (see Section 3.2.3), known to have weak convergence order $\beta = 1$. Under additional regularity (see Talay [1984]) of μ and σ , discarding the remainder R_2 in the higher-order expansion (3.52) can form the basis of a discretization scheme with weak order $\beta = 2$. For us to implement such a scheme, however, we need to concern ourselves with the simulation of the three stochastic double integrals figuring in (3.52). We go through the integrals in order below.

First,

$$\begin{aligned} I_{(1,1)} &\triangleq \int_t^{t+\Delta} \int_t^u dW(s) dW(u) \\ &= \int_t^{t+\Delta} (W(u) - W(t)) dW(u) = \frac{1}{2} (W(t + \Delta) - W(t))^2 - \frac{1}{2} \Delta, \end{aligned} \quad (3.53)$$

where we have used the fact that

$$\int_0^t W(u) dW(u) = \frac{1}{2} W(t)^2 - \frac{1}{2} t,$$

as can be verified by Ito's lemma. Second,

$$\begin{aligned} I_{(0,1)} &\triangleq \int_t^{t+\Delta} \int_t^u ds dW(u) = \int_t^{t+\Delta} (u - t) dW(u) \\ &= \Delta (W(t + \Delta) - W(t)) - \int_t^{t+\Delta} (W(u) - W(t)) du \\ &\triangleq \Delta (W(t + \Delta) - W(t)) - I_{(1,0)}, \end{aligned} \quad (3.54)$$

where we have used the integration-by-parts formula

$$\int_0^t u dW(u) = tW(t) - \int_0^t W(u) du,$$

which follows from applying Ito's lemma to $tW(t)$. In (3.54), the remaining integral $I_{(1,0)}$ on the right-hand-side is the same as the final double integral in (3.52), namely

$$I_{(1,0)} \triangleq \int_t^{t+\Delta} \int_t^u dW(s) du = \int_t^{t+\Delta} (W(u) - W(t)) du.$$

Reversing the order of integration, we get

$$\begin{aligned} I_{(1,0)} &= \int_t^{t+\Delta} \int_t^u dW(s) du \\ &= \int_t^{t+\Delta} \int_u^{t+\Delta} ds dW(u) = \int_t^{t+\Delta} (t + \Delta - u) dW(u) \end{aligned}$$

so we see, from Theorem 1.1.3 and the discussion in Section 1.6 on linear SDEs, that $I_{(1,0)}$ is Gaussian with mean 0 and variance

$$\text{Var}(I_{(1,0)}) = \int_t^{t+\Delta} (t + \Delta - u)^2 du = \frac{1}{3} \Delta^3.$$

The covariance between $I_{(1,0)}$ and $W(t + \Delta) - W(t)$ can be computed as

$$\text{Cov}(I_{(1,0)}, W(t + \Delta) - W(t)) = \int_t^{t+\Delta} (t + \Delta - u) du = \frac{1}{2} \Delta^2.$$

With the results above, we can cast the Taylor expansion (3.52) in the form of a simulation scheme (μ , σ , and their derivatives are to be evaluated at $t = i\Delta$ and $X = \widehat{X}(i\Delta)$),

$$\begin{aligned}\widehat{X}_{i+1} = & \widehat{X}_i + \mu\Delta + \sigma Z_{i,1}\sqrt{\Delta} + \mathcal{L}_0\mu\frac{1}{2}\Delta^2 + \mathcal{L}_1\mu Z_{i,2}\sqrt{\frac{1}{3}\Delta^3} \\ & + \mathcal{L}_0\sigma\left[\Delta Z_{i,1}\sqrt{\Delta} - Z_{i,2}\sqrt{\frac{1}{3}\Delta^3}\right] + \mathcal{L}_1\sigma\left(\frac{1}{2}Z_{i,1}^2\Delta - \frac{1}{2}\Delta\right),\end{aligned}\quad (3.55)$$

where $Z_{i,1}$ and $Z_{i,2}$ are sequences of $\mathcal{N}(0, 1)$ Gaussian variables with pairwise correlation

$$\rho(Z_{i,1}, Z_{i,2}) = \frac{\frac{1}{2}\Delta^2}{\sqrt{\frac{1}{3}\Delta^3}\sqrt{\Delta}} = \sqrt{\frac{3}{4}}.$$

The scheme above is known as the *Milstein scheme*. As mentioned earlier, the scheme has weak convergence order 2 under fairly strong regularity assumptions on μ and σ . We note that in the literature on SDE simulation, the Milstein scheme is often presented in a simplified form with the integral $I_{(1,0)}$ simulated as

$$\mathbb{E}\left(\int_t^{t+\Delta} \int_t^u dW(s) du \middle| W(t), W(t+\Delta)\right) = \frac{1}{2}\Delta(W(t+\Delta) - W(t)),$$

which corresponds to replacing $Z_{i,2}\sqrt{\Delta^3/3}$ with $\frac{1}{2}\Delta Z_{i,1}$ in (3.55). See Kloeden and Platen [2000] for a discussion of why this type of simplification does not affect the weak convergence order. The same source also contains a full discussion of how to extend the Milstein scheme to multi-dimensional SDEs.

3.2.7 Other Second-Order Schemes

The need to explicitly compute derivatives of the functions μ and σ often makes the Milstein scheme inconvenient to apply. High-order simulation schemes that substitute finite difference approximations for derivatives exist, and retain second (or higher) order weak convergence, are surveyed in Kloeden and Platen [2000]. To give an example of such a scheme, consider the scalar case $d = p = 1$ and assume that SDE coefficient functions μ and σ are function of x only. A derivative-free scheme that achieves second-order weak convergence is (from Kloeden and Platen [2000], Chapter 15)

$$\begin{aligned}\widehat{X}_{i+1} = & \widehat{X}_i + \frac{1}{2}\left(\mu(\overline{X}) + \mu(\widehat{X}_i)\right)\Delta \\ & + \frac{1}{4}\left(\sigma(\overline{X}^+) + \sigma(\overline{X}^-) + 2\sigma(\widehat{X}_i)\right)Z_i\sqrt{\Delta} \\ & + \frac{1}{4}\left(\sigma(\overline{X}^+) - \sigma(\overline{X}^-)\right)(Z_i^2\Delta - \Delta)\Delta^{-1/2},\end{aligned}\quad (3.56)$$

where the Z_i 's are a sequence of $\mathcal{N}(0, 1)$ Gaussian variables, and

$$\begin{aligned}\bar{X} &= \hat{X}_i + \mu(\hat{X}_i) \Delta + \sigma(\hat{X}_i) Z_i \sqrt{\Delta}, \\ \bar{X}^\pm &= \hat{X}_i + \mu(\hat{X}_i) \Delta \pm \sigma(\hat{X}_i) \sqrt{\Delta}.\end{aligned}$$

Comparison of (3.56) with the simplified Milstein scheme in the previous section shows that (3.56) avoids derivatives by using additional supporting values \bar{X} and \bar{X}^\pm .

Another, quite different, approach to avoid explicit derivatives applies the classical idea of *Richardson extrapolation* to the Euler scheme. This idea was proposed by Talay and Tubaro [1990] and takes advantage of the fact that, under additional regularity conditions, the error of the Euler scheme can be sharpened beyond (3.30) (with $\beta = 1$) to

$$E\left(g\left(\hat{X}(T)\right)\right) = E(g(X(T))) + c\Delta + O(\Delta^2), \quad (3.57)$$

for a constant c . Let \hat{X}_Δ and $\hat{X}_{2\Delta}$ be estimates of X based on Euler discretizations with time steps of Δ and 2Δ , respectively. Provided that (3.57) holds, we can write

$$2E\left(g\left(\hat{X}_\Delta(T)\right)\right) - E\left(g\left(\hat{X}_{2\Delta}(T)\right)\right) = E(g(X(T))) + O(\Delta^2), \quad (3.58)$$

which is our second-order extrapolation formula. As the Euler scheme is simple to set up, the extrapolation scheme is an attractive alternative to other second-order techniques. In practice, however, the convergence of the Euler scheme may not always be smooth enough to make (3.58) work well. Numerical experiments will nearly always be necessary (as is also the case of the Ito-Taylor schemes, for that matter).

A final word about generation of \hat{X}_Δ and $\hat{X}_{2\Delta}$ in the Richardson extrapolation scheme. To avoid duplication of work, we discretize time in buckets of Δ and generate both \hat{X}_Δ and $\hat{X}_{2\Delta}$ simultaneously, combining time steps in pairs for the purpose of generating $\hat{X}_{2\Delta}$. That is, if we use Gaussian increments of $Z_1\sqrt{\Delta}, Z_2\sqrt{\Delta}, \dots$ for \hat{X}_Δ , we use $(Z_1 + Z_2)\sqrt{\Delta}, (Z_3 + Z_4)\sqrt{\Delta}, \dots$ for $\hat{X}_{2\Delta}$. Not only do we save work by re-using Gaussian draws, we most likely also reduce the statistical error of our Monte Carlo estimate of the difference $2E(g(\hat{X}_\Delta(T))) - E(g(\hat{X}_{2\Delta}(T)))$ by raising correlation between $g(\hat{X}_\Delta(T))$ and $g(\hat{X}_{2\Delta}(T))$. We shall return to this idea in Section 3.3.1.

3.2.8 Bias vs. Monte Carlo Error

When we use an m -step discretization scheme in an n -path Monte Carlo run, we are exposed to two types of errors on the expectation we are trying to evaluate: i) the statistical Monte Carlo error e_s (the standard error); and ii) a bias e_b , originating from the discretization scheme. Raising n will reduce the standard error, but will not affect the bias which can only be

reduced by increasing the number of steps m in the time discretization scheme. Raising m and/or n obviously involves a computational cost, so let us briefly consider explicitly the trade-offs involved in simultaneously reducing bias and standard error.

Assume first that the discretization scheme has weak order β . Proceeding informally, we interpret this as implying

$$e_b = c_b \Delta^\beta,$$

for some constant c_b . Also, we know that the variance of e_s is

$$\text{Var}(e_s) = \frac{c_s}{n},$$

for a constant c_s . The total computing time τ is reasonably assumed to be proportional to nm or, using the fact that $\Delta = T/m$,

$$\tau = n \frac{c_\tau}{\Delta} \quad (3.59)$$

for some constant c_τ . For a given computing budget τ , consider minimizing the total mean-square error (MSE) $c_b^2 \Delta^{2\beta} + \frac{c_s}{n}$. Using (3.59) to eliminate a variable, the optimization problem is

$$\min_{\Delta} \left(c_b^2 \Delta^{2\beta} + \frac{c_s c_\tau}{\tau \Delta} \right).$$

Let Δ^* be the value of Δ at which the minimum MSE is attained. Δ^* is seen to satisfy

$$\Delta^* = C \tau^{-\frac{1}{2\beta+1}}, \quad C \triangleq \left(\frac{c_s c_\tau}{2\beta c_b^2} \right)^{\frac{1}{2\beta+1}}, \quad (3.60)$$

such that the minimum MSE becomes

$$C' \tau^{-\frac{2\beta}{2\beta+1}} \quad (3.61)$$

for yet another constant C' .

Equations (3.60) and (3.61) reveal a number of structural characteristics of Monte Carlo simulation of discretized SDEs. For instance, according to (3.61), the optimal root-mean-square (RMS) error behaves with the computing time τ as

$$\text{RMS} \propto \tau^{-\frac{\beta}{2\beta+1}}. \quad (3.62)$$

The computational cost of working with SDEs that are not explicitly solvable are quantified by (3.62). For an unbiased (that is, exact) SDE simulation scheme, $\beta = \infty$ and the optimal RMS error converges at the rate of $\tau^{-\frac{1}{2}}$, consistent with the results of Section 3.1. However, for an Euler scheme ($\beta = 1$) the RMS error convergence rate is lowered to $\tau^{-\frac{1}{3}}$.

Equation (3.60) in principle tells us how to optimally allocate resources between the competing objectives of a lower bias and a lower standard

error. Let m^* be defined through $\Delta^* = T/m^*$ and let n^* be defined through $\tau = c_\tau n^* m^*/T$. After a few rearrangements, we find the intuitive result

$$\sqrt{n^*} = C''(m^*)^\beta,$$

where C'' is a constant independent of τ . When we increase or decrease our computing budget, it is thus reasonable to allocate resources in such a way that we keep the factor $n^{1/2}m^{-\beta}$ constant. More detailed discussion, as well as asymptotic limit results, can be found in Duffie and Glynn [1995].

3.2.9 Sampling of Continuous Process Extremes

We round off our discussion of path simulation schemes by considering the pricing of options that depend on continuously or high-frequency sampled extremes of an SDE. We focus on the scalar case, with our SDE given as

$$dX(t) = \mu(t, X(t)) dt + \sigma(t, X(t)) dW(t),$$

where both X and W are 1-dimensional (i.e., $p = d = 1$). We also assume that the SDE is Euler-discretized according to

$$\widehat{X}_{i+1} = \widehat{X}_i + \mu(i\Delta, \widehat{X}_i) \Delta + \sigma(i\Delta, \widehat{X}_i) \sqrt{\Delta} Z_i, \quad i = 0, 1, \dots, m-1,$$

with $m\Delta = T$.

On the interval $[0, T]$, let the maximum and minimum values of $X(t)$ be denoted $M_{[0,T]}$ and $m_{[0,T]}$, respectively. That is,

$$M_{[0,T]} \triangleq \max_{0 \leq t \leq T} X(t); \quad m_{[0,T]} \triangleq \min_{0 \leq t \leq T} X(t).$$

To give examples of options that depend on $M_{[0,T]}$ and $m_{[0,T]}$, consider for instance the up-and-out call option we encountered in Section 2.7. With a knock-out barrier of H and a terminal strike of K , the terminal maturity payout can be written as¹¹

$$g(T) = 1_{\{M_{[0,T]} < H\}} (X(T) - K)^+.$$

A *double-barrier knock-out call option* with an upper barrier of H and a lower barrier of h pays

$$g(T) = 1_{\{m_{[0,T]} > h\}} 1_{\{M_{[0,T]} < H\}} (X(T) - K)^+.$$

Finally, a so-called *lookback call option* (see Section 2.7) pays

$$g(T) = (M_{[0,T]} - K)^+.$$

¹¹If $X(t)$ is the logarithm of the asset price, we replace this expression by $g(T) = 1_{\{M_{[0,T]} \leq e^H\}} (e^{X(T)} - K)^+$, and similarly for the other payouts considered.

To price options such as those above, we must provide pathwise estimates of $M_{[0,T]}$ and $m_{[0,T]}$. Given our discretization schemes, natural estimators are

$$\widehat{M}_{[0,T]} = \max(X(0), \widehat{X}_1, \dots, \widehat{X}_m), \quad (3.63)$$

$$\widehat{m}_{[0,T]} = \min(X(0), \widehat{X}_1, \dots, \widehat{X}_m). \quad (3.64)$$

Even in cases where the discretization scheme itself is perfectly unbiased, it is clear that these estimators will underestimate the range of the extremes of $X(t)$, by consistently failing to account for the movement (the “overshoot” and “undershoot”) of X between sample points $i\Delta$, $i = 0, 1, \dots, m$. As a consequence, for each simulated path in an otherwise unbiased discretization scheme, almost surely

$$\widehat{M}_{[0,T]} < M_{[0,T]}, \quad \widehat{m}_{[0,T]} > m_{[0,T]}.$$

As shown in Andersen and Brotherton-Ratcliffe [1996], the bias introduced can be very significant, even if Δ is quite small. For instance, for a 1 year lookback option, Andersen and Brotherton-Ratcliffe [1996] report that even daily sampling produces a 6% price error.

To improve the price estimates of options that depend on continuously sampled extremes, we should alter (3.63) and (3.64) to take into consideration movements between sample dates. This can be accomplished by the Brownian bridge technique introduced in Andersen and Brotherton-Ratcliffe [1996] (see also Broadie et al. [1997]). Let us focus on a particular bucket $[i\Delta, (i+1)\Delta]$ and assume, consistent with the Euler scheme, that $\widehat{X}(t)$, $t \in [i\Delta, (i+1)\Delta]$, is a Gaussian process with conditional moments

$$\begin{aligned} \mathbb{E}(\widehat{X}(t) - \widehat{X}_i \mid \widehat{X}_i) &= \mu(i\Delta, \widehat{X}_i)(t - i\Delta), \quad t \in [i\Delta, (i+1)\Delta]; \\ \text{Var}(\widehat{X}(t) - \widehat{X}_i \mid \widehat{X}_i) &= \sigma(i\Delta, \widehat{X}_i)^2(t - i\Delta), \quad t \in [i\Delta, (i+1)\Delta]. \end{aligned}$$

Assume that we have already simulated \widehat{X}_i and $\widehat{X}((i+1)\Delta)$ by the Euler scheme above. Conditional on *both* \widehat{X}_i and \widehat{X}_{i+1} , the process for $\widehat{X}(t)$, $t \in [i\Delta, (i+1)\Delta]$ is a Gaussian process “pinned” at the levels \widehat{X}_i and \widehat{X}_{i+1} . The resulting process is known as a *Brownian bridge* with diffusion coefficient $\sigma(i\Delta, \widehat{X}_i)$. Let \widehat{M}_i^c (\widehat{m}_i^c) be defined as the continuously sampled maximum (minimum) of $\widehat{X}(t)$ on $[i\Delta, (i+1)\Delta]$. The following lemma is a special case of a result in Andersen and Brotherton-Ratcliffe [1996]:

Lemma 3.2.1.

$$\begin{aligned} \mathbb{P}(\widehat{M}_i^c \leq s \mid \widehat{X}_i, \widehat{X}_{i+1}) &= 1 - \xi_i(s), \quad s > \max(\widehat{X}_i, \widehat{X}_{i+1}), \\ \mathbb{P}(\widehat{m}_i^c \leq s \mid \widehat{X}_i, \widehat{X}_{i+1}) &= \xi_i(s), \quad s < \min(\widehat{X}_i, \widehat{X}_{i+1}), \end{aligned}$$

where

$$\xi_i(s) \triangleq \exp\left(\frac{2(s - \hat{X}_i)(\hat{X}_{i+1} - s)}{\sigma(i\Delta, \hat{X}_i)^2 \Delta}\right).$$

We can use the result of the lemma in a number of ways. Most obviously, we can apply it to sample \widehat{M}_i^c and \widehat{m}_i^c directly, by the inverse transform method (see Section 3.1.1.1). To illustrate, consider for instance sampling \widehat{M}_i^c . Having first drawn \hat{X}_i and \hat{X}_{i+1} by usual means, we draw an additional independent $\mathcal{U}(0, 1)$ uniform variable U_i , and set

$$1 - \xi_i(\widehat{M}_i^c) = U_i$$

or, after a few rearrangements,

$$\widehat{M}_i^c = \frac{1}{2} (\hat{X}_{i+1} + \hat{X}_i) + \frac{1}{2} \sqrt{(\hat{X}_{i+1} - \hat{X}_i)^2 - 2\sigma(i\Delta, \hat{X}_i)^2 \Delta \ln(1 - U_i)}.$$

This procedure can be repeated for $i = 0, 1, \dots, m-1$, giving us the improved estimator for the maximum of X over $[0, T]$,

$$\widehat{M}_{[0, T]}^c = \max(\widehat{M}_0^c, \dots, \widehat{M}_{m-1}^c).$$

For options depending on both the minimum and maximum (such as double barrier options and the double lookback options in He et al. [1998]), the necessary extensions required for joint sampling of minimum and maximum are developed in Andersen [1998].

For barrier options, we note that locating \widehat{M}_i^c and \widehat{m}_i^c directly is typically not necessary, as it suffices to check locally whether the barrier is breached. For an up-and-out knock-out option with barrier H , for each interval it thus suffices to check whether $\widehat{M}_i^c > H$ which, conditional on \hat{X}_i and \hat{X}_{i+1} , happens with likelihood $\xi_i(H)$. So, provided that \hat{X}_i and \hat{X}_{i+1} are both below H , determining whether the barrier was nevertheless breached in $[i\Delta, (i+1)\Delta]$ is a matter of drawing a uniform variable U_i and setting

$$1_{\{\widehat{M}_i^c \geq H\}} = 1_{\{U_i \leq \xi_i(H)\}}. \quad (3.65)$$

This scheme is easily extended to time-dependent barriers and to cases where there are rebates¹². As pointed out in Glasserman and Staum [2001], for

¹²To get the timing of rebate payments right, the exact time that the barrier is breached must, in principle, be located. Andersen and Brotherton-Ratcliffe [1996] list analytical Brownian bridge hitting time results that can be used for this purpose. For reasonably fine discretizations, it will often suffice to set the hitting time to, say, the mid-point of the time bucket where the barrier is known to be breached.

Markov processes and the special case of barrier options with no rebates, one can in fact avoid drawing U_i 's altogether, as

$$\begin{aligned} \mathbb{E} \left(g(\widehat{X}_m) \mathbf{1}_{\{\widehat{M}_{[0,T]}^c < H\}} \middle| \widehat{X}_0, \widehat{X}_1, \dots, \widehat{X}_m \right) \\ = \mathbb{E} \left(g(\widehat{X}_m) \prod_{i=0}^{m-1} \mathbf{1}_{\{\widehat{M}_i^c < H\}} \middle| \widehat{X}_0, \widehat{X}_1, \dots, \widehat{X}_m \right) \\ = g(\widehat{X}_m) \prod_{i=0}^{m-1} \mathbb{E} \left(\mathbf{1}_{\{\widehat{M}_i^c < H\}} \middle| \widehat{X}_i, \widehat{X}_{i+1} \right). \end{aligned}$$

Here,

$$\mathbb{E} \left(\mathbf{1}_{\{\widehat{M}_i^c < H\}} \middle| \widehat{X}_i, \widehat{X}_{i+1} \right) = \begin{cases} 0, & \widehat{X}_i \geq H \text{ or } \widehat{X}_{i+1} \geq H, \\ 1 - \xi_i(H), & \widehat{X}_i < H \text{ and } \widehat{X}_{i+1} < H. \end{cases}$$

In other words, rather than explicitly simulating the indicator functions (3.65), it suffices to adjust the terminal payout by the product of conditional survival probabilities along the path. This scheme is an example of *conditional Monte Carlo*, a variance-reduction scheme discussed in more detail in Section 25.2 and in Boyle et al. [1997]. One potential drawback of the scheme is the fact that we typically need to continue the paths for a longer period of time before a barrier crossing is detected and the path can be stopped.

We round off this section with a few comments. First, we note that the schemes above assume that X is well approximated by a Gaussian process. In some applications the geometric Brownian motion, say, may be a more appropriate model. In Lemma 3.2.1, we can easily accommodate this by simply replacing \widehat{M}_i^c , \widehat{m}_i^c , s , \widehat{X}_i , and \widehat{X}_{i+1} with $\ln \widehat{M}_i^c$, $\ln \widehat{m}_i^c$, $\ln s$, $\ln \widehat{X}_i$, and $\ln \widehat{X}_{i+1}$. Other transformations are handled the same way.

Secondly, it should be pointed out that many real options are, in fact, not sampled continuously but rather at some finite but high frequency, often daily. Running an Euler scheme with daily discretization is obviously computationally inefficient. Fortunately, we can often use our scheme above as part of a Richardson-type interpolation idea. Indeed, it can often be established (see Andersen and Brotherton-Ratcliffe [1996]) that options on process extremes converge as $O(\sqrt{\Delta})$. If we first compute a price estimate \widehat{V}^Δ based on a relatively coarse value of Δ , we can then write

$$\widehat{V}^{\Delta^*} \approx \widehat{V}^c + (\widehat{V}^\Delta - \widehat{V}^c)\Delta^*/\Delta, \quad \Delta^* < \Delta,$$

where \widehat{V}^c is the continuously monitored price computed by the scheme outlined above. We have here implicitly made the assumption that the regular Euler bias is small relative to the bias induced by using the wrong sampling frequency. The idea above is developed further in Chapter V of Andersen [1996], where a number of numerical results can also be found.

And finally, one may wonder whether it is possible to deal with a continuously monitored barrier option in a discrete-time simulation by adjusting the *barrier*, rather than the underlying *process*. As it turns out, this is indeed possible. Specifically, in the Black-Scholes-Merton model with volatility σ , let $V^c(H)$, $V^\Delta(H)$ be the values of a continuously and discretely sampled barrier options with barrier H , respectively. Assuming that the discrete sampling happens on a time grid with spacing Δ , we have the following result from Broadie et al. [1997]:

Theorem 3.2.2. *The following holds,*

$$V^\Delta(H) = V\left(He^{\pm\beta\sigma\sqrt{\Delta}}\right) + o\left(\sqrt{\Delta}\right),$$

where $+$ applies if $H > X(0)$, $-$ applies if $H < X(0)$, and $\beta = \zeta(1/2)/\sqrt{2\pi} \approx 0.5826$, with $\zeta(\cdot)$ being the Riemann zeta function.

According to this result we can price a continuous barrier option by evaluating a discrete barrier option instead (e.g., one where the barrier monitoring takes place only on the simulation dates of the Monte Carlo scheme used), but with the discrete barrier level shifted according to the theorem. Theorem 3.2.2 can also be used to save computation time by, say, turning a barrier option with daily observations into an option with quarterly observations, as the theorem shows that

$$V^{\Delta^*}(H) \approx V^\Delta\left(He^{\pm\beta\sigma(\sqrt{\Delta^*}-\sqrt{\Delta})}\right).$$

While the result of Theorem 3.2.2 is only proved for the log-normal process, practical experience shows that it is robust across a wide variety of models. A similar approach exists for lookback options, see Broadie et al. [1999].

3.2.10 PCA and Bridge Construction of Brownian Motion Paths

3.2.10.1 Brownian Bridge and Quasi-Random Sequences

To close out the section on sample path simulation, let us address alternative ways of generating sample paths of Brownian motion. So far, to produce a sample of the vector $\mathbf{W} = (W(\Delta), W(2\Delta), \dots, W(m\Delta))^\top$, we have relied exclusively on the forward recursion

$$W(i\Delta + \Delta) = W(i\Delta) + Z_i\sqrt{\Delta}, \quad W(0) = 0, \quad (3.66)$$

where Z_0, Z_1, \dots, Z_{m-1} is a sequence of independent standard Gaussian variables. Rather than filling out the elements of \mathbf{W} in order, we may, for instance, rely on a *Brownian bridge (BB) construction* where we first

sample the end-point $W(m\Delta)$, then sample the mid-point¹³ $W(\lfloor m/2 \rfloor \Delta)$ *conditional* on $W(m\Delta)$, and so forth. In executing this scheme, we can use the easily proven result below.

Lemma 3.2.3. *Let $\underline{t} < t < \bar{t}$. Conditional on $W(\underline{t})$ and $W(\bar{t})$, $W(t)$ is Gaussian with moments*

$$\begin{aligned}\mathbb{E}(W(t)|W(\underline{t}) = \underline{w}, W(\bar{t}) = \bar{w}) &= \underline{w} \times \frac{\bar{t} - t}{\bar{t} - \underline{t}} + \bar{w} \times \frac{t - \underline{t}}{\bar{t} - \underline{t}}, \\ \text{Var}(W(t)|W(\underline{t}) = \underline{w}, W(\bar{t}) = \bar{w}) &= \frac{(\bar{t} - t)(t - \underline{t})}{\bar{t} - \underline{t}}.\end{aligned}$$

The BB scheme for construction of \mathbf{W} relies on repeated application of the result in Lemma 3.2.3 to progressively fill in \mathbf{W} in the “bisection” manner described above; consult any Monte Carlo textbook (e.g. Jäckel [2002] or Glasserman [2004]) if further details are required. As is the case for the standard scheme (3.66), a total of m standard Gaussian random variables are needed to construct a single sample of \mathbf{W} by the Brownian bridge scheme, so the latter offers no computational advantage over the former. Why then use the Brownian bridge construction?

One important distinction between the BB construction and (3.66) is the fact that the Brownian bridge assigns different importance to the random numbers used to produce \mathbf{W} . For instance, the very first Gaussian number drawn in the BB technique *alone* determines the end-point $W(m\Delta)$ of the Brownian motion — and thereby establishes a significant part of the overall coarse structure of the path of W . Subsequent random number draws contribute by filling in the details of the W -path, with late draws adding only to the fine-structure of the path. In contrast, with (3.66) the end-point $W(m\Delta)$ is affected equally by the m random numbers Z_0, Z_1, \dots, Z_{m-1} . In most financial problems the coarse shape of the path of W is more critical than finer details, so ultimately the BB technique allows us to identify and isolate the important features of the Brownian motion path. In some variance reduction techniques this can be important, as it allows us to focus computational effort on the random numbers that matter the most. Also, some variance reduction techniques that are known to work particularly well on low-dimensional problems can now be applied to the (low-dimensional) random numbers that contribute most to the sample path.

One relevant technique is *quasi-random sequences* (also known as *low-discrepancy sequences*), a method of generating points on the hypercube that are as “dispersed” as possible. A good survey of the underlying ideas and theory can be found in Jäckel [2002] or Glasserman [2004], with source code available (for the special case of *Sobol sequences*) in Press et al. [1992]; suffice to say that quasi-random sequences can, under some circumstances,

¹³ $\lfloor x \rfloor$ denotes the integer part of a real variable x .

accelerate Monte Carlo convergence substantially¹⁴. It is well-known, however, that the efficacy of quasi-random sequences depend strongly on the problem dimension (here: m , the number of random numbers needed per path), and that the sequences deteriorate in higher dimensions. When quasi-random sequences are combined with BB simulation of the path, however, the (low-dimensional) points of the sequences that are well-distributed can be applied to generate — by the methods in Section 3.1.1 — the Gaussian samples that determine the coarse structure of the paths, whereas the poorly distributed (high-dimensional) parts of the sequence can be relegated to the generation of less important fine-structure details¹⁵. A full account of this idea can be found in Moskowitz and Caflisch [1996].

3.2.10.2 PC Construction

With the Brownian Bridge (BB) construction of Brownian motion, much of the variance of the sample paths \mathbf{W} is explained by the values of the first few (Gaussian) random variables drawn in the path simulation. We recall from Section 3.1.3, however, that the *optimal* way to project the variation of a Gaussian vector onto a low-dimensional set of random variables is done through a principal component (PC) construction, rather than the Brownian bridge. To demonstrate how a PC construction of $\mathbf{W} = (W(\Delta), W(2\Delta), \dots, W(m\Delta))^\top$ would proceed, we first notice that the $m \times m$ variance-covariance matrix Σ of \mathbf{W} has elements

$$\Sigma_{i,j} = \mathbb{E}(W(i\Delta) W(j\Delta)) = \sqrt{\frac{\min(i, j)}{\max(i, j)}}, \quad i, j = 1, 2, \dots, m.$$

As shown in Åkesson and Lehoczy [1998], the eigenvalues of Σ can be found analytically to be

$$\lambda_i = \frac{\Delta}{4} \sin\left(\frac{\pi}{2} \cdot \frac{2i-1}{2m+1}\right)^{-2}, \quad i = 1, \dots, m,$$

where $\lambda_1 > \lambda_2 > \dots > \lambda_m$. Let e_i be the eigenvector associated with λ_i , then it is also known that $e_i = (e_{i,1}, e_{i,2}, \dots, e_{i,m})^\top$, where

$$e_{i,j} = \frac{2}{\sqrt{2m+1}} \sin\left(j\pi \cdot \frac{2i-1}{2m+1}\right), \quad j = 1, \dots, m.$$

From the results in Section 3.1.3, we know that we can write

¹⁴Theoretically from $O(1/\sqrt{N})$ to (nearly) $O(1/N)$. Comparative tests on actual finance problems can be found in, for instance, Brotherton-Ratcliffe [1994], Paskov and Traub [1995], and Joy et al. [1996].

¹⁵Alternatively we can use regular *pseudo-random numbers* for this.

$$\mathbf{W} = \sum_{i=1}^m Z_i \sqrt{\lambda_i} e_i, \quad (3.67)$$

where Z_1, Z_2, \dots, Z_m is a sequence of independent standard Gaussian random variables. This equation constitutes the principal components construction of the Brownian path, and it is characterized by the fact that for any $k \leq m$, the first k terms of (3.67) (that is, $\sum_{i=1}^k Z_i \sqrt{\lambda_i} e_i$) explain as much of the variance of \mathbf{W} as is possible with k Gaussian variables. Even more so than for the Brownian bridge, the PC construction of a Brownian motion thus connects the overall shape of the Brownian path to a few of the Gaussian random variables Z_i , with the remaining random variables contributing only high-frequency details. As explained above, this can be useful in certain variance reduction techniques by allowing us to focus our attention and resources on just a few of the m random variables needed to simulate \mathbf{W} . We note that the PC construction is more expensive to compute than the BB technique (the latter is $O(m)$ whereas the former can be seen from (3.67) to be $O(m^2)$), so the optimality of the PC approach may, in some applications, be outweighed by its lack of speed.

While we developed the BB and PC constructions exclusively on an equidistant time grid, they easily extend to non-equidistant grids. When the grid is non-equidistant, the variance-covariance matrix of \mathbf{W} has elements

$$\Sigma_{i,j} = \mathbb{E}(W(t_i)W(t_j)) = \sqrt{\frac{\min(t_i, t_j)}{\max(t_i, t_j)}}, \quad i, j = 1, 2, \dots, m,$$

and eigenvectors and eigenvalues must then be found numerically, rather than through the analytical results listed earlier. Also, both the BB and PC techniques can easily be extended to the case of multi-variate Brownian motions, see Jäckel [2002].

Finally, for those interested in such matters, we note that in the limit $m \rightarrow \infty$, the PC construction of Brownian motion is known as the *Karhunen-Loeve decomposition*. In the continuous-time limit, the BB representation is sometimes known as a *Haar function* decomposition of Brownian motion.

3.3 Sensitivity Computations

In most finance applications, the fact that options must be dynamically hedged and risk managed requires us not only to produce an estimate of an option price, but also to compute reliable estimates of the sensitivity of the price with respect to the underlying state variables, as well as various other model parameters. In this section, we will present a number of methods for sensitivity computations by Monte Carlo methods. For each method, we use the problem of estimating the stock price delta of options in the Black-Scholes economy as a motivating example. We shall spend much more

time on sensitivity computations (by PDE and Monte Carlo methods) in Part V of this book, often in the context of particular interest rate products. Here we just give a flavor of things to come.

3.3.1 Finite Difference Estimates

3.3.1.1 Black-Scholes Delta

Consider a T -maturity European option on a dividend-free stock S in the Black-Scholes economy. Let the payout function be $g(S(T))$, and assume that the continuously compounded interest rate is a constant r . With $V(S_0)$ denoting the time 0 price of the option given $S(0) = S_0$, we are interested in computing

$$\frac{dV}{dS_0} = \lim_{h \rightarrow 0} \frac{V(S_0 + h) - V(S_0)}{h}. \quad (3.68)$$

In a Monte Carlo setting, we can approximate this derivative (“delta”) by finite difference techniques as follows. First, for some fixed number ε draw random standard Gaussian variables Z and Z_ε , and set (see Section 3.2.1)

$$S(T) = S_0 \exp \left(\left(r - \frac{1}{2}\sigma^2 \right) T + \sigma\sqrt{T}Z \right),$$

$$S_\varepsilon(T) = (S_0 + \varepsilon) \exp \left(\left(r - \frac{1}{2}\sigma^2 \right) T + \sigma\sqrt{T}Z_\varepsilon \right),$$

where σ as always denotes the constant volatility of the stock. We then form the difference

$$\delta = e^{-rT}\varepsilon^{-1} (g(S_\varepsilon(T)) - g(S(T))),$$

such that δ constitutes a single-sample estimate for dV/dS_0 . By generating n independent replications of δ and forming the sample average, we will obtain in the limit $n \rightarrow \infty$ the finite difference ratio

$$\frac{V(S_0 + \varepsilon) - V(S_0)}{\varepsilon}. \quad (3.69)$$

We know from Chapter 2 that this estimate will be biased relative to the true derivative dV/dS_0 by an amount of order $O(\varepsilon^2)$.

We have so far not mentioned whether the standard Gaussian variables Z and Z_ε should be independent or not. To analyze this, we need to consider the variance of the Monte Carlo estimator of (3.69). From Theorem 3.1.2, we know that for a finite number of trials n , the variance of our sample average will decrease as v_ε/n , where

$$\begin{aligned} v_\varepsilon &= \varepsilon^{-2} e^{-2rT} \text{Var}(g(S_\varepsilon(T)) - g(S(T))) \\ &= \varepsilon^{-2} e^{-2rT} [\text{Var}(g(S(T))) + \text{Var}(g(S_\varepsilon(T))) \\ &\quad - 2\text{Cov}(g(S_\varepsilon(T)), g(S(T)))]. \end{aligned}$$

If the random numbers Z and Z_ε are independent, $\text{Cov}(g(S_\varepsilon(T)), g(S(T)))$ will be zero and

$$v_\varepsilon \approx 2\varepsilon^{-2} e^{-2rT} \text{Var}(g(S(T))).$$

Making ε approach zero — as is needed to reduce the bias of the finite difference approximation (3.69) — will cause v_ε grow at a rate of $O(\varepsilon^{-2})$. This is obviously not ideal as our Monte Carlo estimate will be swamped by noise if ε is picked too small. On the other hand, if we set Z and Z_ε to be *identical*, we see that

$$S_\varepsilon(T) = (S_0 + \varepsilon) S(T)/S_0$$

and would expect

$$\text{Cov}(g(S_\varepsilon(T)), g(S(T))) > 0,$$

which would reduce v_ε relative to the independent case. For smooth g , a Taylor expansion in ε shows that

$$\begin{aligned} g(S_\varepsilon(T)) &= g(S(T)) + (S_\varepsilon(T) - S(T)) g'(S(T)) + \dots \\ &= g(S(T)) + \varepsilon S(T)/S_0 \cdot g'(S(T)) + \dots \end{aligned}$$

If derivatives of g are bounded

$$\text{Cov}(g(S_\varepsilon(T)), g(S(T))) = \text{Var}(g(S(T))) + O(\varepsilon^2),$$

and similarly for $\text{Var}(g(S_\varepsilon(T)))$. In other words,

$$v_\varepsilon = e^{-2rT} O(1) \tag{3.70}$$

which is a clear improvement over the earlier $O(\varepsilon^{-2})$ result.

The result (3.70) hinged on the payout function having bounded derivatives. We can, in fact, relax this considerably, to functions that are essentially just continuous in the stock price; see Section 3.3.2.2 for a discussion. For discontinuous payouts, however, (3.70) will not hold. To demonstrate, consider a digital option paying

$$g(S(T)) = 1_{\{S(T) > K\}},$$

for some strike K . With $Z = Z_\varepsilon$ we get (assuming $\varepsilon > 0$ and that the probability measure is P)

$$\begin{aligned} \mathbb{E}\left([g(S_\varepsilon(T)) - g(S(T))]^2\right) &= P(S(T) \leq K < S_\varepsilon(T)) \\ &= P(S(T) \leq K < (1 + \varepsilon/S_0) S(T)) = O(\varepsilon), \end{aligned}$$

compared with the $O(\varepsilon^2)$ result for smooth g .

3.3.1.2 General Case

To generalize the problem considered in the previous section, we consider a setting where a random variable Y depends on a parameter $\alpha \in \mathbb{R}$, in the sense that each value of α uniquely determines a scheme for the generation of Y . The random variable Y will typically represent a (discounted) option payout, and α is typically an initial value of an asset price (as in Section 3.3.1.1) or a parameter in the (vector) equations determining the dynamics of the underlying model. Let

$$V(\alpha) = \mathbb{E}(Y(\alpha)),$$

and consider the problem of determining $dV/d\alpha$.

In the basic finite difference Monte Carlo approximation to $dV/d\alpha$, we use the sample average of one-sided difference coefficients,

$$\bar{\delta}_n = \frac{\bar{Y}_n(\alpha + \varepsilon) - \bar{Y}_n(\alpha)}{\varepsilon},$$

where $\bar{Y}_n(\alpha)$ is the sample average of n realizations of $Y(\alpha)$. In the limit,

$$\lim_{n \rightarrow \infty} \bar{\delta}_n = \frac{V(\alpha + \varepsilon) - V(\alpha)}{\varepsilon} = dV/d\alpha + O(\varepsilon^2).$$

If we instead wish to use a central estimator

$$\bar{\delta}_n^c = \frac{\bar{Y}_n(\alpha + \varepsilon) - \bar{Y}_n(\alpha - \varepsilon)}{2\varepsilon},$$

we get

$$\lim_{n \rightarrow \infty} \bar{\delta}_n^c = \frac{V(\alpha + \varepsilon) - V(\alpha - \varepsilon)}{2\varepsilon} = dV/d\alpha + O(\varepsilon^3),$$

but now need to simulate an extra random variable (that is, $Y(\alpha - \varepsilon)$), increasing the computational cost.

In the generation of $\bar{Y}_n(\alpha \pm \varepsilon)$ and $\bar{Y}_n(\alpha)$, the individual samples of $Y(\alpha + \varepsilon)$, $Y(\alpha - \varepsilon)$, and $Y(\alpha)$ would typically be based on a series of draws of vector-valued support variables Z , with $Y(\alpha) = Y(Z; \alpha)$, and so forth. For instance, in an m -step Euler simulation of an SDE with d Brownian motions, each SDE path (and each outcome of Y) would involve $d \cdot m$ i.i.d. standard Gaussian variables $Z_1, \dots, Z_{d \cdot m}$. The observations in the previous section tell us that to minimize variance we should use the same Z for $Y(\alpha + \varepsilon)$, $Y(\alpha - \varepsilon)$, and $Y(\alpha)$. In practice, this is often easiest to accomplish by simply using the same random number seed (see Section 3.1.1) in otherwise separate computations of each of the quantities $\bar{Y}_n(\alpha + \varepsilon)$, $\bar{Y}_n(\alpha - \varepsilon)$, and $\bar{Y}_n(\alpha)$. Assuming this so-called *common random number scheme* is followed, the variance analysis in Section 3.3.1.1 can be generalized to our setting, and we would expect that either i) $\text{Var}(\bar{\delta}_n) = \text{Var}(\bar{\delta}_n^c) = O(\varepsilon^{-1}n^{-1})$; or ii)

$\text{Var}(\bar{\delta}_n) = \text{Var}(\bar{\delta}_n^c) = O(n^{-1})$. Case ii) essentially requires a.s. continuity¹⁶ of $Y(\alpha)$ with respect to α , as would be the case when Y represents a continuous option payout function. Case i) generally applies when Y represents a discontinuous option payout, such as the digital option considered in Section 3.3.1.1.

If case ii) above applies, the estimator variance is independent of ε , and ε should be picked as small as possible (a matter of machine precision) to minimize the $O(\varepsilon^2)$ and $O(\varepsilon^3)$ biases of $\bar{\delta}_n$ and $\bar{\delta}_n^c$. If the overhead of evaluating $Y(Z; \alpha)$ for given Z is small relative to the cost of generating Z , the central estimator $\bar{\delta}_n^c$ will dominate. For complicated payout functions, however, there may be situations when $\bar{\delta}_n$ is preferable, despite its slower convergence rate in ε . For case i), we must weigh bias against variance in a manner quite similar to the discussion in Section 3.2.8: if ε is small the difference coefficient bias will be small, but the variance of the estimators $\bar{\delta}_n$ and $\bar{\delta}_n^c$ will be high. An RMS minimization similar to the one in Section 3.2.8 is possible; in the interest of brevity, we leave this as an exercise to the reader (see also Glasserman [2004], Chapter 7).

3.3.2 Pathwise Estimate

3.3.2.1 Black-Scholes Delta

Reverting back to the setting of Section 3.3.1.1, let us take another look at the delta definition (3.68):

$$\begin{aligned}\frac{dV}{dS_0} &= \lim_{h \rightarrow 0} \frac{V(S_0 + h) - V(S_0)}{h} \\ &= e^{-rT} \lim_{h \rightarrow 0} \mathbb{E} \left(\frac{g(S_h(T)) - g(S(T))}{h} \right),\end{aligned}\tag{3.71}$$

where we have used the same notation as earlier:

$$S_h(T) = (S_0 + h) S(T)/S_0.\tag{3.72}$$

Under sufficient regularity on g , we can interchange expectation and limit in (3.71), such that simply

$$\begin{aligned}\frac{dV}{dS_0} &= e^{-rT} \mathbb{E} \left(\lim_{h \rightarrow 0} \frac{g(S_h(T)) - g(S(T))}{h} \right) \\ &= e^{-rT} \mathbb{E} \left(g'(S(T)) \frac{dS(T)}{dS_0} \right) \\ &= e^{-rT} \mathbb{E} \left(g'(S(T)) \frac{S(T)}{S_0} \right),\end{aligned}\tag{3.73}$$

¹⁶More precisely, we need uniform integrability in the difference coefficients $[Y(\alpha + \varepsilon) - Y(\alpha)]\varepsilon^{-1}$ and $\frac{1}{2}[Y(\alpha + \varepsilon) - Y(\alpha - \varepsilon)]\varepsilon^{-1}$. See Section 3.3.2.2 for more precise conditions.

where $g'(x) \triangleq dg/dx$, and the last equality follows by the linearity of (3.72).

We can implement the result (3.73) directly in a Monte Carlo trial, by generating samples of $S(T)$ and recording the sample averages of $g'(S(T))S(T)/S_0$. The resulting estimate for dV/dS_0 is a direct and unbiased estimate of the true derivative; it is known as a *pathwise estimate*.

For (3.73) to hold, g should be continuous, but does not necessarily need to be differentiable everywhere (it suffices that g is Lipschitz continuous, as discussed below). A regular call option payout $g(x) = (x - K)^+$, for instance, is non-differentiable at $x = K$, but we simply write

$$g'(x) = 1_{\{x>K\}}$$

and proceed directly with (3.73). For discontinuous payouts¹⁷, however, care must be taken as a direct application of (3.73) will introduce a bias. To demonstrate, consider the case $g(x) = 1_{\{x>K\}}$. Proceeding informally, a literal application of (3.73) results in

$$\frac{dV}{dS_0} = e^{-rT} \mathbb{E} \left(\delta(S(T) - K) \frac{S(T)}{S_0} \right),$$

where $\delta(x)$ is the Dirac delta function. While correct, this result is unsuited for Monte Carlo simulation: no matter how many samples n we draw of $S(T)$, the likelihood of $\delta(S(T) - K)$ being non-zero is zero, and the derivative would almost surely be estimated as 0. The correct result, however, is

$$\begin{aligned} e^{-rT} \mathbb{E} \left(\delta(S(T) - K) \frac{S(T)}{S_0} \right) &= e^{-rT} \mathbb{E} \left(\delta(S(T) - K) \frac{K}{S_0} \right) \\ &= e^{-rT} \frac{K}{S_0} \mathbb{E} (\delta(S(T) - K)) \\ &= e^{-rT} \frac{K}{S_0} \varphi_S(K), \end{aligned}$$

where $\varphi_S(\cdot)$ is the density of $S(T)$.

3.3.2.2 General Case

In a general setting, the technique employed in Section 3.3.2.1 above is known as *infinitesimal perturbation analysis*. A broad overview of the technique can be found in Glasserman [2004], with applications to finance covered in Broadie and Glasserman [1996]. Our treatment follows the latter closely.

Borrowing the notation of Section 3.3.1.2, we again consider estimating $dV/d\alpha$, where $V(\alpha) = \mathbb{E}(Y(\alpha))$. The basic idea of the pathwise derivative estimate is to write

¹⁷Or for the evaluation of, say, the second derivative (gamma) of a call payout.

$$\frac{dV}{d\alpha} = \frac{d}{d\alpha} \mathbb{E}(Y(\alpha)) = \mathbb{E}\left(\frac{d}{d\alpha} Y(\alpha)\right). \quad (3.74)$$

The exchange of expectation and differentiation requires certain regularity conditions to be valid. In practice, the most interesting situation arises when Y represents a (discounted) payout function, such that

$$Y(\alpha) = g(X(\alpha)),$$

where $X(\alpha) = (X_1(\alpha), \dots, X_q(\alpha))^{\top}$ is a q -dimensional random vector of observations (possibly at different dates) of asset prices. In this case, we have the following result, from Broadie and Glasserman [1996]:

Proposition 3.3.1. *For all α in some open interval \mathcal{A} assume that $dX_i/d\alpha$ exists almost surely for all $i = 1, \dots, q$. Suppose that the function g is almost surely differentiable¹⁸ and is Lipschitz, such that*

$$|g(x) - g(y)| \leq k|x - y|$$

for some constant k . Finally, assume that there exists finite-mean random variables β_i , $i = 1, \dots, q$, such that for all $\alpha_1, \alpha_2 \in \mathcal{A}$

$$|X_i(\alpha_2) - X_i(\alpha_1)| \leq \beta_i |\alpha_2 - \alpha_1|.$$

In this case, (3.74) holds.

The first two assumptions of Proposition ensure that the random variable $dY(\alpha)/d\alpha$ exists almost surely, with its value given by the chain rule

$$\frac{d}{d\alpha} Y(\alpha) = \sum_{i=1}^q \frac{\partial g}{\partial X_i} \frac{dX_i}{d\alpha}. \quad (3.75)$$

As we saw earlier, in Section 3.3.2.1, almost sure existence of $dY(\alpha)/d\alpha$ is not sufficient for the pathwise method to yield an unbiased estimator, we also need, roughly speaking, for g to be continuous at the points at which differentiability fails. The last two conditions ensure this, and together imply that Y is almost surely Lipschitz in α :

$$|Y(\alpha_2) - Y(\alpha_1)| \leq \beta_Y |\alpha_2 - \alpha_1|, \quad \beta_Y = k \sum_{i=1}^q \beta_i.$$

As $\alpha^{-1}|Y(h + \alpha) - Y(h)|$ is then bounded by β_Y , where $\mathbb{E}(\beta_Y) < \infty$, the result of Proposition 3.3.1 follows from the dominated convergence theorem. See Broadie and Glasserman [1996] for further details.

¹⁸That is, differentiable everywhere except on some set \mathcal{X} where $P(X(\alpha) \in \mathcal{X}) = 0$.

Remark 3.3.2. If in Proposition 3.3.1 we further assume that $E(\beta_i^2) < \infty$, it follows that $E(Y(h + \varepsilon) - Y(h)) \leq \beta_Y^2 \varepsilon^2$, such that

$$\text{Var}(Y(h + \varepsilon) - Y(h)) = O(\varepsilon^2).$$

We recognize this as case ii) from Section 3.3.1.2 on finite difference estimates, for which we have now made the regularity conditions more precise. In practice, the Lipschitz continuity of g is the critical condition.

Remark 3.3.3. For discontinuous payouts, the pathwise method will yield a biased estimator. As we saw in Section 3.3.2.1, however, for a simple process where the transition density is known, the bias can often be accounted for. We shall see an example of this in Chapter 24. Notice that if the transition density is known, another method for sensitivity simulation — the likelihood ratio method — also applies. See Section 3.3.3 below for details about this method.

3.3.2.3 Sensitivity Path Generation

In Section 3.3.2.1, simulation of $dY(\alpha)/d\alpha$ was straightforward due to the simplicity of the Black-Scholes dynamics. In general, we see from (3.75) that generation of the random variables $dY(\alpha)/d\alpha$ will require us to compute the partial derivatives of the payout with respect to the underlying assets ($\partial g/\partial X_i$), as well as the sensitivities of the assets with respect to the perturbation parameter α ($dX_i/d\alpha$). The latter is normally the most difficult, and we shall outline a general approach here.

For illustration, consider a scalar SDE of the usual form

$$dX(t) = \mu(t, X(t)) dt + \sigma(t, X(t)) dW(t),$$

where X and W are one-dimensional. Let α be a parameter on which $X(t)$ depends (such as $X(0)$ or some parameter of μ or σ). Let $D_\alpha(t)$ denote $dX(t)/d\alpha$. Formally differentiating the SDE with respect to α , we get

$$dD_\alpha(t) = \mu'(t, X(t)) D_\alpha(t) dt + \sigma'(t, X(t)) D_\alpha(t) dW(t), \quad (3.76)$$

where $\mu'(t, x) = \partial\mu(t, x)/\partial x$, and similar for σ' . This SDE can be discretized and simulated in parallel with the simulation of the SDE for $X(t)$ itself. In general, the work associated with this will obviously be more substantial than for the Black-Scholes delta, where we saw that $D_\alpha(t)$ could be recovered as the simple fraction $X(t)/X(0)$ (see equation (3.73)).

A few notes on the technique above. First, some regularity is obviously needed for (3.76) to be meaningful; see Kunita [1990] for some relevant results. Second, extensions to multi-dimensional SDEs are straightforward, although the dimension of the total scheme can be large. For instance, if X is p -dimensional and we wish to compute sensitivities with respect to $X_i(0)$,

$i = 1, \dots, p$, a $p \times p$ system of SDEs for quantities $dX_i(t)/dX_j(0)$ will be required. We shall discuss approximative methods to improve efficiency of such high-dimensional matrix SDEs later, in Chapter 24 (see also Glasserman and Zhao [1999]).

3.3.3 Likelihood Ratio Method

As discussed above, the pathwise derivative method typically applies only to options with sufficiently smooth payouts¹⁹ and can be cumbersome for multi-dimensional SDEs. For processes with explicitly known transition densities, the alternative *likelihood ratio method* can be used. This method applies to discontinuous payout functions and, unlike the pathwise method, requires little knowledge of the payout function and its derivatives, making it convenient for general implementation on a computer. When both methods apply, however, the pathwise derivative method generally is more efficient.

3.3.3.1 Black-Scholes Delta

In the notation of Sections 3.3.1.1 and 3.3.2.1, the Black-Scholes price of a call option can be written

$$\begin{aligned} V(S_0) &= e^{-rT} \mathbb{E} \left((S(T) - K)^+ \mid S(0) = S_0 \right) \\ &= e^{-rT} \mathbb{E} \left(\left(e^{\ln S_0 + (r - \frac{1}{2}\sigma^2)T + \sigma\sqrt{T}Z} - K \right)^+ \right) \\ &= e^{-rT} \mathbb{E} \left(\left(e^{Y(T)} - K \right)^+ \right), \end{aligned}$$

where $Z \sim \mathcal{N}(0, 1)$ and

$$Y(T) \sim \mathcal{N} \left(\ln S_0 + \left(r - \frac{1}{2}\sigma^2 \right) T, \sigma^2 T \right).$$

The density of Y is thereby a function of S_0 :

$$\begin{aligned} P(Y \in dy) &= \frac{dy}{\sqrt{2\pi}\sigma\sqrt{T}} \exp \left(-\frac{1}{2} \left(\frac{y - \ln S_0 - (r - \frac{1}{2}\sigma^2)T}{\sigma\sqrt{T}} \right)^2 \right) \\ &\triangleq \varphi(y; S_0) dy. \end{aligned}$$

Thereby

$$V(S_0) = e^{-rT} \int_{-\infty}^{\infty} (e^y - K)^+ \varphi(y; S_0) dy, \quad (3.77)$$

such that

¹⁹But see Remark 3.3.3.

$$\begin{aligned}
\frac{dV(S_0)}{dS_0} &= e^{-rT} \int_{-\infty}^{\infty} (e^y - K)^+ \frac{\partial \varphi(y; S_0)}{\partial S_0} dy \\
&= e^{-rT} \int_{-\infty}^{\infty} (e^y - K)^+ \frac{\partial \varphi(y; S_0)}{\partial S_0} \frac{\varphi(y; S_0)}{\varphi(y; S_0)} dy \\
&= e^{-rT} \int_{-\infty}^{\infty} (e^y - K)^+ \frac{\partial \ln \varphi(y; S_0)}{\partial S_0} \varphi(y; S_0) dy. \tag{3.78}
\end{aligned}$$

Comparison of (3.78) with (3.77) demonstrates that we can effectively compute the Black-Scholes delta as the price of a security that pays out at time T the amount

$$(e^{Y(T)} - K)^+ l(Y(T)), \tag{3.79}$$

where l is the so-called *log-likelihood ratio* (also known as the *score function*)

$$l(Y(T)) = \frac{\partial \ln \varphi(Y(T); S_0)}{\partial S_0} = \frac{Y(T) - \ln S_0 - (r - \frac{1}{2}\sigma^2)T}{S_0 \sigma^2 T} = \frac{Z}{S_0 \sigma \sqrt{T}}. \tag{3.80}$$

By differentiating the density, rather than the payout itself, the likelihood ratio technique applies to even discontinuous payouts, requiring only that the density is smooth (which is clearly the case here). Notice in particular that the log-likelihood ratio is independent of the payout, allowing us to use the same function $l(Y(T))$ for all European style payout functions $V(T) = g(S(T))$.

3.3.3.2 General Case

As in Section 3.3.2.2, consider now the general case where a random variable $Y(\alpha)$ represents a (deflated) payout function g applied to a vector of random variables $X(\alpha) = (X_1(\alpha), \dots, X_q(\alpha))^{\top}$. Again, α is a parameter with respect to which we wish to compute sensitivities. Let the joint density of $X(\alpha)$ be denoted $f(x; \alpha)$, $x \in \mathbb{R}^q$. We then have

$$V(\alpha) = \int_{\mathbb{R}^q} g(x) f(x; \alpha) dx.$$

Making the reasonable assumption that density $f(x; \alpha)$ is a smooth function of α , we interchange integration and differentiation, such that

$$\frac{\partial V(\alpha)}{\partial \alpha} = \int_{\mathbb{R}^q} g(x) \frac{\partial f(x; \alpha)}{\partial \alpha} dx = \int_{\mathbb{R}^q} g(x) \frac{\partial \ln f(x; \alpha)}{\partial \alpha} f(x; \alpha) dx.$$

As for the Black-Scholes case above, the derivative $\partial V(\alpha)/\partial \alpha$ can thus be computed as the expectation of the payout modified by a log-likelihood ratio:

$$g(x)l(x), \quad l(x) = \frac{\partial \ln f(x; \alpha)}{\partial \alpha}.$$

3.3.3.3 Euler Schemes

In practice, the reliance on explicit knowledge of a transition density can be a considerable obstacle, and may rule out the application of the likelihood ratio method for many complex models. In cases where process dynamics are simulated through a simple time-discretization scheme, the situation is, however, salvageable, as we shall now demonstrate.

For illustration, consider an asset $X(t)$ that follows an SDE of the type

$$dX(t) = \mu(t, X(t); \alpha) dt + \sigma(t, X(t); \alpha) dW(t),$$

where μ and σ are smooth functions, and where α is a parameter. In general, we do not know the exact transition density for $X(t)$. However, suppose now that we use an Euler scheme to simulate on some grid $\{t_i\}_{i=1}^m$, i.e.

$$\widehat{X}(t_{i+1}) = \widehat{X}(t_i) + \mu(t_i, \widehat{X}(t_i); \alpha) (t_{i+1} - t_i) + \sigma(t_i, \widehat{X}(t_i); \alpha) \sqrt{t_{i+1} - t_i} Z_i,$$

where $t_0 = 0$ and Z_0, Z_1, \dots, Z_{m-1} is a sequence of i.i.d. standard Gaussian random variables. Clearly, the transition density for the $\widehat{X}(t_{i+1})$ is now Gaussian,

$$\widehat{X}(t_{i+1}) | \widehat{X}(t_i) \sim \mathcal{N}\left(m_i(\widehat{X}(t_i); \alpha), s_i(\widehat{X}(t_i); \alpha)\right),$$

where, for $i = 0, \dots, m-1$,

$$\begin{aligned} m_i(\widehat{X}(t_i); \alpha) &= \widehat{X}(t_i) + \mu(t_i, \widehat{X}(t_i); \alpha) (t_{i+1} - t_i), \\ s_i(\widehat{X}(t_i); \alpha) &= \sigma(t_i, \widehat{X}(t_i); \alpha)^2 (t_{i+1} - t_i). \end{aligned}$$

Set $\widehat{X} = (\widehat{X}(t_1), \dots, \widehat{X}(t_m))^\top$; the density of this vector is (where $x_0 = X(0)$)

$$f(x_1, \dots, x_m; \alpha) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi} \sqrt{s_{i-1}(x_{i-1}; \alpha)}} \exp\left(-\frac{(x_i - m_{i-1}(x_{i-1}; \alpha))^2}{2s_{i-1}(x_{i-1}; \alpha)}\right). \quad (3.81)$$

Consider some (potentially path-dependent) security V with payout function $g(\widehat{X})$ and time 0 price of $V(\alpha)$. Equipped with (3.81), we can estimate the parameter sensitivity by the likelihood ratio method as

$$\frac{\partial V(\alpha)}{\partial \alpha} = E\left(g(\widehat{X}) \frac{d}{d\alpha} \ln(f(\widehat{X}; \alpha))\right),$$

where

$$\begin{aligned} &\frac{d}{d\alpha} \ln(f(\widehat{X}; \alpha)) \\ &= \frac{d}{d\alpha} \sum_{i=1}^m \left(-\frac{1}{2} \ln(s_{i-1}(x_{i-1}; \alpha)) - \frac{(x_i - m_{i-1}(x_{i-1}; \alpha))^2}{2s_{i-1}(x_{i-1}; \alpha)}\right). \quad (3.82) \end{aligned}$$

This derivative can typically be computed in closed form; if not, one can estimate it by finite differences (as in Su and Randall [2008]). Notice that when $\alpha = X(0)$ (i.e. we are trying to estimate the delta), only s_0 and m_0 will depend on α , simplifying computations.

The idea used above extends easily to vector-valued $X(t)$. In principle, higher-order schemes (e.g. the Milstein scheme) can also be used, although the complexity increases considerably.

3.3.3.4 Some Remarks

The main advantage of the likelihood ratio is the fact that it applies to classes of payouts for which other methods (pathwise methods, finite difference method) do not work well. Moreover, the method is easy and efficient to implement, as the log-likelihood ratio l is independent of the payout and does not — unlike the general pathwise method — require simulation of any quantities other than the vector X itself. As discussed above, the primary drawback of the method is its reliance on explicit knowledge of the process density, ruling out many of the more advanced models (although, as Section 3.3.3.3 demonstrated, there may sometimes be ways around this). Further, the variance of the likelihood ratio method can often be quite big, particularly if the parameter α simultaneously affects multiple stochastic variables. A fuller discussion of this issue, as well as the related issue of absolute continuity, can be found in Glasserman [2004]. We note in passing that the likelihood ratio method is a special case of a body of methods that have emerged from the so-called *Malliavin calculus*; see Fournie et al. [1999] for a survey. Most Malliavin methods other than the basic likelihood ratio method are, however, not particularly attractive due to computational issues²⁰.

We round off this section by noting that the various methods for derivative estimates can often be successfully combined. For instance, while it is common to use either the pathwise method or the likelihood ratio method to compute first order sensitivities (such as delta), second-order sensitivities (such as gamma) are often done by the finite difference method applied to first-order sensitivities, often using fairly sizable shifts of the underlying variables. By combining the pathwise method with the likelihood ratio method, we can also address the fact that the first derivative of many kinked option payouts is discontinuous, allowing us to produce a bias-free estimate of the second derivative. See Fournie et al. [1999] for other examples of combining the pathwise method with the likelihood ratio method.

²⁰Besides, the Malliavin calculus itself, a very technical area of mathematics even for specialists, can be avoided altogether, as Chen and Glasserman [2007b] demonstrate.

3.4 Variance Reduction Techniques

As discussed earlier, the convergence of the Monte Carlo method is quite slow, of order $O(n^{-1/2})$ where n is the number of Monte Carlo samples. While there is little that can be done to improve²¹ the $n^{-1/2}$ order itself, the constant multiplying $n^{-1/2}$ can be affected by a careful choice of the Monte Carlo estimator. Methods to improve numerical efficiency this way are known collectively as *variance reduction techniques*, and constitute a major area of research in the theory of Monte Carlo methods. Our introduction of the topic is limited to a few basic examples. More details are provided later in the book for concrete models and products (see e.g. Chapter 25), and more information can be found in the standard Monte Carlo literature, including Hammersley and Handscomb [1965] and the survey article Boyle et al. [1997].

3.4.1 Variance Reduction and Efficiency

We recall that the goal of the Monte Carlo method is to estimate some quantity μ (e.g., the price of a financial contract) as the sample mean of n i.i.d. random variables Y_1, \dots, Y_n , where each Y_i has expectation $E(Y_i) = \mu$ and variance $\text{Var}(Y_i) = \sigma^2$. From Section 3.1, we know that for large n the standard error of the sample mean $n^{-1} \sum Y_i$ is $\sigma n^{-1/2}$, with the probabilistic error bounds on the estimate of μ being proportional to the standard error.

Suppose now that we have available two sets of i.i.d. sequences $Y_{1,i}$ and $Y_{2,i}$, $i = 1, \dots, n$, where $E(Y_{1,i}) = E(Y_{2,i}) = \mu$, but $\text{Var}(Y_{1,i}) = \sigma_1^2$ and $\text{Var}(Y_{2,i}) = \sigma_2^2$, with $\sigma_1 \neq \sigma_2$. Also suppose that the time it takes on a computer to generate individual samples $Y_{1,i}$ and $Y_{2,i}$ is τ_1 and τ_2 , respectively. Which of the two estimators $n^{-1} \sum Y_{1,i}$ and $n^{-1} \sum Y_{2,i}$ is preferable? To answer this question, assume that we have a large fixed computing time budget τ . The number of replications of $Y_{1,i}$ and $Y_{2,i}$ that can be executed are thus (the integer parts of) τ/τ_1 and τ/τ_2 , respectively. To this correspond sample mean standard errors of

$$\frac{\sigma_1}{\sqrt{\tau/\tau_1}} \text{ and } \frac{\sigma_2}{\sqrt{\tau/\tau_2}},$$

respectively. It follows that, for large τ , the estimator based on the sequence $Y_{1,i}$, $i = 1, \dots, n$, is preferable, if

$$\frac{\sigma_1}{\sqrt{\tau/\tau_1}} < \frac{\sigma_2}{\sqrt{\tau/\tau_2}},$$

or equivalently

²¹As we discussed earlier, the quasi-random Monte Carlo method can theoretically achieve better convergence order than $O(n^{-1/2})$.

$$\sigma_1^2 \tau_1 < \sigma_2^2 \tau_2.$$

For obvious reasons, the product of variance and per-sample computing time is known as the *efficiency* of a Monte Carlo estimator. In devising methods to improve Monte Carlo performance, efficiency should always constitute the measure of comparison. For instance, a high-variance estimator may, in fact, be preferable to a low-variance estimator, provided that the former takes less time to compute than the latter.

Duffie and Glynn [1995] discuss Monte Carlo efficiency in more depth, with additional analysis of the effects of bias (see also Section 3.2.8) and cases where τ_1 and τ_2 are random.

3.4.2 Antithetic Variates

3.4.2.1 The Gaussian Case

A simple and easily implemented variance reduction technique is the method of *antithetic variates*. Assume that we are interested in estimating the expected value of a random variable $Y = G(Z)$, where G is a real-valued function and Z is a q -dimensional vector of independent standard Gaussian random variables. This problem routinely arises in determining the expected value of a function of assets driven by a vector SDE; Z then represents the aggregation of all independent standard Gaussian variables used to produce Brownian motion increments; see for instance Section 3.2.3. For n independent realizations of Z , Z_1, \dots, Z_n , rather than using the regular sample average estimator for $E(Y)$, consider instead using

$$\bar{Y}_n^a = n^{-1} \sum_{i=1}^n \frac{G(Z_i) + G(-Z_i)}{2}.$$

In other words, in addition to the set Z_1, \dots, Z_n of Gaussian samples, we also effectively include the set $-Z_1, \dots, -Z_n$ in the Monte Carlo trial. As $-Z_1, \dots, -Z_n$ itself is a sequence of n independent Gaussian samples, we still must have

$$E(\bar{Y}_n^a) = E(Y),$$

so the antithetic estimator is unbiased. Also, as $G(Z)$ and $G(-Z)$ have identical variance,

$$\text{Var}(\bar{Y}_n^a) = n^{-1} \left[\frac{1}{2} \text{Var}(Y) + \frac{1}{2} \text{Cov}(G(Z), G(-Z)) \right] = \frac{\text{Var}(Y)}{n} \frac{(1 + \rho)}{2},$$

where ρ is the correlation between $G(Z)$ and $G(-Z)$. Recalling that the regular sample average has variance

$$\text{Var}(\bar{Y}_n) = \frac{\text{Var}(Y)}{n},$$

we conclude that $\text{Var}(\bar{Y}_n^a) < \text{Var}(\bar{Y}_n)$ as long as $\rho < 1$ (which is obviously likely).

While use of antithetic variates can always be expected to lower the standard error, it is not necessarily more efficient than regular Monte Carlo, in the sense defined in Section 3.4.1. For instance, if generation of the Z_i 's is of negligible cost relative to the evaluation of $G(Z)$, computation of \bar{Y}_n^a will take about twice as long as the regular sample average \bar{Y}_n . For this case, the results in Section 3.4.1 show that for antithetic variates to constitute an improvement in computational efficiency, we must require that

$$\text{Var}(\bar{Y}_n^a) \leq \frac{1}{2}\text{Var}(\bar{Y}_n),$$

or

$$\rho \leq 0.$$

A sufficient condition for $\rho < 0$ is that G be monotone in all q elements of Z . Given this, we would expect antithetic variates to be most suitable for option payouts that depend monotonically on prices.

3.4.2.2 General Case

While the method of antithetic variates is primarily associated with the idea of changing signs on Gaussian variables, the method can, in fact, be extended to other distributions. At the most basic level, most simulation trials involve a series of uniform draws that are translated to other random variables, using techniques described in Section 3.1.1. In this case, we can focus our attention on estimating the mean of a random variable $Y = H(U)$, where H is a function and U is a vector of independent uniformly distributed random variables. We notice that if $U = (U_1, \dots, U_q)^\top$ is a vector of independent uniform random variables on $[0, 1]$, then so is $\tilde{U} = (1 - U_1, \dots, 1 - U_q)^\top$. The pair $\{U, \tilde{U}\}$ is thereby antithetic (negatively dependent) in the same way as the Gaussian pair $\{Z, -Z\}$ above, and we can estimate the mean of Y as the average of independent samples of the form

$$\frac{H(U) + H(\tilde{U})}{2}.$$

From the discussion above, it follows that if H is monotonic in U , the resulting scheme will exhibit better computational efficiency than regular Monte Carlo.

As an aside, we note that the simple “reflection” of a vector of uniforms advocated above is, as should be obvious, not the only possible way of generating an antithetic sample — for instance, we could have chosen to reflect only select dimensions of the U -vector. A similar observation holds for the case of vector-valued Gaussian variables. The general idea of applying deterministic transformations to a vector-valued sample of random numbers as a way to reduce variance is sometimes known as *systematic sampling*.

3.4.3 Control Variates

3.4.3.1 Basic Idea

While we may need to use Monte Carlo simulation to estimate the unknown mean of a random variable Y , there may be random variables “close” to Y with means that can be computed analytically. It seems reasonable that the additional information about Y revealed by these random variables could be useful in improving our estimate of $E(Y)$. While a number of strategies are possible²², we shall here focus on the so-called *control variate* method.

Formally, let

$$Y^c = (Y_1^c, \dots, Y_q^c)^\top$$

be a vector of *control variates* (or just *controls*), ideally with strong negative or positive correlation to a variable Y . The mean of Y^c is known to be

$$E(Y^c) = \mu^c = (\mu_1^c, \dots, \mu_q^c)^\top.$$

Now, introduce an arbitrary constant vector

$$\beta = (\beta_1, \dots, \beta_q)^\top$$

and consider forming the linear combination

$$X = Y - \beta^\top (Y^c - \mu^c). \quad (3.83)$$

Clearly

$$E(X) = E(Y) - \beta^\top (E(Y^c) - \mu^c) = E(Y),$$

so using Monte Carlo sampling to estimate the mean of X will provide an unbiased estimate of $E(Y)$, regardless of the choice of β .

To analyze the variance of the new variable X , let Σ_{Y^c} be the $q \times q$ covariance matrix of the vector Y^c , and let Σ_{Y,Y^c} be the q -dimensional vector of covariances between Y and the components of Y^c . The variance of X can then be shown to be

$$\text{Var}(X) = \text{Var}(Y) - 2\beta^\top \Sigma_{Y,Y^c} + \beta^\top \Sigma_{Y^c} \beta. \quad (3.84)$$

Whether or not this constitutes an improvement (in the sense that $\text{Var}(X) < \text{Var}(Y)$) is largely a matter of what β is chosen to be. We have the following easily proven lemma.

Lemma 3.4.1. *The function $\text{Var}(X) = \text{Var}(Y) - 2\beta^\top \Sigma_{Y,Y^c} + \beta^\top \Sigma_{Y^c} \beta$ is minimized at*

$$\beta^* = \Sigma_{Y^c}^{-1} \Sigma_{Y,Y^c}$$

²²Other methods include *moment matching* and *importance sampling*. We shall cover the latter strategy shortly; the former is discussed in Boyle et al. [1997], where it is concluded that control variates are superior, at least asymptotically.

with minimum value

$$\min_{\beta} \text{Var}(X) = (1 - R^2)\text{Var}(Y), \quad R^2 \triangleq \frac{\Sigma_{Y,Y^c}^\top \Sigma_{Y^c}^{-1} \Sigma_{Y,Y^c}}{\text{Var}(Y)} \geq 0. \quad (3.85)$$

In the lemma, we recognize the scalar R^2 as the R -squared of a multi-dimensional regression of Y against Y^c . Similarly, the components of the optimal vector β^* are the regression coefficients (the slopes) on the vector Y^c . In practice, we may not know Σ_{Y^c} and Σ_{Y,Y^c} explicitly, in which case we simply replace these with empirical estimates, as obtained by an n -sample Monte Carlo trial. We note that if the random samples used to estimate β^* are the same as those used to estimate $E(X)$, a small bias is typically introduced. This can be circumvented by using separate random numbers for the estimates of β^* and $E(X)$, but in practice this is rarely worth the effort. Nelson [1990], among others, analyzes this issue in more detail.

While the usage of control variates will always lower variance (unless Y and Y^c are perfectly uncorrelated), an improvement of computational efficiency over standard Monte Carlo is, of course, not guaranteed. Consider, for instance, the case where the computational effort involved in generating a single sample of X is $q + 1$ times that of generating Y itself. This will be the case, if i) the effort of drawing random numbers is small relative to computing Y itself; and ii) each of the components of Y^q take about the same time to compute as Y . According to the result in Section 3.4.1, for this special case the control variate method will only entail an increase in efficiency, if

$$(1 - R^2)\text{Var}(Y)(q + 1) < \text{Var}(Y)$$

or

$$1 - R^2 < \frac{1}{q + 1}.$$

As q grows large, this requirement obviously becomes increasingly difficult to satisfy. Rather than indiscriminately adding multiple controls, it is therefore normally best to properly analyze a given problem and use only a few well-chosen variables with strong (negative or positive) correlation to the variable in question.

3.4.3.2 Non-Linear Controls

Our discussion of the control variate method has so far only considered linear controls (3.83), where the modified estimator involves a linear combination of control variates. The resulting estimate of $E(Y)$ are n -point sample averages of the type

$$\bar{Y}_n - \beta^\top (\bar{Y}_n^c - \mu^c).$$

A more general formulation than (3.83) approximates $E(Y)$ with

$$f(\bar{Y}_n^c, \bar{Y}_n) \quad (3.86)$$

for some function f satisfying

$$f(\mu^c, y) = y. \quad (3.87)$$

The requirement (3.87) ensures that $f(\bar{Y}_n^c, \bar{Y}_n)$ approaches $E(Y)$ in the large-sample limit; unlike the regular control variate formulation, however, (3.86) may involve a bias for finite sample sizes.

If f is smooth, a result by Glynn and Whitt [1989] demonstrates that for sufficiently large samples, any non-linear control variate estimator of the type (3.86) is equivalent to an ordinary linear control variate estimator. Still, there may be situations where a non-linear control variate estimator is appropriate, either because i) the sample size is not large enough to justify the result in Glynn and Whitt [1989]; or ii) because the “effective” β weighting of \bar{Y}_n^c implied by f is close to optimal, allowing us to skip the estimation of β^* .

To give an example of non-linear control variates, let us consider the “delta” method of Clewlow and Carverhill [1994]. To state the basic idea, consider the estimate of

$$V(0) = E(g(X(T))),$$

where $X(t)$ is a p -dimensional vector process and $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is a smooth function. Assume that all components of X are martingales, as is the case when X represents assets deflated by a numeraire. We recall from Section 1.7 that, under certain regularity conditions, we have

$$V(T) = V(0) + \int_0^T \sum_{i=1}^p V_{x_i}(t) dX_i(t),$$

where we use the notation $V_{x_i}(t)$ from Section 1.7 to denote, informally, $V_{x_i}(t) = \partial V(t)/\partial X_i(t)$. On a simulation time line $\{t_j\}_{j=1}^m$, we can write, in the style of an Euler scheme,

$$V(T) \approx V(0) + \sum_{j=1}^m \sum_{i=1}^p V_{x_i}(t_{j-1}) (X_i(t_j) - X_i(t_{j-1})).$$

As the zero-mean quantity

$$\sum_{j=1}^m \sum_{i=1}^p V_{x_i}(t_{j-1}) (X_i(t_j) - X_i(t_{j-1}))$$

is likely to have high correlation to $V(T)$, we can consider using it as a control variate. One obstacle is the fact that the derivatives $V_{x_i}(t)$ are likely

to be unknown (as the function $V(t)$ is unknown). Often, however, we can provide an inspired guess for these derivatives, based on perhaps a simpler model or on regression information. The former idea is outlined in Clewlow and Carverhill [1994], and the latter shall be discussed further in Chapter 25. In any case, the resulting scheme ends up effectively using the increments $X_i(t_j) - X_i(t_{j-1})$ as controls, with non-constant weights $V_{x_i}(t_{j-1})$ being functions of the X_i themselves.

3.4.4 Importance Sampling

3.4.4.1 Basic Idea

The basic idea of the *importance sampling method* is to use a measure shift to reduce variance. For a given measure P , consider estimating

$$\mu = E^P(Y), \quad (3.88)$$

where Y is a scalar random variable. Let \widehat{P} be a measure equivalent to P . From the Radon-Nikodym theorem in Chapter 1, we have

$$\mu = E^{\widehat{P}}(Y/R), \quad (3.89)$$

where R is the Radon-Nikodym derivative

$$R = d\widehat{P}/dP, \quad E^P(R) = 1.$$

While (3.88) and (3.89) are both valid expressions for μ , it is possible that the variance of Y/R under measure \widehat{P} is lower than the variance of Y under P , making (3.89) potentially more efficient for Monte Carlo purposes. As an extreme case, consider setting (assuming $Y > 0$ a.s.)

$$R = \frac{Y}{E^P(Y)} = \frac{Y}{\mu}. \quad (3.90)$$

In this case

$$Y/R = \mu$$

and non-random, implying that the measure shift from P to \widehat{P} has removed *all* variance. The problem with the “perfect” choice (3.90) is obviously that we do not know μ — if we did, there would be no need to estimate it by Monte Carlo methods. Nevertheless, we may be able to provide a good guess for μ , allowing us to use (3.90) in an approximate sense.

3.4.4.2 Density Formulation

Importance sampling methods are often most conveniently (and most intuitively) treated in terms of probability densities, so let us cast the description

of Section 3.4.4.1 in such terms. Specifically, let us assume that Y can be represented as $g(X)$, where $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is a well-behaved function and X is p -dimensional with probability density $f : \mathbb{R}^p \rightarrow \mathbb{R}$. We then write

$$\mu = \mathbb{E}^P(g(X)) = \int_{\mathbb{R}^p} g(x) f(x) dx,$$

to which corresponds a regular Monte Carlo estimator

$$\bar{\mu}_n = \frac{1}{n} \sum_{i=1}^n g(X_i),$$

where the X_i are independent samples of X , drawn from the density f . Let $h : \mathbb{R}^p \rightarrow \mathbb{R}$ be another density, satisfying the continuity requirement that $h(x) > 0$ whenever $f(x) > 0$. We can then also represent μ as

$$\mu = \int_{\mathbb{R}^p} g(x) \frac{f(x)}{h(x)} h(x) dx,$$

which we can interpret as

$$\mu = \mathbb{E}^{\widehat{P}} \left(g(X) \frac{f(X)}{h(X)} \right),$$

where \widehat{P} is a measure under which X has density $h(x)$. Comparison to the results above identifies the so-called *likelihood ratio* $l(x) = f(x)/h(x)$ as the Radon-Nikodym derivative $dP/d\widehat{P}$ (or $1/R$) governing the shift from P to \widehat{P} . If now X_1, \dots, X_n are independent draws from h (and *not* f), the importance sampling Monte Carlo estimator for μ takes the form

$$\bar{\mu}_n^h = \frac{1}{n} \sum_{i=1}^n g(X_i) \frac{f(X_i)}{h(X_i)}.$$

Let us investigate under which circumstances importance sampling will lead to an improvement in variance. We have

$$\begin{aligned} \text{Var}(\bar{\mu}_n^h) &= \frac{1}{n} \left[\mathbb{E}^{\widehat{P}} \left(g(X)^2 \frac{f(X)^2}{h(X)^2} \right) - \mu^2 \right] \\ &= \frac{1}{n} \left[\mathbb{E}^P \left(g(X)^2 \frac{f(X)}{h(X)} \right) - \mu^2 \right], \end{aligned}$$

and

$$\text{Var}(\bar{\mu}_n) = \frac{1}{n} \left[\mathbb{E}^P \left(g(X)^2 \right) - \mu^2 \right].$$

Hence, importance sampling will lower variance, provided that

$$\mathbb{E}^P \left(g(X)^2 \frac{f(X)}{h(X)} \right) < \mathbb{E}^P \left(g(X)^2 \right).$$

Choosing the importance sampling density $h(x)$ wisely is key to the efficiency of the importance sampling. As an extreme, suppose we could set

$$h(x) = Cf(x)g(x), \quad (3.91)$$

where the constant C is dictated by the need for $h(x)$ to integrate to 1:

$$C^{-1} = \int_{\mathbb{R}^p} g(x)f(x) dx = \mu.$$

In this case,

$$\mathbb{E}^P \left(g(X)^2 \frac{f(X)}{h(X)} \right) = C^{-1} \mathbb{E}^P (g(X)) = \mu^2$$

and

$$\text{Var} (\bar{\mu}_n^h) = 0.$$

This replicates a similar argument in Section 3.4.4.1 (see equation (3.90)), and is equally useless in practice: to compute (3.91) we need to normalize by the constant $1/\mu$, where μ is the quantity that we are trying to estimate in the first place. Nevertheless, (3.91) provides some useful practical guidance: a good choice of likelihood density will sample in proportion to f and g . That is, values of X where both the density $f(X)$ and the payout $g(X)$ are high should be assigned a high value of $h(X)$ (high “importance”), and values of X where either $f(X)$ or $g(X)$ (or both) are low should be assigned a low value of $h(X)$ (low “importance”). This rule is often particularly easy and efficient to apply to situations where $g(X)$ is significant only for a set $X \in \mathcal{A}$, where $P(X \in \mathcal{A})$ is small. Such rare-event problems are a classical application of importance sampling; we give a simple example in Section 3.4.4.5. Related applications to barrier options can be found later, in Chapter 25, with more such examples in Boyle et al. [1997].

3.4.4.3 Importance Sampling and SDEs

Consider now a dynamic setting where we are given a P-measure SDE

$$dX(t) = \mu(t, X(t)) dt + \sigma(t, X(t)) dW(t), \quad (3.92)$$

where X is p -dimensional and W is d -dimensional. We wish to evaluate

$$\mathbb{E}^P (g(X(T)))$$

for a real-valued function g . To shift measure, we introduce the density process

$$d\varsigma(t) = -\varsigma(t)\theta(t)^\top dW(t), \quad \varsigma(0) = 1, \quad (3.93)$$

for some adapted d -dimensional process $\theta(t)$, sufficiently regular to make $\varsigma(\cdot)$ a martingale (see Chapter 1). Let

$$\varsigma(t) = E_t^P \left(\frac{d\hat{P}}{dP} \right),$$

for a new measure \hat{P} . By the Girsanov theorem, under \hat{P} ,

$$dX(t) = [\mu(t, X(t)) - \sigma(t, X(t))\theta(t)] dt + \sigma(t, X(t)) d\hat{W}(t), \quad (3.94)$$

where \hat{W} is a Brownian motion in \hat{P} . Also, by the Radon-Nikodym theorem,

$$E^P(g(X(T))) = E^{\hat{P}} \left(\frac{g(X(T))}{\varsigma(T)} \right). \quad (3.95)$$

In a Monte Carlo setting, rather than simulating (3.92) (using methods from Section 3.2) and computing the sample mean of $g(X(T))$, we can instead jointly simulate (3.93) and (3.94) and compute the sample mean of $g(X(T))/\varsigma(T)$. The validity of this approach is independent of the choice of θ in (3.93), and we can use θ as a parameter to minimize the variance of $g(X(T))/\varsigma(T)$ under \hat{P} .

To find the optimal choice for θ , define

$$u(t, X(t)) = E_t^P(g(X(T))), \quad t \leq T,$$

and consider setting

$$\varsigma(t)u(0, X(0)) = u(t, X(t)).$$

By Ito's lemma,

$$d\varsigma(t) = -\varsigma(t)\theta(t)^\top dW(t), \quad \varsigma(0) = 1,$$

where

$$\theta(t) = -u(t, X(t))^{-1} \sigma(t, X(t))^\top \frac{\partial u(t, X(t))}{\partial x}, \quad (3.96)$$

with $\partial u(t, X(t))/\partial x$ being a p -dimensional vector of partial derivatives $\{\partial u(t, X(t))/\partial x_i\}$. The choice for θ in (3.96) is optimal as we have

$$g(X(T))/\varsigma(T) = u(0, X(0)) = E^P(g(X(T))),$$

which is non-random with zero variance. As in earlier examples, the optimal choice for $\theta(t)$ cannot be applied directly as it requires knowledge of $E_t^P(g(X(T)))$ for all t , knowledge which we never possess in practice. In many applications, however, we can often make an educated guess for u , based perhaps on either a simpler SDE than (3.92) or on a simpler payout function than g . We shall see an example of this in Chapter 25; another application can be found in Schoenmakers and Heemink [1997].

3.4.4.4 More on SDE Path Simulation

Let us consider an alternative point of view about SDE simulations, where we assume that the SDE (3.92) is simulated by an m -dimensional Euler scheme (or similar), such that we can write (see also Section 3.4.2.1) for some function $G : \mathbb{R}^{p \times m} \rightarrow \mathbb{R}$,

$$g(X(T)) = G(Z_1, \dots, Z_m),$$

where the Z_i are independent p -dimensional Gaussian vectors. With the Gaussian density of Z_i being denoted $\phi(z)$, $z \in \mathbb{R}^p$, the independence of the Z_i 's allows us to write

$$\mathbb{E}^P(g(X(T))) = \int_{\mathbb{R}^{p \times m}} G(z_1, \dots, z_m) \prod_{i=1}^m \phi(z_i) dz, \quad z \triangleq (z_1, \dots, z_m).$$

If we apply a change of measure that preserves independence of Z_i but alters the common marginal density from $\phi(z)$ to $h(z)$, the likelihood ratio is easily seen to be

$$l(z) = \prod_{i=1}^m \frac{\phi(z_i)}{h(z_i)},$$

such that

$$\mathbb{E}^P(g(X(T))) = \mathbb{E}^{\widehat{P}} \left(G(Z_1, \dots, Z_m) \prod_{i=1}^m \frac{\phi(Z_i)}{h(Z_i)} \right).$$

It is understood that the Z_i used to advance the SDE simulation under \widehat{P} are drawn from the density h , rather than ϕ .

To give a concrete example of a measure shift, assume for simplicity that $p = 1$ and consider shifting the means of the Z_i from zero to some scalar²³ μ , but retaining unit variance. For this, we must set

$$h(z_i) = \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{1}{2}(z_i - \mu)^2 \right),$$

whereby

$$l(z) = l(z; \mu) = \exp \left(-\mu \sum_{i=1}^m z_i + \frac{m}{2} \mu^2 \right). \quad (3.97)$$

Here μ is a free variable, which can be set to minimize the variance of the term

$$G(Z)l(Z; \mu), \quad Z \triangleq (Z_1, \dots, Z_m)$$

²³It is also straightforward to introduce a measure shift that moves the means of the Z_i to *different* means μ_i , $i = 1, \dots, m$.

under \widehat{P} . Sometimes this minimization problem can be handled analytically (see Section 3.4.4.5), but most often numerical methods are required. Examples of how to perform this minimization by Monte Carlo simulation can be found in, for example, Su and Fu [2002] and Capriotti [2007]. The approach in Capriotti [2007] (called *least-squares importance sampling*) is particularly straightforward, as the optimization problem is here cast as a least-squares regression problem for which well-known numerical schemes exist such as, e.g., the Levenberg-Marquardt routine in Press et al. [1992]. Both Su and Fu [2002] and Capriotti [2007] point out that, when computing variance, it is advantageous to cast the problem back into the original probability measure P by using

$$E^{\widehat{P}}(G(Z)^2 l(Z; \mu)^2) = E^P(G(Z)^2 l(Z; \mu)).$$

Let us finally note that the measure transformation employed above is a special case of so-called *exponential twisting* (also known as *Esscher transform*), under which a density $f(x)$, $x \in \mathbb{R}$, is transformed into

$$f_\theta(x) = e^{\theta x - \gamma(\theta)} f(x),$$

where θ is a twisting parameter and γ is the *cumulant-generating function*

$$\gamma(\theta) = \ln \left(\int_{\mathbb{R}} e^{\theta x} dx \right).$$

For a standard Gaussian variable, $\gamma(\theta) = \theta^2/2$, demonstrating that the shift of mean employed above is indeed a special case of exponential twisting. We notice that exponential twisting is often a very convenient starting point when working with parametric families of Radon-Nikodym derivatives.

3.4.4.5 Rare Event Simulation and Linearization

For illustrative purposes, consider finally the problem of estimating by Monte Carlo

$$P(Z > c), \quad (3.98)$$

where $Z \sim \mathcal{N}(0, 1)$ is standard Gaussian under the measure P , and c is a big number. In ordinary Monte Carlo, we write

$$P(Z > c) = E^P(1_{\{Z>c\}})$$

and use the sample mean estimator

$$P(Z > c) \approx \frac{1}{n} \sum_{i=1}^n 1_{\{Z_i > c\}},$$

where Z_1, \dots, Z_n are independent standard Gaussian samples. We notice that

$$\begin{aligned}\text{Var}^P(1_{\{Z>c\}}) &= E^P((1_{\{Z>c\}})^2) - E^P(1_{\{Z>c\}})^2 \\ &= E^P(1_{\{Z>c\}}) - E^P(1_{\{Z>c\}})^2 \\ &= P(Z > c)(1 - P(Z > c)),\end{aligned}$$

with sample mean estimator variance being n times smaller. Consider now introducing a probability measure that shifts the mean of Z from 0 to μ . The likelihood ratio is seen from (3.97) to be

$$l(z) = e^{-\mu z + \mu^2/2},$$

such that

$$P(Z > c) = E^{\widehat{P}}(e^{-\mu Z + \mu^2/2} 1_{\{Z>c\}}),$$

where $Z \sim \mathcal{N}(\mu, 1)$ in the measure \widehat{P} . A Monte Carlo estimator for this is then

$$\frac{1}{n} \sum_{i=1}^n e^{-\mu(Z_i + \mu) + \mu^2/2} 1_{\{Z_i + \mu > c\}},$$

where Z_1, \dots, Z_n are again independent standard Gaussian samples. Notice that we have added to the Z_i the mean μ to reflect the shift of measure from P to \widehat{P} . As for variance, we have

$$\begin{aligned}\text{Var}^{\widehat{P}}(e^{-\mu Z + \mu^2/2} 1_{\{Z>c\}}) &= E^{\widehat{P}}(e^{-2\mu Z + \mu^2} (1_{\{Z>c\}})^2) - P(Z > c)^2 \\ &= E^{\widehat{P}}(e^{-2\mu Z + \mu^2} 1_{\{Z>c\}}) - P(Z > c)^2 \\ &= E^P(e^{-\mu Z + \mu^2/2} 1_{\{Z>c\}}) - P(Z > c)^2 \\ &= e^{\mu^2} P(Z > c + \mu) - P(Z > c)^2,\end{aligned}\tag{3.99}$$

where the last equation follows from the properties of the standard Gaussian density. The choice of μ that minimizes the variance under \widehat{P} is the solution to

$$\min_{\mu} e^{\mu^2} P(Z > c + \mu).$$

Differentiating with respect to μ and setting the resulting expression to zero shows that the variance is minimized at μ^* , where

$$2\mu^* [1 - \Phi(c + \mu^*)] - \phi(c + \mu^*) = 0,\tag{3.100}$$

with Φ and ϕ being the standard Gaussian distribution function and density, respectively. This expression can be solved for μ^* with the aid of a numerical root solver. Alternatively, we can use the fact that c is large to rely on the asymptotic approximation

$$1 - \Phi(c + \mu^*) \approx \frac{\phi(c + \mu^*)}{c + \mu^*},$$

which leads to

$$2\mu^* \frac{\phi(c + \mu^*)}{c + \mu^*} \approx \phi(c + \mu^*) \Rightarrow \mu^* \approx c. \quad (3.101)$$

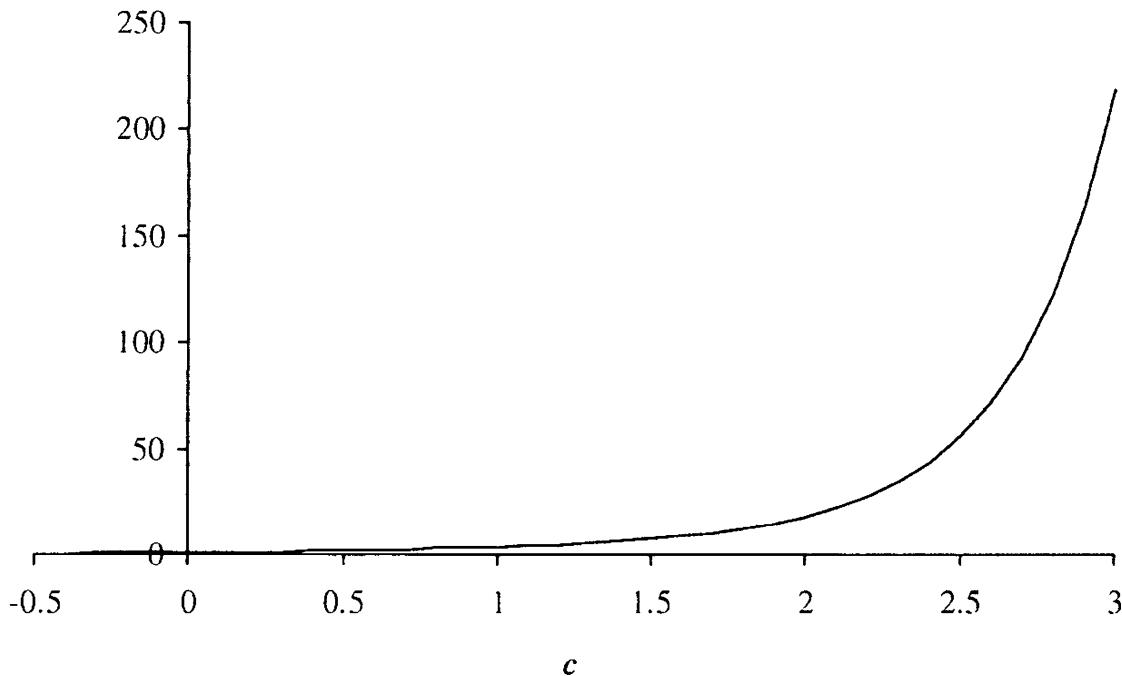
Note that this implies that the probability of Z exceeding c in measure \widehat{P} is approximately $\frac{1}{2}$. This is an intuitive result²⁴, consistent with the discussion at the end of Section 3.4.4.2.

To measure the efficacy of importance sampling, we can use (3.99) to define a variance efficiency ratio as

$$\frac{P(Z > c)(1 - P(Z > c))}{e^{\mu^2} P(Z > c + \mu) - P(Z > c)^2}. \quad (3.102)$$

Figure 3.1 graphs this ratio when μ is set to c , as prescribed in (3.101). For large c , the improvements to variance associated with using importance sampling can be seen to be extremely significant.

Fig. 3.1. Variance Ratio



Notes: The figure graphs the ratio (3.102), with μ set according to (3.101).

It is also illustrative to consider the multi-variate extension to the problem above. Here, we are interested in estimating

$$E^P(1_{\{X>c\}}),$$

²⁴For a somewhat more accurate approximation to μ^* , see Jäckel [2004].

where c is a p -dimensional constant and X is a p -dimensional vector of Gaussian random variables with mean 0 and covariance matrix Σ . Let C be the Cholesky decomposition of Σ , such that

$$\mathbb{E}^P(1_{\{X>c\}}) = \mathbb{E}^P(1_{\{CZ>c\}}) = \mathbb{E}^P(1_{\{Z>c'\}}), \quad c' \triangleq C^{-1}c,$$

where Z is a p -dimensional vector of independent standard Gaussian variables. Let us introduce a measure \widehat{P} where the mean of Z has been shifted to μ , a p -dimensional vector. Following the same steps as for the univariate case, we have

$$\mathbb{E}^P(1_{\{Z>c'\}}) = \mathbb{E}^{\widehat{P}}\left(\exp\left(-\mu^\top Z + \frac{1}{2}\mu^\top \mu\right)1_{\{Z>c'\}}\right),$$

with variance

$$\begin{aligned} \text{Var}^{\widehat{P}}\left(\exp\left(-\mu^\top Z + \frac{1}{2}\mu^\top \mu\right)1_{\{Z>c'\}}\right) \\ = e^{\mu^\top \mu} \mathbb{E}^P(1_{\{Z>c'+\mu\}}) - P(X > c)^2. \end{aligned}$$

A direct optimization of this expression in μ involves multi-dimensional Gaussian integrals, so we wish to resort to approximations. We can use the arguments of Section 3.4.4.2 to argue that the optimal importance sampling density should be proportional to

$$1_{\{z>c'\}} \exp\left(-\frac{1}{2}z^\top z\right), \quad (3.103)$$

since $\exp(-z^\top z/2)$ is proportional to the normal density. Following the idea in Glasserman et al. [1999], we can choose μ such that the location of the peak of an $\mathcal{N}(\mu, I)$ distribution coincides with the peak of (3.103). In other words, we approximate the optimal μ as the value μ^* of z that solves

$$\max_z \left\{ 1_{\{z>c'\}} \exp\left(-\frac{1}{2}z^\top z\right) \right\} = \min_{z>c'} \{z^\top z\}. \quad (3.104)$$

If we assume, say, that all components of c' are larger than 0, then obviously

$$\mu^* = c' = C^{-1}c,$$

consistent with the approximative univariate result (3.101).

We note that the idea behind (3.104) is not limited to situations where we evaluate expectations of an indicator function. For instance, suppose we, as in Section 3.4.4.4, wish to estimate

$$\mathbb{E}^P(G(Z)),$$

for a smooth function $G : \mathbb{R}^p \rightarrow \mathbb{R}$. Restricting ourselves again to the class of measure shifts that only move the mean of Z , the approximately optimal mean shift μ solves

$$\max_z \left\{ G(z) \exp \left(-\frac{1}{2} z^\top z \right) \right\}$$

or, if $G(Z)$ is strictly positive,

$$\max_z \left\{ w(z) - \frac{1}{2} z^\top z \right\}, \quad w(z) \triangleq \ln(G(z)).$$

The first-order condition for the optimum is

$$\nabla w(\mu^*) = (\mu^*)^\top, \quad (3.105)$$

where ∇ is the gradient operator, $\nabla = (\partial/\partial z_1, \dots, \partial/\partial z_p)$ (row vector). This is a fixed-point condition that can be solved by numerical methods. The result (3.105) is exact if w is linear in its argument; the method above can thus be seen as a linearization through a first-order Taylor approximation. Glasserman et al. [1999] demonstrate that, under some conditions, (3.105) satisfies a certain asymptotic optimality property.

3.5 Some Notes on Bermudan Security Pricing

As alluded to in the beginning of this chapter, one drawback of Monte Carlo methods is the difficulty associated with the pricing of securities with early exercise rights. We demonstrated earlier in the chapter that Monte Carlo path generation runs *forward* in time, making direct application of dynamic programming and backward induction (see Chapters 1 and 2) impossible. Indeed, until the early 1990's, it was generally believed that Monte Carlo techniques were inherently incompatible with the pricing of early exercise rights. In the last decade, however, this belief has been overturned, with the advent of several different techniques for Monte Carlo pricing of options with early exercise rights. Most of these techniques are rather advanced and a detailed description will be postponed until later in this book, when the interest rate modeling foundation has been properly laid and the details of callable interest rate securities have been covered. For now, we only provide a brief discussion of certain generic principles, with additional details to be filled in later, in Chapters 18 and 19, among others. We start by establishing some notation and reminding the reader of some basic results from Chapter 1.

3.5.1 Basic Idea

For the remainder of this section, we consider the pricing of a Bermudan security C , with a payout function²⁵ $U(t) = U(t, x(t))$, where $x(t)$ is a

²⁵For many exotic interest rate options, the function $U(t, x(t))$ may actually not be known in closed form. We deal with this complication in Chapter 18.

p -dimensional vector of Markovian state variables. The allowed discrete set of exercise dates is denoted $\mathcal{D} = \{T_1, T_2, \dots, T_B\}$, with $T_B = T$ being the terminal maturity of C . We fix a numeraire N , assumed to be a function of $x(t)$, $N(t) = N(t, x(t))$. From Section 1.10, we recall that

$$C(0) = N(0) \sup_{\tau \in \mathcal{T}} \mathbb{E}^N \left(\frac{U(\tau)}{N(\tau)} \right),$$

where \mathbb{E}^N denotes expectation in the measure Q^N induced by the numeraire N , and \mathcal{T} is the set of stopping time strategies taking values in \mathcal{D} . More generally, we write

$$C(t) = N(t) \sup_{\tau \in \mathcal{T}(t)} \mathbb{E}_t^N \left(\frac{U(\tau)}{N(\tau)} \right), \quad (3.106)$$

where $\mathcal{T}(t)$ is the set of stopping time strategies in \mathcal{D} for which $\tau \geq t$. We also recall that when $t \in (T_{i-1}, T_i]$, we have

$$C(t) = N(t) \mathbb{E}_t^N \left(N(T_i)^{-1} \max(H_i(T_i), U(T_i)) \right), \quad (3.107)$$

where the hold value (see Section 1.10) $H_i(T_i)$ is defined as

$$H_i(T_i) = N(T_i) \mathbb{E}_{T_i}^N \left(N(T_{i+1})^{-1} C(T_{i+1}) \right).$$

Notice that (3.107) establishes that the optimal exercise strategy, as seen from time t , is

$$\tau^* = \inf \{T_i \geq t : U(T_i) \geq H_i(T_i)\}. \quad (3.108)$$

3.5.2 Parametric Lower Bound Methods

Assuming that we are able to simulate the p -dimensional vector $x(t)$ through time, it follows from (3.106) that a lower bound for $C(0)$ can be computed by Monte Carlo method through any exogenous guess for the optimal exercise strategy τ^* . One fairly intuitive approach involves a user-supplied specification of a parametric stopping rule, $\tau(\alpha) \in \mathcal{T}$, where $\alpha \in A \subset \mathbb{R}^m$ is an m -dimensional parameter vector. Defining $x = (x(T_1), \dots, x(T_B))$, for a given value of α , we have the following algorithm.

1. Generate n independent paths $x^{(k)}$, $k = 1, \dots, n$. For path k , let $\tau^{(k)}(\alpha)$ be the exercise date suggested by the parametric stopping rule.
2. For each path k , set $U^{(k)} = U(\tau^{(k)}(\alpha), x(\tau^{(k)}(\alpha)))$ and $C_\alpha^{(k)} = N(\tau^{(k)}(\alpha), x(\tau^{(k)}(\alpha)))^{-1} U^{(k)}$.
3. Return $\bar{C}_\alpha(0) = N(0)n^{-1} \sum_{k=1}^n C_\alpha^{(k)}$ as our estimate for the Bermudan option value $C(0)$.

Let $C_\alpha(0) \triangleq \mathbb{E}^N(\bar{C}_\alpha(0))$. As $\tau(\alpha)$ in general will be sub-optimal, it is clear that

$$C_\alpha(0) \leq C(0). \quad (3.109)$$

To get as close as possible to $C(0)$, it is, of course, preferable to use the value $\alpha^* \in A$ for which $C_\alpha(0)$ is optimized, i.e.

$$\alpha^* = \operatorname{argsup}_{\alpha \in A} C_\alpha(0).$$

A tempting way of estimating $C_{\alpha^*}(0)$ would be to modify step 3 in the algorithm above to

3a. Return $\bar{C}_{\alpha^*}(0) = \sup_{\alpha \in A} \bar{C}_\alpha(0)$.

Leaving aside the question of how one might execute the optimization in Step 3a, we notice that the estimator $\bar{C}_{\alpha^*}(0)$ is biased high relative to $C_{\alpha^*}(0)$:

$$\mathbb{E}^N(\bar{C}_{\alpha^*}(0)) \geq \sup_{\alpha \in A} C_\alpha(0). \quad (3.110)$$

This inequality states that the expected value of the maximum over α must be at least as large as the maximum over α of expected values, a consequence of Jensen's inequality. We may interpret the bias of $\mathbb{E}^N(\bar{C}_{\alpha^*}(0))$ as a *perfect foresight bias*: by using in-sample information to estimate α^* , we effectively "cheat" by making the optimum specific to the same n samples that are also used to determine the option value.

The combination of inequalities (3.109) and (3.110) shows that the quantity $\bar{C}_{\alpha^*}(0)$ from Step 3a has an *indeterminate* bias relative to the true option price. As a bias is generally inevitable when using parametric exercise strategies, in practice it is preferable to at least know its sign. To accomplish this, we can retain the estimated value of α^* found as described above, but draw a *separate* set of Monte Carlo paths when pricing the option. That is, we replace Step 3a with the following two steps:

3b. Set $\hat{\alpha}^* = \operatorname{argsup}_{\alpha \in A} \bar{C}_\alpha(0)$.

4b. Draw a fresh set of n_2 independent paths for x and N , with α locked at the value $\hat{\alpha}^*$. Return $\bar{C}_{\hat{\alpha}^*}(0) = N(0)/n_2 \sum_{k=1}^{n_2} C_{\hat{\alpha}^*}^{(k)}$, where the $C_{\hat{\alpha}^*}^{(k)}$ are computed on the new set of paths.

As the parameter $\hat{\alpha}^*$ will a.s. never equal α^* , it follows that

$$\mathbb{E}^N(\bar{C}_{\hat{\alpha}^*}(0)) \leq C_{\alpha^*}(0) \leq C(0),$$

i.e. we are now assured that $\bar{C}_{\hat{\alpha}^*}(0)$ is low-bound estimator.

3.5.3 Parametric Lower Bound: An Example

What constitutes a good parametric exercise rule is strongly instrument-specific and typically requires case-by-case analysis. Even for simple Markov models and standard option payouts, the topology of exercise and continuation regions can be highly complicated (see e.g. Broadie and Detemple [1997]), so this exercise is by no means straightforward. As a first approximation, however, one can always attempt to use a simple rule based on outright “moneyness” of the underlying option payout, as in Andersen [2000a]. According to this rule, one sets $\alpha = (h_1, h_2, \dots, h_B)$ (i.e. $m = B$) and writes

$$\tau(\alpha) = \inf \{T_i : U(T_i, x(T_i)) > h_i\}. \quad (3.111)$$

That is, exercise of the option takes place when it is sufficiently deep in the money, with the term “sufficiently deep” quantified through unknown trigger thresholds $h_i \geq 0$, $i = 1, \dots, B$.

While α is B -dimensional, finding its optimal value is not truly a B -dimensional optimization problem. Rather, due to the Markov assumption on $x(t)$, we may decompose it into a series of $B - 1$ *one-dimensional* optimization problems. Specifically, working backwards in time, suppose that the optimal values of $h_{j+1}, h_{j+2}, \dots, h_B$ are known. We then find the optimal value of h_j , by optimizing on $\bar{C}_\alpha(0)$, but subject to the constraint that exercise is *not* allowed to take place before time j . As $h_{j+1}, h_{j+2}, \dots, h_B$ are assumed known — and h_1, \dots, h_{j-1} do not come into play — the only variable involved in this optimization is h_j . The algorithm starts with $j = B - 1$ and the known²⁶ boundary condition $h_B = 0$.

A couple of comments on the algorithm above are in order. First, we notice that $U(T_i, x(T_i)) > h_i$ can be replaced with any one-parameter boolean function $g(T_i, x(T_i); h_i)$ without affecting the basic algorithm — see Andersen [2000a] for some examples. Second, if in such a boolean function $g(T_i, x(T_i); h_i)$ each h_i is allowed to be q -dimensional, the optimization problem reduces to $B - 1$ q -dimensional optimization problems. And third, for a finite-path simulation, the objective functions in each of the $B - 1$ optimization problems will not be smooth; consequently, the optimization is best performed by an iterative search rather than a derivative-based method. Andersen [2000a] uses the golden section search (see Press et al. [1992]), but simpler strategies based on, say, outright sorting are also possible.

3.5.4 Regression-Based Lower Bound

According to (3.108), an approximation for the optimal exercise strategy can always be constructed through an estimate for hold values $H_i(T_i)$ at all $i = 1, 2, \dots, B - 1$. In our Markov setting, we know that

²⁶At the last possible exercise date, we would, of course, always exercise the option if it is in-the-money.

$$H_i(T_i) = q_i(x(T_i)) = q_i(x_1(T_i), \dots, x_p(T_i))$$

for a set of $B - 1$ functions $q_i : \mathbb{R}^p \rightarrow \mathbb{R}$, $i = 1, 2, \dots, B - 1$; the problem of estimating hold values is equivalent to the problem of estimating the functions q_i .

From Section 3.5.1, we know that

$$q_i(x) = N(T_i, x) E^N(C(T_{i+1})/N(T_{i+1}, x(T_{i+1})) | x(T_i) = x), \quad (3.112)$$

which can be interpreted as the *regression* of $C(T_{i+1})/N(T_{i+1})$ on the Markov state variables $x(T_i)$. Several authors — including Carriére [1996], Longstaff and Schwartz [2001], and Tsitsiklis and Roy [2001] — have used this observation to suggest that $q_i(x)$ be estimated by a linear combination of exogenously specified (basis) functions of $x(T_i)$, with least-squares regression on Monte Carlo paths used to determine the best weights for these functions. That is, we fundamentally assume that

$$q_i(x) = \sum_{j=1}^d \beta_{i,j} \psi_j(x), \quad (3.113)$$

for a set of d basis-functions $\psi_j : \mathbb{R}^p \rightarrow \mathbb{R}$, $j = 1, 2, \dots, d$. Setting $\beta_i = (\beta_{i,1}, \dots, \beta_{i,d})^\top$ and $\psi(x) = (\psi_1(x), \dots, \psi_d(x))^\top$, we can rewrite (3.113) as $q_i(x) = \psi(x)^\top \beta_i$ or, from (3.112),

$$\begin{aligned} E^N(N(T_i, x(T_i)) C(T_{i+1})/N(T_{i+1}, x(T_{i+1})) | x(T_i) = x) \\ = E^N(\psi(x(T_i))^\top | x(T_i) = x) \beta_i. \end{aligned}$$

This, in turn, implies that

$$\Omega_i = \Psi_i \beta_i \Rightarrow \beta_i = \Psi_i^{-1} \Omega_i, \quad (3.114)$$

where Ω_i is the d -dimensional vector

$$\Omega_i = E^N\left(\psi(x(T_i)) \frac{N(T_i, x(T_i))}{N(T_{i+1}, x(T_{i+1}))} C(T_{i+1})\right),$$

and Ψ_i is the $d \times d$ matrix

$$\Psi_i = E^N\left(\psi(x(T_i)) \psi(x(T_i))^\top\right).$$

The rationale for rewriting (3.113) into the seemingly more convoluted representation (3.114) is that the latter leads naturally to the algorithm for a least-squares estimation of β_i : one simply replaces the expectations in Ψ_i and Ω_i with sample averages $\bar{\Psi}_i$ and $\bar{\Omega}_i$ computed on a set of Monte Carlo paths. That is, one uses²⁷

²⁷In practice, a direct solution of linear equations in this fashion can be suboptimal if the matrix $\bar{\Psi}_i$ is ill-conditioned. Instead, one would use either truncated singular value decomposition (TSVD) or Tikhonov regularization to find $\hat{\beta}_i$. We return to this issue in Chapter 18.

$$\widehat{\beta}_i = \overline{\Psi}_i^{-1} \overline{\Omega}_i \quad (3.115)$$

as the sample estimate.

We shall discuss the details of (and many variations on) the regression approach later, in Chapter 18. For now, let us just notice that computation of $\overline{\Omega}_i$ requires estimation of $C(T_{i+1})$, which naturally encourages running the estimation of the $\widehat{\beta}_i$ backwards in i , starting from $i = B - 1$. We also notice that the success of the regression approach depends critically on the choice and number of basis functions ψ_j . We give specific advice on this topic in Chapters 18 and 19.

3.5.5 Upper Bound Methods

Given a martingale M in measure Q^N , we recall from Section 1.10.2 that an M -specific upper bound $C_M(0)$ for a Bermudan option can always be constructed as

$$C_M(0) = N(0) \left\{ M(0) + E^N \left(\max_{t \in \mathcal{D}} \left(\frac{U(t, x(t))}{N(t, x(t))} - M(t) \right) \right) \right\} \geq C(0). \quad (3.116)$$

Let $\mathbf{M} = (M(T_1), \dots, M(T_B))^T$. As long as $M(\cdot)$ can be simulated along with the vector $x(\cdot)$, (3.116) suggests the following Monte Carlo algorithm:

1. Generate n independent paths $x^{(k)}$, $\mathbf{M}^{(k)}$, $k = 1, \dots, n$. For path k , let $\gamma^{(k)}$ be the maximum value of $U(T_i, x(T_i))/N(T_i, x(T_i)) - M(T_i)$ over $i = 1, 2, \dots, B$.
2. Return $\overline{C}_M(0) = N(0)\{M(0) + n^{-1} \sum_{k=1}^n \gamma^{(k)}\}$.

For the upper bound method to be practically useful, we would want the gap $E^N(\overline{C}_M(0)) - C(0)$ to be small. This, in turn, requires that we specify the martingale $M(t)$ to be “reasonable”. More specifically, from results in Section 1.10.2, we would like $M(t)$ to represent the martingale component of a good approximation to the supermartingale $C(t)$. For instance, if we happen to think that $C(t)$ is well-approximated by some known function v of time and x ,

$$C(t) = v(t, x(t)), \quad (3.117)$$

we would set

$$M(t) = \sum_{j=1}^p \frac{\partial v(t, x(t))}{\partial x_j} (dx_j(t) - E_t^N(dx_j(t))). \quad (3.118)$$

As an example, suppose that W^N is a (possibly vector-valued) Brownian motion in Q^N and $dx_j(t) = \mu_j(t, x(t)) dt + \sigma_j(t, x(t)) dW^N(t)$, in which case we could use

$$M(t) = \sum_{j=1}^p \frac{\partial v(t, x(t))}{\partial x_j} \sigma_j(t, x(t)) dW^N(t). \quad (3.119)$$

Occasionally, a natural analytical guess for the function v may exist, but most often the only estimate for v is given only implicitly, through a low-bound estimator based on an approximation of the optimal exercise strategy. A completely generic algorithm to turn a guess for the optimal exercise strategy into a proxy for $M(t)$ is developed in Andersen and Broadie [2004]; we shall discuss this algorithm in detail in Chapter 18. If the evolution of $x(t)$ is described by an SDE, a regression approach to estimation of the terms multiplying $dW^N(t)$ in (3.119) can be found in Belomestny et al. [2007].

3.5.6 Confidence Intervals

Suppose that we simultaneously apply a lower bound and an upper bound method to provide two sample estimates \bar{C}_{lo} and \bar{C}_{up} , with

$$\mathbb{E}^N(\bar{C}_{\text{lo}}) \leq C(0) \leq \mathbb{E}^N(\bar{C}_{\text{up}}). \quad (3.120)$$

Let us assume that the sample standard errors on \bar{C}_{lo} and \bar{C}_{up} have been computed as s_{lo} and s_{up} , respectively. For a sufficiently large number of Monte Carlo trials, we can then use the central limit theorem from Section 3.1 to set up a confidence interval

$$[\bar{C}_{\text{lo}} - u_{\gamma/2} \cdot s_{\text{lo}}, \bar{C}_{\text{up}} + u_{\gamma/2} \cdot s_{\text{up}}],$$

where $\Phi(u_{\gamma/2}) = 1 - \gamma/2$. It is clear from (3.120) that the likelihood of this interval bracketing the true price $C(0)$ is *at least* $1 - \gamma$. It is also clear that this confidence interval will not shrink to zero — even in the limit of an infinite number of samples where $s_{\text{lo}} \rightarrow 0$ and $s_{\text{up}} \rightarrow 0$ — unless \bar{C}_{lo} and \bar{C}_{up} simultaneously achieve the unlikely feat of being perfectly unbiased estimators for $C(0)$.

Finally, let us note that any number inside the interval $[\bar{C}_{\text{lo}}, \bar{C}_{\text{up}}]$ can reasonably be used as an estimator for $C(0)$. To the extent that we have reason to believe²⁸ that \bar{C}_{lo} and \bar{C}_{up} have roughly opposite biases, a natural estimator is $(\bar{C}_{\text{lo}} + \bar{C}_{\text{up}})/2$.

3.5.7 Other Methods

The methods described so far are those that, in our opinion, are most useful for practical Monte Carlo pricing of interest rate options with early exercise rights. Several other methods, however, have been proposed in the literature, some of which have interesting theoretical properties. We highlight the *random tree* method (Broadie and Glasserman [1997]) which builds a random

²⁸There is some evidence that the upper-bound method in Andersen and Broadie [2004] produces a bias that is often roughly opposite of that of the low-bound method from which the martingale M in (3.116) is extracted.

non-recombining lattice by Monte Carlo methods; backward induction arguments are then used to construct high- and low-biased estimators, both of which are convergent to the true price as the number of Monte Carlo paths are increased. The drawback of the method is its computational complexity which increases exponentially in the number of exercise dates (B), ruling out its practical usage for many realistic applications. Broadie and Glasserman [2004] suggest a recombining *stochastic mesh* method that grows only linearly in the number of exercise weights; in its basic form, this method requires explicit knowledge of transition densities as it relies on likelihood ratios to set weights on nodes in the mesh (see Section 3.3.3). As discussed in Glasserman [2004], the concept of stochastic meshes can, however, be broadened to include several other methods, include the regression approach in Section 3.5.4.

3.A Appendix: Constants for Φ^{-1} Algorithm

a_0	2.50662823884	c_0	0.3374754822726147
a_1	-18.61500062529	c_1	0.9761690190917186
a_2	41.39119773534	c_2	0.1607979714918209
a_3	-25.44106049637	c_3	0.0276438810333863
b_0	-8.47351093090	c_4	0.0038405729373609
b_1	23.08336743743	c_5	0.0003951896511919
b_2	-21.06224101826	c_6	0.0000321767881768
b_3	3.13082909833	c_7	0.0000002888167364
		c_8	0.0000003960315187

Fundamentals of Interest Rate Modeling

The purpose of this brief chapter is twofold. First, we introduce notations to characterize prices and yields of basic fixed income market securities. In addition to providing the foundation for a more expansive discussion of fixed income markets (which we shall undertake in Chapter 5), this part of the chapter serves to identify and characterize a number of probability measures that are of fundamental importance in models for the term structure of interest rates. A brief discussion of measures used in a two-currency setting is also provided.

In the second part of the chapter, we discuss general characteristics of models with dynamics driven by vector-valued Brownian motions. This analysis leads to the fundamental class of Heath-Jarrow-Morton (HJM) (see Heath et al. [1992]) models of continuously compounded forward rates. Among other special cases, we discuss in some detail tractable HJM models with Gaussian volatility structure, and provide some results for the case where such models are Markovian. These discussions continue in Chapters 10 through 12 where we consider one- and multi-factor short rate models, and in the Chapter 13 where we introduce the important class of quasi-Gaussian HJM models with local and stochastic volatility.

4.1 Fixed Income Notations

4.1.1 Bonds and Forward Rates

As in earlier chapters, let $P(t, T)$ denote the time t price of a zero-coupon bond (also known as a *discount bond*) delivering for certain \$1 at maturity $T \geq t$. Suppose we are interested in purchasing at some future time T a zero-coupon bond maturing at $T + \tau$, $\tau > 0$. At time $t < T$, the price of such a bond can be locked in by i) purchasing at time t one $(T + \tau)$ -maturity zero-coupon bond; and ii) selling short (“shorting”) $P(t, T + \tau)/P(t, T)$

T -maturity zero-coupon bonds. The time t cost of executing this strategy is zero,

$$-1 \cdot P(t, T + \tau) + P(t, T + \tau)/P(t, T) \cdot P(t, T) = 0,$$

but a flow of

$$-P(t, T + \tau)/P(t, T)$$

will take place at time T as the T -maturity short position matures. This is compensated by an inflow of \$1 at time $T + \tau$. In other words, our trading strategy effectively fixes the time T purchase price of the $(T + \tau)$ -maturity bond at

$$P(t, T, T + \tau) \triangleq P(t, T + \tau)/P(t, T), \quad \tau > 0,$$

a quantity known as the time t *forward price* for the zero-coupon bond spanning $[T, T + \tau]$.

It is often convenient to characterize a forward bond price by a discount rate. One such rate is the *continuously compounded forward yield* $y(t, T, T + \tau)$, defined by

$$e^{-y(t, T, T + \tau)\tau} = P(t, T, T + \tau). \quad (4.1)$$

The time between the maturity of the forward bond and the expiry of the forward contract, i.e. τ , is often called the *tenor* of the forward bond or the forward yield. In the definition of the continuously compounded yield lies an implicit, and idealized, assumption of continuous reinvestment of investment proceeds. Most actual market quotes, however, are based on discrete-time compounding of proceeds. Accordingly, we define a *simple forward rate* $L(t, T, T + \tau)$ as

$$1 + \tau L(t, T, T + \tau) = 1/P(t, T, T + \tau). \quad (4.2)$$

Again, τ is the *tenor* of the forward rate. For an arbitrary set of dates $T = T_0 < T_1 < T_2 < \dots < T_n$, notice that forward bond prices can be recovered from forward rates by simple compounding,

$$P(t, T_n)/P(t, T) = \prod_{i=1}^n \frac{1}{1 + (T_i - T_{i-1}) L(t, T_{i-1}, T_i)}.$$

Unless we state otherwise, throughout this book we shall typically make the assumption that spot rates $L(T, T, T + \tau)$ are the *Libor (London Interbank Offered Rate) rates* quoted in the interbank market. Libor rates are quoted on values of τ ranging from one week ($\tau = 1/52$)¹ to 12 months ($\tau = 1$), and form the basis for a number of floating-rate derivative contracts, such as interest rate swaps and Eurodollar futures. We shall examine these securities in more detail in Chapter 5.

¹Note that in reality the calculation of year fractions τ are governed by fairly complicated market conventions. A brief discussion of this topic can be found in Appendix 5.A.

In the limit $\tau \downarrow 0$,

$$L(t, T, T + \tau) \rightarrow f(t, T),$$

where the quantity $f(t, T)$ is the time t *instantaneous forward rate* to time T . We think of $f(t, T)$ as the forward rate spanning $[T, T + dT]$, observed at time t . The relation between instantaneous forward rates and bond prices is given by the continuous compounding formula

$$P(t, T, T + \tau) = \exp \left(- \int_T^{T+\tau} f(t, u) du \right), \quad (4.3)$$

such that

$$f(t, T) = - \frac{\partial \ln P(t, T)}{\partial T}, \quad (4.4)$$

and, from (4.1),

$$y(t, T, T + \tau) = \tau^{-1} \int_T^{T+\tau} f(t, u) du, \quad f(t, T) = \lim_{\tau \downarrow 0} y(t, T, T + \tau).$$

We also notice the relationship

$$f(t, T) = \frac{\partial (y(t, t, T)(T - t))}{\partial T} = y(t, t, T) + (T - t) \frac{\partial y(t, t, T)}{\partial T}.$$

The quantity

$$r(t) \triangleq f(t, t) \quad (4.5)$$

is an \mathcal{F}_t -measurable random variable known as the *short rate* or sometimes the *spot rate*. Loosely speaking, we can think of $r(t)$ as the overnight rate in effect at time t .

Finally, let us note that interest rates of various flavors and related quantities are typically quoted in percentage points or sometimes in *basis points*, where 1 basis point = 1/100 of one percent.

4.1.2 Futures Rates

Through the market for Eurodollar futures (see Chapter 5), investors can enter into securities that will pay at time T an amount of

$$1 - L(T, T, T + \tau). \quad (4.6)$$

At time 0, a Eurodollar futures contract can be entered into at no upfront cost, but with an implicit obligation of the holder to pay at time T per unit of notional

$$1 - F(0, T, T + \tau)$$

in return for the payout (4.6). Here, $F(t, T, T + \tau)$ is the time t *simple futures rate* for the period $[T, T + \tau]$. Importantly, the futures rate is *marked to market* (or *resettled*) each day, with the day's change in the futures rate immediately credited to or debited from the contract holder's account with the futures exchange. Specifically, after holding the contract for a period of $\Delta = 1$ day, the futures contract holder would thus experience a cash flow of

$$(1 - F(\Delta, T, T + \tau)) - (1 - F(0, T, T + \tau)) \\ = - (F(\Delta, T, T + \tau) - F(0, T, T + \tau)).$$

Continuing the mark-to-market process to maturity shows that the total amount of cash flow received by the holder on $[0, T]$ is

$$- (F(T, T, T + \tau) - F(0, T, T + \tau)) = - (L(T, T, T + \tau) - F(0, T, T + \tau)) \quad (4.7)$$

where we have used the fact that $F(T, T, T + \tau)$ must equal $L(T, T, T + \tau)$ to avoid a *delivery arbitrage*.

The fact that the net cash flow payment (4.7) on a Eurodollar futures contract has been made incrementally on a daily basis has important valuation consequences, and causes the futures rate to differ from the forward rate defined earlier. For instance, under a scenario of rising interest rates, the holder of a Eurodollar futures contract must make payments to the futures exchange. As rates are rising, the contract holder will be faced with a high-rate — and thus unfavorable — borrowing environment for funding these payments. Conversely, when interest rates fall, the reinvestment of received funds will take place at increasingly low rates. Due to the adverse behavior of funding costs and reinvestment gains, we would expect the purchaser of a Eurodollar futures contract to pay less for these instruments than for a comparable instrument without daily mark-to-market. Consequently, we would expect the futures rate to be *above* the corresponding forward rate. We shall quantify this effect in Section 4.5.1 and, with more advanced models, in Section 16.8.

We notice that we can define *instantaneous futures rates* $q(t, T)$ in the same fashion as we defined instantaneous forward rates:

$$q(t, T) = \lim_{\tau \downarrow 0} F(t, T, T + \tau).$$

4.1.3 Annuity Factors and Par Rates

Most fixed income securities involve multiple cash flows taking place on a pre-set schedule of dates, often referred to as a *tenor structure*,

$$0 \leq T_0 < T_1 < \dots < T_N.$$

Given a tenor structure, for any two integers k, m satisfying $0 \leq k < N$, $m > 0$, and $k + m \leq N$, we can define an *annuity factor* $A_{k,m}$ by

$$A_{k,m}(t) = \sum_{n=k}^{k+m-1} P(t, T_{n+1}) \tau_n, \quad \tau_n = T_{n+1} - T_n. \quad (4.8)$$

Annuity factors provide for compact notation when pricing coupon-bearing securities. For instance, a security making m coupon payments of $c\tau_n$ at all $T_{n+1}, n = k, \dots, k + m - 1$, is easily seen to have time t value of

$$cA_{k,m}(t), \quad t \leq T_k.$$

If the security also involves a back-end return of notional at time T_{k+m} (as is the case for a regular coupon-bearing bond), the pricing expression is

$$cA_{k,m}(t) + P(t, T_{k+m}), \quad (4.9)$$

where we assume that the bond has been normalized to have a unit notional. The time t forward price to T_k of the security (4.9) is

$$cA_{k,m}(t)/P(t, T_k) + P(t, T_{k+m})/P(t, T_k);$$

the value of the coupon c for which this expression equals 1 is known as the *forward par rate* or, when used in the context of swap pricing, as the *forward swap rate*. With $S_{k,m}(t)$ denoting the time t swap rate, we apparently have

$$S_{k,m}(t) = \frac{P(t, T_k) - P(t, T_{k+m})}{A_{k,m}(t)}, \quad t \leq T_k. \quad (4.10)$$

From the definition of $L(t, T_n, T_{n+1})$ in (4.2), a little thought shows that the numerator of the expression for $S_{k,m}(t)$ can be expanded into a weighted sum of forward rates, leading to the alternative expression

$$S_{k,m}(t) = \frac{\sum_{n=k}^{k+m-1} \tau_n P(t, T_{n+1}) L_n(t)}{A_{k,m}(t)}, \quad t \leq T_k, \quad (4.11)$$

where we have introduced the useful shorthand

$$L_n(t) \triangleq L(t, T_n, T_{n+1}).$$

It follows that the forward swap rate can be loosely interpreted as a weighted average of simple forward rates on the specified tenor structure. We note for the future that the time T_k is sometimes referred to as the *fixing date*, or *expiry*, of the swap rate $S_{k,m}$, while the length of the corresponding swap, $T_{k+m} - T_k$, is sometimes called the *tenor* of the swap rate.

4.2 Fixed Income Probability Measures

As discussed in Chapter 1, selecting an equivalent martingale measure is largely a matter of choosing a *numeraire*, an asset price process used to

re-normalize the prices of other traded securities. For later reference, this section lists and names a number of important numeraires and measures used in fixed income pricing. Throughout the section, we assume that the market is complete, and we use $V(t)$ to denote the time t price of a derivative security making an \mathcal{F}_T -measurable payment of $V(T)$.

4.2.1 Risk Neutral Measure

The numeraire defining the risk-neutral measure Q is the continuously compounded money market account $\beta(t)$, satisfying the locally deterministic SDE

$$d\beta(t) = r(t)\beta(t) dt, \quad \beta(0) = 1, \quad (4.12)$$

where $r(t)$ is the short rate, $r(t) = f(t, t)$. Solving this equation yields

$$\beta(t) = e^{\int_0^t r(u) du}.$$

From the results of Chapter 1, in the absence of arbitrage the numeraire-deflated process $V(t)/\beta(t)$ must be a martingale, implying the derivative security valuation formula

$$V(t)/\beta(t) = E_t^Q (V(T)/\beta(T)), \quad t \leq T, \quad (4.13)$$

or equivalently

$$V(t) = E_t^Q \left(e^{-\int_t^T r(u) du} V(T) \right). \quad (4.14)$$

If we apply (4.14) to the special case of $V(T) = 1$, we obtain a fundamental bond pricing formula. We highlight the importance of this result by listing it in a lemma.

Lemma 4.2.1. *In the absence of arbitrage, the time t price $P(t, T)$ of a T -maturity zero-coupon bond is*

$$P(t, T) = E_t^Q \left(e^{-\int_t^T r(u) du} \right). \quad (4.15)$$

It follows from Lemma 4.2.1 that specification of the dynamics of $r(t)$ under Q suffices to determine the prices of discount bonds at all times and maturities. Models that are based on such a direct specification of $r(t)$ dynamics are known as *short rate models* and are the subject of Chapters 10 through 12. Notice the resemblance between expressions (4.3) and (4.15); if $r(t)$ is deterministic, the two expressions will agree as $r(u) = f(t, u)$, $u \geq t$. If $r(t)$ is random, one may wonder whether this result will hold in expectation. The answer to this is negative, i.e.

$$f(t, u) \neq E_t^Q (r(u)), \quad (4.16)$$

provided r is random. We prove this in the section below. Under certain idealized conditions, however, equality holds in (4.16) provided $f(t, u)$ is replaced by the *futures* rate $q(t, u)$. The exact result is as follows.

Lemma 4.2.2. Assume that mark-to-market takes place continuously. Under regularity conditions on the short rate $r(\cdot)$ — it suffices that $r(\cdot)$ is positive and bounded — the futures rate $F(\cdot, T, T + \tau)$ is a Q-martingale, and

$$F(t, T, T + \tau) = E_t^Q(L(T, T, T + \tau)). \quad (4.17)$$

Proof. Over a small interval $[t, t + dt]$, we have earlier shown that the cash proceeds from a futures contract are proportional to

$$dF(t, T, T + \tau) = F(t + dt, T, T + \tau) - F(t, T, T + \tau).$$

Suppose that we hold the futures contract up to some arbitrary horizon $t < \mathcal{T} \leq T$ at which point we exit (e.g., by selling the futures contract). Deflating all cash proceeds from this strategy with the numeraire $\beta(t)$ and integrating provides us with the time t value of the futures contract as

$$V_{\text{fut}}(t) = \beta(t) E_t^Q \left(\int_t^{\mathcal{T}} \beta(s)^{-1} dF(s, T, T + \tau) + \beta(\mathcal{T})^{-1} V_{\text{fut}}(\mathcal{T}) \right).$$

As it is always costless to enter into a futures contract, $V_{\text{fut}}(t) = V_{\text{fut}}(\mathcal{T}) = 0$ by definition, so for arbitrary $t < \mathcal{T} \leq T$ we must have (since $\beta(t)$ is positive)

$$E_t^Q \left(\int_t^{\mathcal{T}} \beta(s)^{-1} dF(s, T, T + \tau) \right) = 0. \quad (4.18)$$

Provided that $\beta(s)^{-1}$ is almost surely positive (which is the case if r is bounded), the fact that (4.18) holds for arbitrary horizons $t < \mathcal{T} \leq T$ shows that

$$E_t^Q(dF(s, T, T + \tau)) = 0, \quad t \leq s \leq T,$$

which demonstrates that F is a Q-martingale. The result (4.17) then immediately follows. \square

Equation (4.17) states that the futures rate is the Q-expectation of the time T spot rate $L(T, T, T + \tau)$. A similar relation must then hold for the instantaneous futures rate, i.e.

$$q(t, u) = E_t^Q(r(u)) \quad (4.19)$$

as stated earlier².

Lemma 4.2.2 was first proven by non-probabilistic methods in Cox et al. [1981], who employed a direct, and quite instructive, hedging argument to

²This result should not be confused with the classical *expectations hypothesis* which states that futures (or sometimes forward) rates are unbiased estimators of future spot rates, in the *real-life* probability measure P: $E_t^P(r(T)) = q(t, T)$. The expectations hypothesis amounts to a strong assumption about the market price of risk (see Chapter 1), whereas equation (4.19) is a preference-free arbitrage relationship.

show the result. The assumption of continuous resettlement in the lemma may appear idealized, but the difference between daily and continuous settlement is quite small, as shall be demonstrated in Chapter 16. Explicit modeling of discrete resettlement is nevertheless quite straightforward, and basically involves shifting measure, from the risk-neutral measure to the so-called spot measure, defined below. We return to this issue in Chapter 16.

4.2.2 T -Forward Measure

The T -forward measure Q^T was introduced in Jamshidian [1991b] (see also Geman et al. [1995]), and uses a T -maturity zero-coupon bond as the numeraire asset. As is customary, we let $E^T(\cdot)$ denote expectations in measure Q^T , such that

$$V(t)/P(t, T) = E_t^T(V(T)/P(T, T)), \quad t \leq T.$$

As obviously $P(T, T) = 1$, this expression simplifies to the convenient form

$$V(t) = P(t, T)E_t^T(V(T)). \quad (4.20)$$

Comparison of (4.20) and (4.14) shows that shifting to the T -forward measure in a sense decouples the expectation of the terminal payout $V(T)$ from that of the numeraire. As we shall see, this is often very convenient when we attempt to construct analytical formulas for prices of certain simple interest rate derivatives. From the results of Section 1.3, we note that the explicit connection between the risk-neutral and T -forward measures is given by the density

$$E_t^Q \left(\frac{dQ^T}{dQ} \right) = \frac{P(t, T)/P(0, T)}{\beta(t)}. \quad (4.21)$$

As $P(t, T + \tau)$ is the price of a traded asset, from the definition of the T -forward measure it follows that forward bond prices

$$P(t, T, T + \tau) = P(t, T + \tau)/P(t, T)$$

are martingales in the T -forward measure. We highlight a related result for forward rates below.

Lemma 4.2.3. *In the absence of arbitrage the forward Libor rate $L(t, T, T + \tau)$ is a martingale under $Q^{T+\tau}$, such that*

$$L(t, T, T + \tau) = E_t^{T+\tau}(L(T, T, T + \tau)), \quad t \leq T. \quad (4.22)$$

Proof. By definition (see (4.2))

$$L(t, T, T + \tau) = \tau^{-1} (P(t, T)/P(t, T + \tau) - 1).$$

As $P(t, T)/P(t, T + \tau)$ is a martingale under $Q^{T+\tau}$, so is $L(t, T, T + \tau)$. The result follows. \square

Taking the limit $\tau \downarrow 0$ and setting $T = u$ yields

$$f(t, u) = E_t^u(f(u, u)) = E_t^u(r(u)), \quad (4.23)$$

which should be compared to the result (4.19).

4.2.3 Spot Measure

When working with a multitude of forward rates on a tenor structure $0 = T_0 < T_1 < \dots < T_N$, it is often convenient to introduce a numeraire that can be extended to arbitrary horizons by compounding. While the continuously compounded money market account β would accomplish this, working with a continuously compounded numeraire is inherently awkward in a setting with a discrete tenor structure. As an alternative, we can introduce a discrete-time equivalent of the continuously compounded money market account to be the value of the following trading strategy. At time 0, \$1 is invested in $1/P(0, T_1)$ T_1 -maturity discount bonds, returning the amount

$$1/P(0, T_1) = 1 + \tau_0 L(0, 0, \tau_0)$$

at time T_1 . This amount is then reinvested (“rolled”) at time T_1 in T_2 -maturity bonds, returning

$$1/P(0, T_1) \cdot 1/P(T_1, T_2) = (1 + \tau_0 L(0, 0, T_1))(1 + \tau_1 L(T_1, T_1, T_2))$$

at time T_2 . Repeating this re-investment strategy at each date in the tenor structure gives rise to an asset price process $B(t)$, where $B(0) = 1$ and

$$B(t) = \prod_{n=0}^i (1 + \tau_n L_n(T_n)) P(t, T_{i+1}), \quad T_i < t \leq T_{i+1}, \quad (4.24)$$

where we used the already introduced short-hand notation $L_n(t) = L(t, T_n, T_{n+1})$. The process $B(t)$ is effectively a rolling certificate of deposit, and can be interpreted as a discrete-time equivalent of $\beta(t)$. $B(t)$ will approach $\beta(t)$ as the time spacing of the tenor structure is made increasingly fine.

The measure induced by $B(t)$ is known as the *spot measure* (or sometimes *spot Libor measure*), denoted Q^B . With $E^B(\cdot)$ denoting expectations in this measure we have

$$V(t) = E_t^B \left(V(T) \frac{B(t)}{B(T)} \right),$$

where

$$\begin{aligned} \frac{B(t)}{B(T)} &= \prod_{n=i+1}^j (1 + \tau_n L_n(T_n))^{-1} \frac{P(t, T_{i+1})}{P(T, T_{j+1})}, \\ T_i < t \leq T_{i+1}, \quad T_j < T \leq T_{j+1}. \end{aligned}$$

The similarity between the discrete and continuous money market accounts makes the spot Libor measure resemble the risk-neutral measure in many ways. For example, as we recall from Lemma 4.2.2, the risk-neutral measure is characterized by the fact that a continuously resettled futures rate is a martingale. In close parallel, the futures rate that is marked to market (resettled) discretely on dates T_0, \dots, T_N turns out to be a martingale in the spot Libor measure. We show this in Section 16.8.

4.2.4 Terminal and Hybrid Measures

One advantage of the spot measure over an arbitrary T -forward measure is the fact that the numeraire asset $B(t)$ will remain alive throughout the span of the tenor structure $\{T_n\}_{n=0}^N$. This property of $B(t)$ is necessary for the valuation of securities which may mature (randomly) at any date in the tenor structure. Securities of this type include for instance barrier options (such as range accruals) and options with early exercise rights (such as Bermudan swaptions). On the other hand, if we pick the T -forward measure corresponding to the *last* maturity in the tenor structure, $T = T_N$, this also yields a numeraire asset — the T_N -maturity zero-coupon bond — that is certain to remain alive at all dates in the tenor structure. The measure induced by $P(t, T_N)$ (Q^{T_N}) is often referred to as the *terminal measure*. For a security V maturing at a date $T \leq T_N$ we get, from the usual martingale pricing formula,

$$V(t) = P(t, T_N) E_t^{T_N} (V(T)/P(T, T_N)), \quad t \leq T \leq T_N. \quad (4.25)$$

In (4.25) it is useful to notice that $V(T)/P(T, T_N)$ is the time T_N proceeds of rolling at time T the security payout $V(T)$ into a zero-coupon bond maturing at time T_N , effectively aligning the maturity of the numeraire and the cash flow date of the underlying asset. As an alternative, $V(T)$ could be rolled into the spot numeraire asset $B(T)$, leading to a T_N payout of $V(T)/B(T) \cdot B(T_N)$. This gives rise to the equivalent formula

$$V(t) = P(t, T_N) E_t^{T_N} (V(T)B(T_N)/B(T)), \quad t \leq T \leq T_N. \quad (4.26)$$

We note that this formula can also be derived from the basic relationship between the measures Q^B and Q^{T_N} by simply noting that

$$P(T, T_N) E_T^{T_N} (B(T_N)/B(T)) = B(T) E_T^B (B(T_N)/B(T)/B(T_N)) = 1,$$

such that, by iterated conditional expectations³,

³The *law of iterated conditional expectations*, sometimes known as the *tower rule*, states that for an \mathcal{F}_T -measurable random variable X , $E(E(X|\mathcal{F}_s)|\mathcal{F}_t) = E(X|\mathcal{F}_t)$, where $t \leq s \leq T$.

$$\begin{aligned}
V(t) &= P(t, T_N) \mathbb{E}_t^{T_N} (V(T)/P(T, T_N)) \\
&= P(t, T_N) \mathbb{E}_t^{T_N} \left(V(T) \mathbb{E}_T^{T_N} (B(T_N)/B(T)) \right) \\
&= P(t, T_N) \mathbb{E}_t^{T_N} (V(T) B(T_N)/B(T)),
\end{aligned}$$

as before.

As mentioned, equations (4.25) and (4.26) effectively involve reinvestment of the proceeds $V(T)$ to align cash payment with the numeraire $P(t, T_N)$. If the numeraire expires before the derivative security, we can apply the same reinvestment idea, but this time to the numeraire asset. Consider for instance a derivative security maturing at time T_N (paying $V(T_N)$), and suppose we wish to extend the T -forward measure to price this option. For instance, we can define a numeraire asset as follows:

$$\tilde{P}(t, T) = \begin{cases} P(t, T), & t \leq T, \\ B(t)/B(T), & t > T. \end{cases}$$

This asset corresponds to an investment strategy where we i) at time 0 purchase the T -maturity zero-coupon bond; and ii) at time T invest the proceeds from the zero-coupon bond (\$1) in the spot measure numeraire asset (4.24). Letting $\tilde{\mathbb{Q}}^T$ denote the measure induced by $\tilde{P}(t, T)$, we can write

$$\begin{aligned}
V(t) &= \tilde{P}(t, T) \tilde{\mathbb{E}}_t^T \left(V(T_N)/\tilde{P}(T_N, T) \right) \\
&= \tilde{P}(t, T) \tilde{\mathbb{E}}_t^T (V(T_N) B(T)/B(T_N)), \quad T < T_N, \quad t < T_N,
\end{aligned}$$

where $\tilde{\mathbb{E}}^T$ is the expectations operator for the measure $\tilde{\mathbb{Q}}^T$. If also $t \leq T$, this expression becomes (compare to (4.26))

$$V(t) = P(t, T) \tilde{\mathbb{E}}_t^T (V(T_N) B(T)/B(T_N)), \quad t \leq T < T_N,$$

which, in effect, uses $B(T)/B(T_N)$ to discount $V(T_N)$ back to time T . The equivalent result in the T -forward measure is

$$V(t) = P(t, T) \mathbb{E}_t^T (B(T) \mathbb{E}_T^B (V(T_N)/B(T_N))).$$

Notice that if V matures at time T , rather than at T_N , we simply have

$$V(t) = P(t, T) \tilde{\mathbb{E}}_t^T (V(T)) = P(t, T) \mathbb{E}_t^T (V(T)),$$

which is obvious from the definition of the numeraire $\tilde{P}(t, T)$.

The measure $\tilde{\mathbb{Q}}^T$ is by construction a hybrid between the spot measure and the T -forward measure. Obviously, many other such measures exist, corresponding to different reinvestment strategies of expiring numeraire assets.

4.2.5 Swap Measures

Being a linear combination of zero-coupon bonds (see (4.8)), an annuity factor $A_{k,m}(t)$ on a tenor structure qualifies as a numeraire asset. The measure $Q^{k,m}$ induced by this numeraire is known as a *swap measure* or an *annuity measure*. In the absence of arbitrage we have

$$V(t) = A_{k,m}(t) E_t^{k,m} (V(T)/A_{k,m}(T)),$$

where $E^{k,m}(\cdot)$ denotes expectation under $Q^{k,m}$.

Lemma 4.2.4. *In the absence of arbitrage, the forward swap rate $S_{k,m}(t)$ is a martingale in measure $Q^{k,m}$.*

Proof. By definition

$$S_{k,m}(t) = \frac{P(t, T_k) - P(t, T_{k+m})}{A_{k,m}(t)}.$$

As the numeraire deflated assets $P(t, T_k)/A_{k,m}(t)$ and $P(t, T_{k+m})/A_{k,m}(t)$ must both be martingales, so must be their difference. \square

As we shall see later, swap measures are very useful for analytical manipulations of price formulas for options on swaps.

4.3 Multi-Currency Markets

While this book is primarily dedicated to the study of single-currency interest rate derivatives, occasionally it will be necessary to consider certain effects associated with trading in a multi-currency economy. For instance, in Chapter 6 we touch upon issues of yield curve constructions in non-domestic currencies, and in Chapter 16 we discuss the important practical case where a derivative pays out in a foreign currency, but has a payout function that depends on one or more domestic interest rate variables. This brief section provides background material and notation required for these and other cross-currency applications.

4.3.1 Notations and FX Forwards

We consider two economies, a “domestic” economy and a “foreign” economy. Let $P_d(t, T)$ and $P_f(t, T)$ denote time t zero-coupon bond prices in the domestic and foreign economies, respectively. So, $P_f(t, T)$ (say) is the time t price, in foreign currency, of one unit of foreign currency delivered for certain at time T . Translation of values in foreign currency to domestic currency takes place at a foreign exchange (FX) rate of $X(t)$, measured in units of domestic currency per unit of foreign currency. In other words, the value \tilde{P}_d to a domestic investor of one foreign zero-coupon bond is

$$\tilde{P}_d(t, T) = X(t)P_f(t, T).$$

The quantity

$$X_T(t) = \frac{\tilde{P}_d(t, T)}{P_d(t, T)} = X(t) \frac{P_f(t, T)}{P_d(t, T)}$$

is known as the *forward FX rate* to time T . The name is motivated by the following arbitrage strategy:

- Buy one foreign zero-coupon bond, at a cost of $\tilde{P}_d(t, T)$ in domestic currency.
- Finance the purchase by selling short domestic zero-coupon bonds on a notional of $\tilde{P}_d(t, T)/P_d(t, T)$.

With no outlay at time t , the strategy will generate a net cash flow at time T of one unit of foreign currency and $-X_T(t)$ units of domestic currency, such that the trading strategy in effect has locked in a time t a future time T exchange rate of $X_T(t)$.

4.3.2 Risk Neutral Measures

Let $\beta_d(t)$ and $\beta_f(t)$ be the continuously compounded money market accounts in the domestic and foreign economies, respectively. $\beta_d(t)$ and $\beta_f(t)$ induce two separate risk-neutral measures, denoted Q^d and Q^f ; let us investigate how these measures are related. If $g(T)$ is a random payout at time T made in foreign currency, in a complete market the value (in units of foreign currency) of this payout to a foreign investor is, from standard principles,

$$V_f(t) = \beta_f(t)E_t^f(g(T)\beta_f(T)^{-1}), \quad (4.27)$$

where E_t^f denotes expectations in the foreign risk-neutral measure Q^f . For a domestic investor, the payout of $g(T)$ must be translated to domestic currency units at a rate of $X(T)$, making the effective domestic payout function $g(T)X(T)$. Thereby,

$$V_d(t) = \beta_d(t)E_t^d(g(T)X(T)\beta_d(T)^{-1}), \quad (4.28)$$

where E_t^d denotes expectations in measure Q^d . Importantly, the expressions in (4.27) and (4.28) are linked by the spot exchange rate, as the absence of a cross-currency arbitrage dictates that

$$V_d(t) = X(t)V_f(t),$$

or

$$\beta_d(t)E_t^d(g(T)X(T)\beta_d(T)^{-1}) = X(t)\beta_f(t)E_t^f(g(T)\beta_f(T)^{-1}). \quad (4.29)$$

We use this result to establish the following lemma:

Lemma 4.3.1. *The domestic and foreign risk-neutral probability measures Q^d and Q^f are related by the density process*

$$E_t^d \left(\frac{dQ^f}{dQ^d} \right) = \frac{\beta_f(t)X(t)}{\beta_d(t)X(0)}, \quad t \geq 0.$$

Proof. For an \mathcal{F}_T -measurable variable $Y(T) = g(T)X(T)\beta_d(T)^{-1}$ satisfying regularity conditions, a rearrangement of the basic relation (4.29) yields

$$E_t^d(Y(T)) = X(t) \frac{\beta_f(t)}{\beta_d(t)} E_t^f \left(\frac{Y(T)}{X(T)} \frac{\beta_d(T)}{\beta_f(T)} \right),$$

From the results of Section 1.3, the density relating measures Q^d and Q^f is then as given in the lemma. \square

With $\beta_f(t)X(t)/\beta_d(t)$ being a martingale in the domestic risk-neutral measure, we note that if $X(t)$ is an Ito process, it must take the form

$$dX(t) = X(t)(r_d(t) - r_f(t)) dt + \sigma_X(t)^\top dW(t),$$

where $r_d(t) - r_f(t)$ is the spread between domestic and foreign short rates, $W(t)$ is a (vector-valued) Q^d -Brownian motion, and $\sigma_X(t)$ is some adapted stochastic process satisfying regularity conditions.

4.3.3 Other Measures

Having established the Radon-Nikodym derivative relating the domestic and foreign risk-neutral measures, relations between various other domestic and foreign probability measures are easy to establish. For instance, the following result is easily proven the same way as Lemma 4.3.1.

Lemma 4.3.2. *Let $E_t^{T,d}$ denote expectations in the domestic T -forward probability measure. The domestic and foreign T -forward probability measures $Q^{T,d}$ and $Q^{T,f}$ are related by the density process*

$$E_t^{T,d} \left(\frac{dQ^{T,f}}{dQ^{T,d}} \right) = \frac{P_f(t, T)P_d(0, T)X(t)}{P_d(t, T)P_f(0, T)X(0)} = \frac{X_T(t)}{X_T(0)}, \quad t \geq 0.$$

We highlight the fact that the forward FX rate $X_T(t)$ is a $Q^{T,d}$ -martingale satisfying $X_T(T) = X(T)$. For an \mathcal{F}_T -measurable variable $Y(T)$, we thereby have the convenient expression

$$E_t^{T,d}(Y(T)) = X_T(t) E_t^{T,f} \left(\frac{Y(T)}{X(T)} \right),$$

where $E_t^{T,f}$ denotes expectations in the foreign T -forward probability measure.

4.4 The HJM Analysis

Having defined notations and established basic arbitrage relationships, let us turn to assigning dynamics to the many quantities we have introduced so far. We shall here follow the Heath et al. [1992] (Heath-Jarrow-Morton, or HJM) approach, where all information in the economy is assumed to originate with a finite number of Brownian motions. The resulting class of models is quite broad, and in much of the rest of this book we shall deal with ways to reduce the general HJM model to specific, and tractable, special cases. For now, however, we concentrate on a general analysis, although we keep our treatment fairly informal.

4.4.1 Bond Price Dynamics

In the HJM framework, we concern ourselves with the modeling of how an entire continuum of T -indexed bond prices $P(\cdot, T)$ jointly evolves over time, starting from a known condition $P(0, T)$. We consider models of finite horizon, i.e. with $T \in [0, \mathcal{T}]$, $\mathcal{T} < \infty$, and specialize to a filtration generated by a d -dimensional Brownian motion. We assume that a risk-neutral measure Q exists and is unique. Let $W(t)$ be an adapted d -dimensional Q -Brownian motion, and define deflated bond values as $P_\beta(t, T) = P(t, T)/\beta(t)$, where $\beta(t)$ as always is the continuously rolled money market account. In the absence of arbitrage, $P_\beta(t, T)$ is a martingale in the risk-neutral measure, and the martingale representation theorem then implies that

$$dP_\beta(t, T) = -P_\beta(t, T)\sigma_P(t, T)^\top dW(t), \quad t \leq T, \quad (4.30)$$

where $\sigma_P(t, T) = \sigma_P(t, T, \omega)$ is a d -dimensional stochastic process adapted to the filtration generated by W . We assume that $\sigma_P(t, T)$ is regular enough for $P_\beta(t, T)$ to be a square-integrable martingale. Also, as the bond $P(t, T)$ must equal \$1 at $t = T$ (“pull to par”), we impose the consistency condition

$$\sigma_P(T, T) = 0.$$

Using (4.12) and Ito’s lemma, it follows from (4.30) that

$$dP(t, T)/P(t, T) = r(t) dt - \sigma_P(t, T)^\top dW(t), \quad (4.31)$$

where $r(t)$ is the short rate process. Equation (4.31) defines the class of d -dimensional HJM models.

Another application of Ito’s lemma shows that forward bond prices $P(t, T, T + \tau) = P(t, T + \tau)/P(t, T)$ must satisfy

$$\begin{aligned} dP(t, T, T + \tau)/P(t, T, T + \tau) &= -[\sigma_P(t, T + \tau) - \sigma_P(t, T)]^\top \sigma_P(t, T) dt \\ &\quad - [\sigma_P(t, T + \tau) - \sigma_P(t, T)]^\top dW(t). \end{aligned} \quad (4.32)$$

In the T -forward measure \mathbb{Q}^T , $P(t, T, T + \tau)$ is a martingale (see Section 4.2.2), and

$$dP(t, T, T + \tau)/P(t, T, T + \tau) = -[\sigma_P(t, T + \tau) - \sigma_P(t, T)]^\top dW^T(t), \quad (4.33)$$

where $W^T(t)$ is a \mathbb{Q}^T -Brownian motion. Comparison of (4.32) and (4.33) shows that

$$dW^T(t) = dW(t) + \sigma_P(t, T) dt \quad (4.34)$$

which by Girsanov's theorem identifies the density process for the measure shift between \mathbb{Q}^T and \mathbb{Q} in the HJM setting:

$$\varsigma(t) = \mathbb{E}_t^{\mathbb{Q}} \left(\frac{d\mathbb{Q}^T}{d\mathbb{Q}} \right) = \mathcal{E} \left(- \int_0^t \sigma_P(u, T)^\top dW(u) \right), \quad (4.35)$$

or

$$d\varsigma(t)/\varsigma(t) = -\sigma_P(t, T)^\top dW(t).$$

This result could, of course, have been established from the first principles as well — see equation (4.21).

4.4.2 Forward Rate Dynamics

Traditionally, HJM models are stated in terms of instantaneous forward rates, rather than bond prices. Besides eliminating the need to consider the short rate r , this also reveals a number of fundamental properties of the class of HJM models. By Ito's lemma, in measure \mathbb{Q} ,

$$d \ln P(t, T) = O(dt) - \sigma_P(t, T)^\top dW(t),$$

where for convenience we have omitted writing out the drift term. Differentiating the right- and left-hand sides of this equation with respect to T , we get from equation (4.4),

$$df(t, T) = \mu_f(t, T) dt + \sigma_f(t, T)^\top dW(t),$$

where

$$\sigma_f(t, T) = \frac{\partial}{\partial T} \sigma_P(t, T), \quad (4.36)$$

and $\mu_f(t, T)$ is listed below.

Lemma 4.4.1. *The process for $f(t, T)$ in the T -forward measure is*

$$df(t, T) = \sigma_f(t, T)^\top dW^T(t). \quad (4.37)$$

In the risk-neutral measure, the process is

$$\begin{aligned} df(t, T) &= \sigma_f(t, T)^\top \sigma_P(t, T) dt + \sigma_f(t, T)^\top dW(t) \\ &= \sigma_f(t, T)^\top \int_t^T \sigma_f(t, u) du dt + \sigma_f(t, T)^\top dW(t). \end{aligned} \quad (4.38)$$

Proof. The SDE (4.37) follows directly from the martingale relation (4.23). The risk-neutral process (4.38) then can be derived from the relations (4.34) and (4.35), with the second equality following from (4.36). \square

The equation (4.38) is often considered to be the main result of Heath et al. [1992]. It demonstrates that an HJM model is fully specified once the forward rate diffusion coefficients $\sigma_f(t, T)$ have been specified for all t and T . Note that HJM models take initial forward rates $f(0, T)$ as exogenous inputs, ensuring that these models are automatically consistent with discount bond prices at time 0. This is true irrespective of the choice of $\sigma_f(t, T)$, which can be set freely (subject to regularity conditions) from either empirical analysis, or from a calibration to market prices of fixed income derivatives.

While it is convenient that HJM models are automatically calibrated to initial bond prices, a number of other features of the general HJM model are less attractive. Particularly problematic is the sheer dimensionality of the model: to describe the time t state of a discount bond curve spanning $[t, T]$, we need to keep track of a continuum of forward rates $\{f(t, u), t \leq u \leq T\}$. By Lemma 4.4.1 the forward rate curve follows an infinite-dimensional diffusion process, leaving us with an infinite number of state variables to diffuse. In practice, the implementation of an HJM model will require either making special assumptions about the σ_f process that permit a finite-dimensional Markovian representation of the forward rate curve; or moving from infinitesimal forward rates to continuously compounded forward rates that span time-buckets of finite length. Chapters 10 through 13 and Section 4.5.2 below give examples of the former idea, and Andersen [1995] discusses the latter approach in a Monte Carlo setting. An idea closely related to the discussion in Andersen [1995] is to build a model around a finite set of simple (Libor) forward rates on a fixed tenor structure. This approach has a number of computational and theoretical advantages, and is the subject of Chapter 14. For now, we note that *any* arbitrage-free interest rate model set in a filtration generated exclusively by Brownian motions must be a special case of an HJM model. In particular, any such model must correspond to a particular choice of $\sigma_f(t, T)$.

4.4.3 Short Rate Process

As discussed earlier, specification of a short rate process is, in principle, sufficient to completely specify a full yield curve model. In the HJM framework, it follows from (4.38) that the short rate $r(t)$ in measure Q is

$$r(t) = f(t, t) = f(0, t) + \int_0^t \sigma_f(u, t)^\top \int_u^t \sigma_f(u, s) ds du + \int_0^t \sigma_f(u, t)^\top dW(u).$$

The process for $r(t)$ is generally not Markovian, as can be seen by focusing on the path-dependent term

$$D(t) = \int_0^t \sigma_f(u, t)^\top dW(u)$$

for which we must have

$$\begin{aligned} D(T) &= D(t) + \int_t^T \sigma_f(u, T)^\top dW(u) \\ &\quad + \left\{ \int_0^t \sigma_f(u, T)^\top dW(u) - \int_0^t \sigma_f(u, t)^\top dW(u) \right\}. \end{aligned} \quad (4.39)$$

Thereby

$$E^Q(D(T)|D(t)) \neq E_t^Q(D(T))$$

unless the bracketed term in (4.39) is either non-random, or a deterministic function of $D(t)$ (which is generally not the case).

An interesting area of investigation concerns the conditions under which either $r(t)$ is outright Markov⁴ or, less restrictively, can be written as

$$r(t) = h(t, x(t)),$$

for a deterministic function h and a finite-dimensional Markovian vector of state variables $x(t)$. Definitive results are given in Björk [2001], building on earlier (and considerably less abstract) work by Jamshidian [1991b], Cheyette [1991], and Ritchken and Sankarasubramanian [1995]. Section 4.5.2 and Chapter 13 list some of the results of these papers.

4.5 Examples of HJM Models

4.5.1 The Gaussian Model

In the HJM bond price dynamics (4.31), we now assume that $\sigma_P(t, T)$ is a bounded (d -dimensional) deterministic function of t and T . It follows from (4.32) and (4.33) that forward bond prices are then log-normally distributed in both Q and Q^T . The forward rate process in Q is

$$df(t, T) = \sigma_f(t, T)^\top \sigma_P(t, T) dt + \sigma_f(t, T)^\top dW(t), \quad \sigma_f(t, T) = \frac{\partial}{\partial T} \sigma_P(t, T), \quad (4.40)$$

which implies that $r(T) = f(T, T)$ is Gaussian with Q -moments

$$\begin{aligned} E_t^Q(f(T, T)) &= \int_t^T \sigma_f(u, T)^\top \sigma_P(u, T) du, \\ \text{Var}_t^Q(f(T, T)) &= \int_t^T \sigma_f(u, T)^\top \sigma_f(u, T) du. \end{aligned}$$

⁴In the sense that the time t expectations of functionals of $r(T)$ only require knowledge of $r(t)$ itself. In this case the process for $r(t)$ must be a diffusion characterized by an SDE $dr(t) = \mu_r(t, r(t))dt + \sigma_r(t, r(t))^\top dW(t)$.

The simple form of the Gaussian HJM model makes it quite tractable, permitting analytical price formulas for a number of European options and futures contracts⁵. While the Gaussian HJM model suffers from the drawback of allowing negative forward and spot rates, analytical results derived in the model are often very useful in gaining a deeper understanding of a given contract, even if ultimately a more realistic model will be required for serious pricing purposes. Indeed, results derived for the Gaussian HJM model can often be used as a starting point for development of closed-form approximations in other models; we shall see many examples of this later in the book.

For illustration, we list a few select analytical results below. More formulas can be found in numerous sources, including Chapter II in Andersen [1996], Jamshidian [1991b], and Jamshidian [1993], to name a few.

Proposition 4.5.1 (Option on Zero-Coupon Bond). *Consider a European call option paying at maturity T the amount*

$$V(T) = (P(T, T^*) - K)^+, \quad T^* > T.$$

In the Gaussian HJM model (4.40), we have

$$V(t) = P(t, T^*)\Phi(d_+) - P(t, T)K\Phi(d_-), \quad (4.41)$$

where

$$d_{\pm} = \frac{\ln(P(t, T^*) / (KP(t, T))) \pm v/2}{\sqrt{v}},$$

$$v = \int_t^T |\sigma_P(u, T^*) - \sigma_P(u, T)|^2 du.$$

Proof. In the T -forward measure \mathbb{Q}^T we have, from (4.20),

$$\begin{aligned} V(t) &= P(t, T)\mathbb{E}_t^T \left((P(T, T^*) - K)^+ \right) \\ &= P(t, T)\mathbb{E}_t^T \left((P(T, T, T^*) - K)^+ \right). \end{aligned}$$

From the discussion in Section 4.4.1 we know that $P(t, T, T^*)$ is a \mathbb{Q}^T -martingale characterized by the SDE

$$dP(t, T, T^*)/P(t, T, T^*) = -[\sigma_P(t, T^*) - \sigma_P(t, T)]^\top dW^T(t),$$

where W^T is a d -dimensional \mathbb{Q}^T -Brownian motion. As this is just a GBM process with time-dependent coefficients, the Black-Scholes-Merton results in Section 1.9 apply and lead to (4.41). \square

⁵Indeed, we have already used this model in an equity context — see Section 1.9.3.2.

Proposition 4.5.2 (Caplet). Consider a European call option paying at $T + \tau$ the amount (a caplet)

$$V(T + \tau) = \tau (L(T, T, T + \tau) - K)^+, \quad \tau > 0.$$

In the Gaussian HJM model (4.40), we have

$$V(t) = \tau P(t, T + \tau) (L(t, T, T + \tau) \Phi(d_+) - K \Phi(d_-)),$$

$$d_{\pm} = \frac{\ln(L(t, T, T + \tau)/K) \pm v/2}{\sqrt{v}}, \quad v = \int_t^T |\sigma_P(u, T + \tau) - \sigma_P(u, T)|^2 du.$$

Proof. From Lemma 4.2.3 we know that $L(t, T, T + \tau)$ is a martingale in the $(T + \tau)$ -forward measure. An application of Ito's lemma to the definition (4.2) reveals that

$$dL(t, T, T + \tau)/L(t, T, T + \tau) = [\sigma_P(t, T + \tau) - \sigma_P(t, T)]^\top dW^T(t),$$

and the result follows immediately along the same lines as in the proof of Proposition 4.5.1. \square

Proposition 4.5.3 (Futures Rate). In the Gaussian HJM model (4.40), futures rates are given by

$$F(t, T, T + \tau) = \tau^{-1} \left((1/P(t, T, T + \tau)) e^{\Omega(t, T)} - 1 \right), \quad (4.42)$$

where

$$\Omega(t, T) = \int_t^T [\sigma_P(u, T + \tau) - \sigma_P(u, T)]^\top \sigma_P(u, T + \tau) du.$$

Proof. From Lemma 4.2.2,

$$\begin{aligned} F(t, T, T + \tau) &= E_t^Q (L(T, T, T + \tau)) \\ &= \tau^{-1} E_t^Q (1/P(T, T + \tau) - 1) \\ &= \tau^{-1} E_t^Q (G(T) - 1), \end{aligned} \quad (4.43)$$

where we have introduced an auxiliary variable

$$G(t) \triangleq P(t, T)/P(t, T + \tau) = 1/P(t, T, T + \tau).$$

Ito's lemma shows that (see also (4.32)) in measure Q

$$\begin{aligned} dG(t)/G(t) &= [\sigma_P(t, T + \tau) - \sigma_P(t, T)]^\top \sigma_P(t, T + \tau) dt \\ &\quad + [\sigma_P(t, T + \tau) - \sigma_P(t, T)]^\top dW(t), \end{aligned}$$

such that

$$E_t^Q(G(T)) = G(t)e^{\Omega(t,T)} = (1/P(t, T, T + \tau)) e^{\Omega(t,T)},$$

where Ω is as given above. The result of Proposition 4.5.3 then follows directly from (4.43). \square

In any rational model $\Omega(t, T) \geq 0$, such that $F(t, T, T + \tau) \geq L(t, T, T + \tau)$, consistent with the qualitative discussion in Section 4.1.2. As shown in Chapter II of Andersen [1996], the spread (also known as futures *convexity*) between futures and forward rates can be decomposed into two components: i) a term originating from the mark-to-market mechanism of a futures contract; and ii) a term originating from the fact that a futures contract — unlike a regular forward rate agreement — pays out the rate at the date it settles (at time T) rather than one period ahead (at time $T + \tau$). Andersen [1996], Chapter II, additionally contains a number of numerical examples examining typical futures-forward spreads, and also investigates the pricing of options on futures rates.

Section 16.8 looks in detail into pricing interest rate futures under more advanced models.

4.5.2 Gaussian HJM Models with Markovian Short Rate

Although quite tractable, the Gaussian HJM model generally does not allow for a finite-dimensional Markovian representation, and typically does not imply Markov-diffusive behavior of the short rate. As shown in Carverhill [1994], the short rate can be made Markovian, however, by imposing certain conditions on the deterministic forward rate volatility function $\sigma_f(t, T)$. To explore this, first recall from Section 4.4.3 the relation

$$r(t) = f(0, t) + \int_0^t \sigma_f(u, t)^\top \int_u^t \sigma_f(u, s) ds du + \int_0^t \sigma_f(u, t)^\top dW(u),$$

where now σ_f is deterministic. Consider imposing the special choice

$$\sigma_f(t, T) = g(t)h(T), \quad (4.44)$$

where h is a positive real function and $g : \mathbb{R} \rightarrow \mathbb{R}^{d \times 1}$ can take any sign. For this case we have

$$\sigma_P(t, T) = \int_t^T \sigma_f(t, u) du = g(t) \int_t^T h(u) du,$$

and

$$\begin{aligned} r(t) &= f(0, t) + h(t) \int_0^t g(u)^\top g(u) \left(\int_u^t h(s) ds \right) du + h(t) \int_0^t g(u)^\top dW(u) \\ &\triangleq f(0, t) + h(t) \int_0^t m_f(t, u) du + h(t) \int_0^t g(u)^\top dW(u). \end{aligned} \quad (4.45)$$

Importantly, the term

$$D(t) = \int_0^t \sigma_f(u, t)^\top dW(u) = h(t) \int_0^t g(u)^\top dW(u)$$

is now Markov, since

$$D(T) = h(T) \int_0^T g(u)^\top dW(u) = \frac{h(T)}{h(t)} D(t) + h(T) \int_t^T g(u)^\top dW(u),$$

which should be compared to the general (non-Markov) expression (4.39).

To show that the short rate is Markovian, we differentiate (4.45) with respect to t , yielding

$$\begin{aligned} dr(t) &= \frac{\partial f(0, t)}{\partial t} dt + h'(t) \left(\int_0^t m_f(t, u) du + \frac{D(t)}{h(t)} \right) dt \\ &\quad + h(t) \frac{\partial}{\partial t} \int_0^t m_f(t, u) du dt + h(t) g(t)^\top dW(t) \\ &= \frac{\partial f(0, t)}{\partial t} dt + h'(t) \left(\frac{r(t) - f(0, t)}{h(t)} \right) dt \\ &\quad + h(t)^2 \int_0^t g(u)^\top g(u) du dt + h(t) g(t)^\top dW(t) \\ &= \left(\frac{\partial f(0, t)}{\partial t} - \frac{h'(t)}{h(t)} f(0, t) + h(t)^2 \int_0^t g(u)^\top g(u) du + \frac{h'(t)}{h(t)} r(t) \right) dt \\ &\quad + h(t) g(t)^\top dW(t), \end{aligned} \tag{4.46}$$

where the second equality follows from rearrangement of (4.45). This leads to the following result.

Proposition 4.5.4. *In the d -dimensional Gaussian HJM model, when (4.44) holds the short rate satisfies an SDE of the type*

$$dr(t) = (a(t) - \varkappa(t)r(t)) dt + \sigma_r(t)^\top dW(t),$$

where $\varkappa : \mathbb{R} \rightarrow \mathbb{R}$ and $\sigma_r : \mathbb{R} \rightarrow \mathbb{R}^{d \times 1}$ are deterministic functions of time, and

$$\begin{aligned} a(t) &= \frac{\partial f(0, t)}{\partial t} + \varkappa(t)f(0, t) + \int_0^t e^{-2 \int_u^t \varkappa(s) ds} \sigma_r(u)^\top \sigma_r(u) du \\ &= \frac{\partial f(0, t)}{\partial t} + \varkappa(t)f(0, t) + \int_0^t \sigma_f(u, t)^\top \sigma_f(u, t) du. \end{aligned}$$

Proof. First, by way of defining \varkappa and σ_r , we set

$$h(T) = e^{-\int_0^T \varkappa(s) ds}; \quad g(t) = e^{\int_0^t \varkappa(s) ds} \sigma_r(t),$$

such that $h'(t)/h(t) = -\varkappa(t)$ and

$$\sigma_f(t, T) = e^{-\int_t^T \varkappa(s) ds} \sigma_r(t).$$

The result of Proposition 4.5.4 then follows directly by insertion into (4.46).

□

4.5.3 Log-Normal HJM Models

To avoid the negative forward rates inherent in Gaussian HJM models, it is tempting to consider forward rate specifications of the type

$$\sigma_f(t, T) = f(t, T)\sigma(t, T), \quad (4.47)$$

where $\sigma(t, T)$ is deterministic and bounded. In the T -forward measure

$$df(t, T) = f(t, T)\sigma(t, T)^\top dW^T(t)$$

such that $f(t, T)$ is log-normally distributed. While avoiding negative rates, the specification (4.47) has severe technical problems: in Q , forward rates will explode to infinity with non-zero probability. Attempts to apply the valuation formula (4.15) will thereby result in all zero-coupon bond prices being zero, implying obvious arbitrage opportunities. To suggest a rationale for the exploding rates, consider the Q -dynamics

$$df(t, T) = \left(f(t, T)\sigma(t, T)^\top \int_t^T f(t, u)\sigma(t, u)du \right) dt + f(t, T)\sigma(t, T)^\top dW(t).$$

Loosely speaking, the drift-term is proportional to forward rates *squared*, which, in the light of the linear growth condition in Theorem 1.6.1, may cause us to suspect problems with the existence of a non-exploding solution. Morton [1988] confirms this rigorously.

One solution to the explosion problem involves enforcing a strict upper bound on $\sigma_f(t, T)$, as in

$$\sigma_f(t, T) = \min(f(t, T), M)\sigma(t, T),$$

where M is a large positive constant. For the one-factor case ($d = 1$) Heath et al. [1992] demonstrate that this specification will ensure non-negative forward rates⁶ and will prevent rate explosions. Nevertheless, the model is clearly awkward in its dependence on the arbitrary constant M . A more satisfying solution is discussed in Chapter 14, where we show that the explosion problem can be circumvented by working with simply — rather than continuously — compounded forward rates. A related issue in short rate models is also discussed in Chapter 11.

⁶To see this, notice that $M > 0$ guarantees that $df(t, T) = 0$ if $f(t, T)$ should ever reach 0.

Fixed Income Instruments

At this point, we have established the mathematical and numerical prerequisites needed for the remaining part of the book, much of which is devoted to the development of models for fixed income derivatives. Before delving into the modeling exercise, this final foundational chapter provides a tour of actual fixed income markets as well as an overview of the types of products traded. The simpler (and more liquid) of these products will typically serve as calibration targets to parameterize the models we develop; others (the more complicated and illiquid ones) will constitute the contracts that our model are ultimately meant to price and hedge. Throughout the chapter — and, indeed, this book — our focus is on the securities tied to the so-called *Libor rate*; this will include essentially all high-end exotic securities as well as more basic instruments such as FRAs, caps, and swaptions. Our priorities dictate that we leave out government, corporate, and mortgage bonds, as well as the derivatives associated with these types of securities. A discussion of these classes of securities, along with many more details on the organization and workings of fixed income markets, can be found in specialist literature, such as Fabozzi and Modigliani [1996], Fabozzi [1985], Fabozzi [2001], and Fabozzi and Fabozzi [1989].

5.1 Fixed Income Markets and Participants

At the most fundamental level, interest rates determine the economic cost of borrowing and lending, and as such define present values of future cash flows. In general, cash flows occurring at different times are discounted at different rates, reflecting market fluctuations in demand for money and risk preferences of market participants. The dependence of interest rates on time is described by the so-called *term structure of interest rates*, easily visualized as a curve that assigns a particular interest rate (or, equivalently, a discount factor) to each future date.

For a given entity, the cost of borrowing money will depend on its credit quality. Governments of developed countries, perceived to have virtually no possibility of default, issue bonds at comparatively low interest rates that reflect this perception. While the market in government debt is vast, corporations typically find it more convenient to use and originate fixed income instruments linked to rates that are more reflective of their own financing costs (i.e., credit quality). By far the most common of such reference rates is the London Interbank Offered rate, commonly known as the *Libor rate*. The Libor rate is a filtered average of bank estimates of rates at which they can borrow for a given term in the *interbank money market*, i.e. the wholesale market in which banks provide unsecured short-term credit to each other. Libor rates are quoted for multiple deposit maturities ranging from one day to one year, and are set every business day by averaging polling results from a number of large banks. Libor rates are available for deposits in different currencies, so that there is a USD-Libor rate, a EUR-Libor rate, and so on.

While Libor rates are probably the most used reference rates for interest rate contracts, there are other important rates to be aware of. For example, in the United States, banks are required to hold certain balances (“Federal funds”) with the Federal Reserve, the central bank of the US. If a bank does not have sufficient balances, it can borrow them from another bank that has an excess on its account. The overnight interest rate charged in this case is called the (*effective*) *Federal funds rate*¹, or sometimes simply the Fed funds rate. This rate is often considered the best available proxy for a risk-free USD rate, in part because the Fed funds rate is normally the contractual rate used to accrue interest on posted collateral², as explained in Piterbarg [2010]. It is worth noting that the Fed funds rate used to be closely linked to the overnight Libor rate, with the spread between the two in the single basis points. However, in the subprime crisis of 2007–2009 the two have diverged significantly; the implications of this for interest rate curve construction are discussed in Section 6.5.3. Instruments linked to averages of the (*effective*) Fed fund rate over different terms are actively traded, giving rise to a term structure of Fed funds linked rates.

A special feature of the US public debt markets gives rise to another set of rates. In particular, interest on bonds issued by states and other local governments of the US is often free of the federal tax. The Bond Market Association, a trade association of the bond industry, publishes the *BMA rate* (or *BMA index*) which is the estimate of borrowing by such municipalities.

¹The target rate, set by the Federal reserve, is aptly called the *target* Fed funds rate.

²To mitigate credit risk, many derivatives transactions require posting of collateral (normally cash or Treasury bonds) in the amount of the current mark-to-market. ISDA [2005] contains a detailed description of collateral agreements; according to ISDA [2009], in 2009 about 65% of all OTC derivatives transactions involved such agreements.

There is a well-developed market in interest rate derivatives that are linked to the BMA rate.

The Euro and GBP markets do not have the same mechanism as the US does for Federal funds, but overnight rates that are proxies for risk-free borrowing in those currencies do exist. They are called *Eonia* (Euro OverNight Index Average) in the Eurozone and *Sonia* (Sterling OverNight Index Average) in Great Britain, and are computed as averages of all *actual* overnight lending/borrowing transactions by qualifying banks weighted by the size of the transactions. We emphasize that these rates reflect the actual transactions that have happened, in contrast to Libor which reflect banks' estimates of rates at which borrowing (for a given term) might take place. In the crisis of 2007–2009 there have been serious concerns about the integrity of the Libor rate and whether it really reflected the actual cost of funding for banks, and even some calls to scrap the Libor rate altogether. While the Libor rate has survived the crisis, the importance of overnight rates has increased dramatically, with the market in FedFunds/Eonia/Sonia linked derivatives, most importantly in *overnight index swaps*, or OIS, of various maturities growing dramatically. As with the Fed funds rate, Eonia and Sonia have diverged significantly from the corresponding Libor rates during market turbulence, and the decoupling continues to persist. As with the Fed funds rate, the implications of these developments on interest rate curve construction are discussed in Section 6.5.3.

Interest rates change day-to-day in response to changing macroeconomic and market conditions. With the cost of borrowing and lending money affecting all aspects of the economy, it is no surprise that a vast market in derivatives on interest rates has developed. Motivations of participants are diverse, ranging from locking in the cost of financing to pure speculation.

The fixed income market can be broadly split into two (overlapping) segments: the *exchange* market and the *over-the-counter*, or *OTC*, market. Contracts linked to the level of interest rates are traded on many securities exchanges. The exchanges attract all types of investors, including market makers, hedgers and speculators; see Hull [2006] for details on all. As of March 2008, notional amounts outstanding were \$26 trillion in exchange traded interest rate futures, and \$45 trillion in exchange traded interest rate options. While these are impressive numbers, far more fixed income derivatives trade in OTC markets than in exchange markets: as of December 2007, the notional amounts outstanding of OTC interest rate derivatives amounted to \$393 trillion³. The OTC market can loosely be visualized as a network of banks that trade with each other under terms governed by agreements spelled out by the trade organization International Swaps and Derivatives Association (ISDA). Central to OTC markets are the *interest*

³All figures from the report “Semiannual OTC derivatives statistics at end-December 2007” by Bank for International Settlements, available from www.bis.org.

rate dealers, banks with trading desks specializing in fixed income trading. The dealers provide liquidity in various types of securities, and are typically the most sophisticated players in the market. The dealers trade either on their own account or on behalf of customers such as *financial institutions* and *corporates*.

Financial institutions include mortgage companies (organizations that originate, package or service residential and commercial mortgage loans), pension funds, mutual funds, insurance companies, hedge funds, and other entities whose primary activities are related to financial markets. Financial institutions seek to either make money directly by engaging in trading activities (hedge funds), or to hedge their exposures (mortgage originators or servicers), or to achieve superior returns on their investments (pension funds, insurance companies). Among financial institutions, an important role is played by *issuers*, companies that issue structured notes for private and public placement. Structured notes deliver appealing return profiles to investors, returns that are essentially financed by selling options back to issuers. Issuance of increasingly complicated structured notes drives the exotic end of the fixed income markets.

Corporates are companies with primary activities not directly linked to fixed income markets, but whose operational results may be affected by the interest rate environment. For instance, many companies raise funds by borrowing from banks or by selling bonds, and are therefore affected by the prevailing levels of interest rates. Corporates often seek to lock in favorable interest rates for borrowing money, to hedge their interest rate exposures, to transform their liabilities from one type (e.g., a fixed rate liability) to another (a floating rate liability), or to design custom borrowing schemes around their expected future borrowing needs.

5.2 Certificates of Deposit and Libor Rates

Having identified the main types of market participants, we now proceed to define the universe of securities that this book will cover. For technical precision, we shall occasionally need to refer to the risk-neutral measure Q , as well as its associated expectation operator $E = E^Q$ and its numeraire $\beta(t)$.

We start with the *certificate of deposit* (or *CD*), a deposit of money for a pre-specified term at a pre-specified interest rate. Terms may range from one week to one year or more, with the most popular being a 3 month or a 6 month term, depending on the currency of the deposit. If 1 (dollar) is deposited at time T for a period of τ years, then the amount of capital to be returned at time $T + \tau$ is given by⁴

⁴As was mentioned earlier, the computation of τ from given start- and end-dates will involve certain formal day counting rules, see Appendix 5.A.

$$1 + \tau L,$$

where L is, by definition, the interest rate for the CD. The rate is quoted as a simple rate, i.e. a rate with the compounding frequency equal to the term of the deposit. Notice that the average value of L for CDs quoted in the interbank market will, by definition, be equal to the (spot) Libor rate for tenor τ . Spot Libor rates for various tenors are calculated daily and are published by major news services such as Bloomberg or Reuters. As mentioned above, Libor serves as the primary reference rate in fixed income markets.

If $P(T, T + \tau)$ is the (Libor-based) discount factor to date $T + \tau$ as observed at T , then the discounted value of receiving $1 + \tau L$ at time $T + \tau$ should be equal to 1 at time T , i.e.

$$1 = P(T, T + \tau)(1 + \tau L).$$

In particular, recalling the definition (4.2) of $L(t, T, T + \tau)$, the rate L paid on the CD is a simple spot rate

$$L = L(T, T, T + \tau) = \frac{1}{\tau} \left(\frac{1}{P(T, T + \tau)} - 1 \right). \quad (5.1)$$

5.3 Forward Rate Agreements (FRA)

A certificate of deposit allows a market participant to lock in an interest rate for a given period of time, effective immediately. Many market participants, however, find it convenient to lock in interest rates for a given period of time that starts in the future. Contracts that provide such a rate guarantee are known as *forward contracts* or, in a fixed income context, *forward rate agreements* (FRAs). An FRA for the period $[T, T + \tau]$ is a contract to exchange fixed rate payment (agreed at the initiation of the contract) against a payment based on the time T spot Libor rate of tenor τ . While all payments on an FRA are exchanged at, or near⁵, time T , the contract is structured so that the payments are made in $T + \tau$ dollars.

Formally, consider the origination at time t , $t \leq T$, of a unit notional FRA contract with a rate of k . Ignoring payment delays, from the perspective of the fixed rate payer the net payment at time T will be

$$V_{\text{FRA}}(T) = \tau(L(T, T, T + \tau) - k) / (1 + \tau L(T, T, T + \tau)),$$

with the (contractually specified) factor $1/(1 + \tau L(T, T, T + \tau))$ applied to roll the payment to the future date $T + \tau$. We note that

⁵Typical market conventions call for a two business day payment delay, see Appendix 5.A for more details.

$$1 / (1 + \tau L(T, T, T + \tau)) = P(T, T + \tau)$$

so, by the fundamental pricing result (4.13), the value of this contract at time t is equal to

$$V_{\text{FRA}}(t) = \beta(t) E_t (\beta(T)^{-1} \tau (L(T, T, T + \tau) - k) P(T, T + \tau))$$

(recall that $\beta(\cdot)$ is the money market account). Substituting (5.1) we obtain

$$V_{\text{FRA}}(t) = \beta(t) E_t (\beta(T)^{-1} (1 - P(T, T + \tau) - \tau k P(T, T + \tau))).$$

Since $P(\cdot, T + \tau)$ is a traded asset, its price deflated by the numeraire $\beta(\cdot)$ is a martingale. Thus

$$\begin{aligned} V_{\text{FRA}}(t) &= P(t, T) - P(t, T + \tau) - \tau k P(t, T + \tau) \\ &= \tau P(t, T + \tau) \left(\frac{P(t, T) - P(t, T + \tau)}{\tau P(t, T + \tau)} - k \right). \end{aligned} \quad (5.2)$$

Most often, FRAs are issued at no cost to either party at the time of origination. The value of k that makes the FRA contract have value 0 at the contract initiation time t is given by the *forward Libor rate* (see (4.2)),

$$k = L(t, T, T + \tau) = \frac{P(t, T) - P(t, T + \tau)}{\tau P(t, T + \tau)}.$$

Thus, a forward Libor rate has the financial interpretation of being a break-even rate on an FRA contract in interbank markets.

5.4 Eurodollar Futures

FRAs, being forward contracts on Libor rates, allow market participants to either lock in favorable rates for future periods, or to speculate on the future direction of rates. FRAs trade in the OTC market, and are open only to institutions that participate in this market. Alternatively, *futures contracts* on Libor rates are available on a number of international exchanges, including the Chicago Mercantile Exchange (CME), London International Financial Futures and Options Exchange (LIFFE), and Marché à Terme International de France (MATIF). The CME interest rate futures contract on a three-month spot Libor rate on US dollar denominated deposits is called the *Eurodollar futures* or, simply, *ED futures* contract.

At maturity T , an ED futures contract is settled at

$$100 \times (1 - L(T, T, T + \tau)).$$

The *futures rate* $F(t, T, T + \tau)$ at time t (see (4.1.2)) is defined to be the rate such that the *quoted futures price* at time t is equal to⁶

⁶So, if the futures rate is 5%, the quoted futures price is 95.

$$100 \times (1 - F(t, T, T + \tau)).$$

As is the case for all futures contracts, ED futures are settled (marked to market) daily. Confusing matters somewhat, the actual amount of money that is settled between holders of the long and the short positions in an ED future is determined by the daily change in the *actual futures price* defined by

$$N_{ED} \times \left[1 - \frac{1}{4} F(t, T, T + \tau) \right],$$

where N_{ED} is the notional principal of the contract (\$1,000,000 for the CME's ED futures). In particular, for 1 basis point (0.01%) increase in the rate $F(t, T, T + \tau)$, the CME contract buyer pays $1,000,000 \times 0.25 \times 0.0001 = 25$ dollars to the seller.

As explained in Chapter 4, futures rates $F(t, T, T + \tau)$ are generally different from forward Libor rates $L(t, T, T + \tau)$. The problem of computing the difference, the *ED convexity adjustment*, is considered in Section 16.8.

Unlike FRAs, for which the deposit period is negotiated between two parties, ED futures are standardized. Available contracts expire on four specific dates, one each in March, June, September and December, over the next ten years. Such standardization increases liquidity in each particular contract.

5.5 Fixed-for-Floating Swaps

A *swap* is a generic term for an OTC derivative in which two counterparties agree to exchange one stream of cash flows against another stream. These streams are called the *legs* of the swap. A *plain vanilla fixed-for-floating interest rate swap* (a *plain vanilla swap*, or just a *swap* if there is no confusion) is a swap in which one leg is a stream of fixed rate payments and the other a stream of payments based on a floating rate, most often Libor. The legs are denominated in the same currency, have the same notional, and expire on the same date. Payment streams are made on a pre-defined schedule of contiguous time intervals, known as *periods*. Typically, the floating rate is observed (or *fixed*) at the beginning of each period, with both fixed and floating rate coupons being paid out at the end of the period. A plain-vanilla swap is economically equivalent⁷ to a multi-period FRA, and serves the same purpose in the market as regular FRAs. Between interest rate dealers

⁷This is true up to subtle but potentially important discounting issues. As we have pointed out in Section 5.3, the net payment of an FRA is *contractually* discounted using Libor rate from $T + \tau$ to T , whereas in a swap, the net payment for a given period is discounted at the money market account rate from the end to the beginning of the accrual period. The two types of discounting can in fact be different in the presence of discounting-index *basis*, see Sections 6.5.2 and 6.5.3.

and financial institutions, swaps of different maturities are often traded to adjust interest risk positions of the parties involved, or to simply make bets on future direction of interest rates. Swaps are also used by corporates, often in conjunction with bond or note issuance, to transform fixed rate obligations into floating ones, or vice versa.

To formally define a fixed-floating swap, one specifies a tenor structure, i.e. an increasing sequence of maturity times, normally spaced roughly equidistantly (see Section 4.1.3)

$$0 \leq T_0 < T_1 < T_2 < \dots < T_N, \quad \tau_n = T_{n+1} - T_n. \quad (5.3)$$

In a fixed-floating swap with fixed rate k , one party (the fixed rate payer) pays simple interest based on the rate k in return for simple interest payments computed from the Libor rate fixing on date T_n , for each period $[T_n, T_{n+1}]$, $n = 0, \dots, N - 1$. The payments are exchanged at the end of each period, i.e. at time T_{n+1} . In practice, the payments are netted, and only their difference changes hands. From the perspective of the fixed rate payer, the net cash flow of the swap at time T_{n+1} is therefore given by (on a unit notional)

$$\tau_n (L_n(T_n) - k), \quad L_n(t) = L(t, T_n, T_{n+1}),$$

for $n = 0, \dots, N - 1$. Dates when the Libor rates are observed are typically called *fixing dates*; dates when payments occur are called *payment dates*.

By the fundamental valuation result (4.13), the value of a swap is equal to the expected discounted value of its (netted) payments. Specifically, the value to the fixed rate payer of a unit notional fixed-floating swap at time t , $0 \leq t \leq T_0$, is given by⁸

$$\begin{aligned} V_{\text{swap}}(t) &= \beta(t) \sum_{n=0}^{N-1} \tau_n E_t \left(\beta(T_{n+1})^{-1} (L_n(T_n) - k) \right) \\ &= \beta(t) \sum_{n=0}^{N-1} \tau_n E_t \left(\beta(T_n)^{-1} (L_n(T_n) - k) P(T_n, T_{n+1}) \right). \end{aligned}$$

Using the definition of Libor rates $L_n(T_n)$,

$$V_{\text{swap}}(t) = \beta(t) \sum_{n=0}^{N-1} E_t \left(\beta(T_n)^{-1} (1 - P(T_n, T_{n+1}) - \tau_n k P(T_n, T_{n+1})) \right).$$

For each n , $P(\cdot, T_n)$ is a traded asset, so its price deflated by the numeraire $\beta(\cdot)$ is a martingale. Hence

$$V_{\text{swap}}(t) = \sum_{n=0}^{N-1} (P(t, T_n) - P(t, T_{n+1}) - \tau_n k P(t, T_{n+1})).$$

⁸This is a somewhat idealized expression. See Appendix 5.A for more details on market day counting conventions and related topics.

Recalling the definition of $L_n(t)$, this can be rewritten as

$$V_{\text{swap}}(t) = \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) (L_n(t) - k).$$

An important observation is that a vanilla fixed-floating swap can be valued on date t using only the term structure of interest rates observed on that date. In particular, swap values are not affected by the *dynamics* of interest rates, only their current levels.

The swap valuation formula above can be rewritten as follows,

$$V_{\text{swap}}(t) = \left(\sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) \right) \left(\frac{\sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) L_n(t)}{\sum_{n=0}^{N-1} \tau_n P(t, T_{n+1})} - k \right).$$

Using the definitions (4.8), (4.10) and (4.11) from Chapter 4:

$$A(t) \triangleq A_{0,N}(t) = \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}), \quad (5.4)$$

$$S(t) \triangleq S_{0,N}(t) = \frac{\sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) L_n(t)}{\sum_{n=0}^{N-1} \tau_n P(t, T_{n+1})}, \quad (5.5)$$

we obtain the convenient formula

$$V_{\text{swap}}(t) = A(t) (S(t) - k). \quad (5.6)$$

The quantity $A(\cdot)$ is the *annuity* of the swap (or its *PVBP*, for Present Value of a Basis Point), and the quantity $S(t)$ is the *forward swap rate*. Clearly, $S(t)$ is the value of the fixed rate that makes the swap have value 0 to both parties at time t ; S is consequently often referred to as a *par* or *break-even* rate.

For plain-vanilla swaps, the fixed rate and the swap notional are constant through time. More general swaps are, however, not bound by such restrictions and both the fixed rate and the notional may vary from period to period. A non-standard swap with a notional schedule $\{q_n\}_{n=0}^{N-1}$ (non-constant but deterministic) and a fixed rate schedule $\{k_n\}_{n=0}^{N-1}$ has the value

$$\begin{aligned} V_{\text{genswap}}(t) &= \beta(t) \sum_{n=0}^{N-1} \tau_n q_n E_t \left(\beta(T_{n+1})^{-1} (L_n(T_n) - k_n) \right) \\ &= \sum_{n=0}^{N-1} \tau_n q_n P(t, T_{n+1}) (L_n(t) - k_n). \end{aligned}$$

Certain general swaps have dedicated names, such as *amortizing swaps* (notional decreases with time) and *accreting swaps* (notional increases with time).

As we mentioned in Section 5.1, swaps linked to overnight rates (Fed-Funds/Eonia/Sonia) have recently become more popular. Among them the overnight index swap (OIS) is probably the most liquid, and is defined as a swap that pays a *compounded* overnight rate against fixed rate payments. To write down its definition, let us assume that a tenor structure (5.3) is given, and denote by $\{t_{n,i}\}_{i=1}^{K_n}$ the collection of all business days in the period $[T_n, T_{n+1})$, so that $T_n = t_{n,1} < \dots < t_{K_n} < T_{n+1}$. Then the net payment of the OIS with fixed rate k at time T_{n+1} is given by

$$\tau_n (\bar{L}_n - k),$$

where the floating rate \bar{L}_n for the n -th period of OIS is given by

$$\bar{L}_n = \frac{1}{\tau_n} \left(\prod_{i=1}^{K_n-1} (1 + (t_{i+1} - t_i)L(t_i, t_i, t_{i+1})) - 1 \right). \quad (5.7)$$

Here we used the notation $L(t_i, t_i, t_{i+1})$ to denote the overnight rate. Equating the overnight rate with the short rate, we can use a more mathematically convenient (although not exactly correct) expression

$$\bar{L}_n = \frac{1}{\tau_n} \left(e^{\int_{T_n}^{T_n+1} r(t) dt} - 1 \right). \quad (5.8)$$

5.6 Libor-in-Arrears Swaps

Allowing the fixed rate and the notional to vary through time is not the only way to generalize a swap. For a *Libor-in-arrears swap*, Libor rates are observed (fixed) at the end of each period rather than at the beginning. Thus, a value of a Libor-in-arrears payer swap is equal to

$$V_{\text{LIA}}(t) = \beta(t) \sum_{n=0}^{N-1} \tau_n E_t \left(\beta(T_{n+1})^{-1} (L_{n+1}(T_{n+1}) - k) \right).$$

Interestingly, this seemingly innocuous modification makes the value of a swap model-dependent, in contrast to the standard fixed-floating swap. We will discuss pricing of in-arrears swaps in Chapter 16.

Libor-in-arrears swaps are popular in upward-sloping interest rate curve environments, i.e. when long-tenor rates are higher than shorter-tenor ones. In such a scenario, the break-even fixed rate on the Libor-in-arrears swap tends to look more “attractive” than that of a standard fixed-floating swap, thus increasing the desirability of the swap to those seeking to receive fixed rate payments.

5.7 Averaging Swaps

Libor rates are not restricted to being observed on either the start date or the end of the pay period. A popular example is the *averaging swap*, i.e. a swap where the floating rate is determined as an average of Libor rate observations taken at regular intervals over each coupon period. For example, let $\{(t_{n,i}^f, t_{n,i}^s, t_{n,i}^e)\}_{i=1}^{K_n}$ be a collection of date triplets (fixing, start and end date) that define the rates to be used in calculating the payment in period n . Defining a set of weights $w_{n,i}$, $i = 1, \dots, K_n$, the floating rate \bar{L}_n for the period $[T_n, T_{n+1}]$ may be defined as

$$\bar{L}_n = \sum_{i=1}^{K_n} w_{n,i} L(t_{n,i}^f, t_{n,i}^s, t_{n,i}^e).$$

For the fixed rate swap payer, the averaging swap value is therefore

$$V_{\text{average}}(t) = \beta(t) \sum_{n=0}^{N-1} \tau_n E_t \left(\beta(T_{n+1})^{-1} (\bar{L}_n - k) \right). \quad (5.9)$$

As a rule, the weights $w_{n,i}$ sum up to 1, $\sum_{i=1}^{K_n} w_{n,i} = 1$; the weights usually reflect the number of days (using the appropriate day counting conventions) that a given rate $L(t_{n,i}^f, t_{n,i}^s, t_{n,i}^e)$ is supposed to be in effect. Computation of the valuation expression (5.9) can be done using techniques similar to those required for in-arrears swaps; see Chapter 16 for details.

Swaps linked to the average of the Federal funds rate are common examples of an averaging swap. Particularly noteworthy is the *Fed funds/Libor basis swap* which pays the average of the Fed funds rate (over a given period) against a payment based on a Libor rate for that period. This instrument is an example of a *floating-floating single-currency basis swap*, i.e., a swap that exchanges payments based on two different floating rates in the same currency. Closely related to Fed funds basis swaps are the *Fed funds futures* contracts traded on the Chicago Board of Trade (CBOT) exchange. These contract uses the 30 day running average of the Federal funds rate for settlement.

Remark 5.7.1. Going forward, in our product descriptions we shall normally assume that all cash flows pay at the end of the periods in which they fix. While this is common practice, as we have just seen the “pay-in-arrears” rule can be broken at will depending on the client’s needs — the only (self-evident) restriction is that payments should be fixed by the time they are made.

5.8 Caps and Floors

A firm with liabilities funded at a floating (i.e., Libor) rate is naturally concerned with the possibility that interest rates, and thus its interest rate

payments, may increase in the future. One way to immunize against this risk is to pay fixed on a fixed-floating interest rate swap, in effect turning floating rate payments into fixed ones. While this will guarantee a fixed rate for funding payments for the duration of the swap, it will also mean forgoing the possibility of benefiting from a potential future drop in rates. An *interest rate cap* is a security that allows one to benefit from low floating rates yet be protected from high rates. Similarly, for an investor with assets earning a floating rate, a low-rate scenario is unfavorable. An *interest rate floor* is an instrument designed to protect against low interest rates yet allow the holder to benefit from high rates.

Formally, a cap is a strip of *caplets*, call options on successive Libor rates, and a floor is a strip of *floorlets*, put options on successive Libor rates. We encountered caplets already in Section 4.5.1 and recall that this instrument pays

$$\tau_n (L_n(T_n) - k)^+$$

per unit notional at time T_{n+1} . Similarly, a floorlet pays

$$\tau_n (k - L_n(T_n))^+$$

per unit notional at time T_{n+1} . Then, N -period caps and floors have values at time t of

$$V_{\text{cap}}(t) = \beta(t) \sum_{n=0}^{N-1} \tau_n E_t \left(\beta(T_{n+1})^{-1} (L_n(T_n) - k)^+ \right),$$

$$V_{\text{floor}}(t) = \beta(t) \sum_{n=0}^{N-1} \tau_n E_t \left(\beta(T_{n+1})^{-1} (k - L_n(T_n))^+ \right).$$

By switching to the T_{n+1} -forward measure (see Section 4.2.2) for the n -th caplet/floorlet, the valuation formulas can be written in a more convenient form

$$V_{\text{cap}}(t) = \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) E_t^{T_{n+1}} \left((L_n(T_n) - k)^+ \right),$$

$$V_{\text{floor}}(t) = \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) E_t^{T_{n+1}} \left((k - L_n(T_n))^+ \right).$$

By Lemma 4.2.3, the Libor rate $L_n(\cdot)$ is a martingale under the T_{n+1} -forward measure. Hence, caplets/floorlets can be priced using “vanilla” models⁹, such as the log-normal Black model (see Remark 1.9.4).

⁹By a *vanilla* model we mean a model that specifies the dynamics (or just the terminal distribution) of only a single rate, or at most a few rates, in contrast to term structure models that specify consistent dynamics for the entire term structure of interest rates. Often vanilla models are borrowed from equity or FX modeling; having the underlying rate a martingale makes such borrowing painless. We discuss vanilla models in Chapters 7, 8, 9 and 17.

The OTC market in caps/floors is very liquid. While individual caplets/floorlets are not traded, caps/floors are available in a number of maturities. This allows the volatility information for individual forward Libor rates to be extracted from market quotes for caps/floors of different maturities, at least in principle¹⁰. Once extracted, these volatilities may be combined with the volatilities observed from European swaption quotes (see below), to form a set of market inputs to which interest rate models for exotics are calibrated.

5.9 Digital Caps and Floors

Digital caps and floors work like regular caps and floors, except that the n -th digital caplet pays

$$\tau_n \times 1_{\{L_n(T_n) > k\}}.$$

Similarly, the n -th digital floorlet pays

$$\tau_n \times 1_{\{L_n(T_n) < k\}}.$$

Digital caps and floors provide a leveraged way to bet on the future direction of interest rates, more so than through standard caps and floors.

5.10 European Swaptions

Caps and floors have an asymmetric exposure to interest rates, a characteristic used by both hedgers and speculators. A similar exposure profile is provided by options on swaps, the so-called *European swaptions*. A European swaption gives the holder a right, but not an obligation, to enter a swap at a future date at a given fixed rate. A *payer* swaption is an option to pay the fixed leg on a fixed-floating swap; a *receiver* swaption is an option to receive the fixed leg.

Assuming the underlying swap starts on the expiry date T_0 of the option (a typical situation), the payoff for a payer swaption at time T_0 then equals

$$V_{\text{swaption}}(T_0) = (V_{\text{swap}}(T_0))^+ = \left(\sum_{n=0}^{N-1} \tau_n P(T_0, T_{n+1}) (L_n(T_0) - k) \right)^+. \quad (5.10)$$

The value at an intermediate time t , $t < T_0$, must then equal

$$\begin{aligned} V_{\text{swaption}}(t) &= \beta(t) E_t (\beta(T_0)^{-1} V_{\text{swaption}}(T_0)) \\ &= \beta(t) E_t \left(\beta(T_0)^{-1} \sum_{n=0}^{N-1} \tau_n P(T_0, T_{n+1}) (L_n(T_0) - k) \right)^+, \end{aligned}$$

¹⁰This “volatility bootstrap” is by no means trivial; we discuss it in Section 16.2.

which, using (5.6), can be rewritten in the more compact form

$$V_{\text{swaption}}(t) = \beta(t) E_t \left(\beta(T_0)^{-1} A(T_0) (S(T_0) - k)^+ \right). \quad (5.11)$$

Moreover, switching to the annuity measure, also known as the *swap measure*, Q^A from Section 4.2.5, the swaption value can be expressed as

$$V_{\text{swaption}}(t) = A(t) E_t^A (S(T_0) - k)^+, \quad (5.12)$$

with the forward swap rate $S(\cdot)$ being a martingale in the swap measure Q^A ; see Lemma 4.2.4.

It is evident from (5.12) that a payer European swaption is a call option — and a receiver European swaption is a put option — on the forward swap rate, struck at the fixed rate of the swap. Hence, swaptions could be priced using a vanilla model (see footnote 9), such as the Black model or similar. Conversely, values of European swaptions can be translated into market-implied distributional characteristics of forward swap rates, a topic discussed at length in Section 7.1.2. In particular, it is universal practice to quote swaption prices in terms of *implied* Black volatilities, i.e. volatilities that recover market price when used in the Black formula. In some markets (e.g., the US), it is also common to quote implied *Gaussian* volatilities, defined in the same way with regard to a Gaussian (rather than log-normal) model for the distribution of interest rates, see (7.16).

The market in swaptions is very liquid, with many different option maturities and swap underlyings actively traded. To characterize the full universe of traded instruments, given a tenor structure (5.3) we consider swaptions of different expiries $\{T_n\}_{n=0}^{N-1}$ that can be exercised into swaps that start at T_n and cover m periods¹¹, i.e. their last payment date is T_{n+m} . For a convenient way to denote the various swaptions, recall definitions (4.8), (4.10) and (4.11) and introduce

$$A_{n,m}(t) = \sum_{i=n}^{n+m-1} \tau_i P(t, T_{i+1}), \quad (5.13)$$

$$S_{n,m}(t) = \frac{\sum_{i=n}^{n+m-1} \tau_i P(t, T_{i+1}) L_i(t)}{\sum_{i=n}^{n+m-1} \tau_i P(t, T_{i+1})}, \quad (5.14)$$

for $n = 0, \dots, N-1$, $m = 1, \dots, N-n$. Then the value of the (n, m) -swaption (a short-hand for an “ m -period swaption with expiry T_n ”) is equal to

$$A_{n,m}(t) E_t^{n,m} \left((S_{n,m}(T_n) - k)^+ \right),$$

where $E_t^{n,m}$ denotes time t expectation in the appropriate swap measure, $Q^{n,m}$. Note that in trader parlance, a (vanilla) T_n -maturity European swaption on a swap that runs from T_n to T_m is said to be a “ T_n into $T_{m+n} - T_n$ ”

¹¹A bit confusingly, such a swaption is often said to have *tenor* $T_{n+m} - T_n$, a characterization it inherits from the underlying swap rate.

swaption. For instance, a 5 year option on a 10 year swap would be a “5-into-10” (or “5y-into-10y”, or simply “5y10y”) swaption.

Clearly, when $m = 1$, the (n, m) -swaption reduces to a caplet (or floorlet) on the Libor rate $L_n(\cdot)$, so caplets and floorlets can be thought of as one-period swaptions. Whenever in this book swaptions are discussed or used, caplets and floorlets are thereby implicitly included. Collectively, all (n, m) -swaptions constitute *the swaption grid*.

Market quotes on swaptions, typically in terms of implied volatilities, in the swaption grid provide the most readily-available information on the volatility structure of interest rates. As swaptions in the grid cover overlapping sections of the term structure of interest rates, extracting clean volatility information from market quotes is a non-trivial exercise that forms the foundation for calibration of models used for exotic interest rate derivative pricing. We will have much to say about such *volatility calibration* later on.

While options on plain-vanilla swaps comprise the bulk of the liquid (“vanilla”) interest rate market, options on general swaps (i.e. on swaps with non-constant notional and fixed rates) also trade and are properly treated as exotic derivatives. Often, general swaps can be decomposed into baskets of standard swaps, in which case options on general swaps become *basket options*. Valuation of basket options requires information on the co-dependence structure of securities in the basket, information that is not readily available from the vanilla options markets. We demonstrate how to handle this complication in Section 19.4.

5.10.1 Cash-Settled Swaptions

The swaption contract discussed in the previous section involves *physical settlement*, in the sense that an actual interest rate swap is entered into, should the option be exercised at its expiry. Physically-settled swaptions are also known as *swap-settled swaptions*. An economically equivalent swaption contract is one that instead settles into a cash payment equal to the PV (present value) of the swap as observed at time T_0 . Indeed, for both types, the swaption payoff (for a payer) is given by

$$A(T_0)(S(T_0) - k)^+, \quad (5.15)$$

see (5.10) and (5.11). In the European markets, a third variety of swaptions is common, the so-called *cash-settled* swaptions. For this type of option, rather than entering into a swap, the option holder will receive a cash payout upon exercise. The settlement amount is calculated by a formula similar to (5.15), except the annuity $A(\cdot)$ is not calculated by (5.4), but instead by discounting fixed rate payments at the swap rate $S(T_0)$. Specifically,

$$V_{\text{css}}(T_0) = a(S(T_0))(S(T_0) - k)^+,$$

where

$$a(x) = \sum_{n=0}^{N-1} \frac{\tau_n}{\prod_{i=0}^n (1 + \tau_i x)}.$$

Notice that the cash settlement mechanism ensures a well-defined present value of the option payout, as long as the swap rate $S(T_0)$ is observable. In contrast, the value of exercise of a physically settled swaption — the computation of which requires knowledge of a strip of discount factors — may be estimated differently by different dealers, due to bid-ask spread effects and differences in curve building technology (see Chapter 6). Technically, however, the cash settlement mechanism induces certain valuation complications, and cash-settled swaptions cannot, strictly speaking, be considered vanilla options that can be priced using, e.g., a Black-type formula¹². This follows from the fact that in the measure associated with the deflator $X(t) = a(S(t))$, the swap rate $S(\cdot)$ is *not* a martingale, and certain drift adjustments are required. We discuss valuation of cash-settled swaptions in Section 16.6.12. As they are the most liquidly-traded OTC interest rate options in the European market, cash-settled swaptions still could (and should) be used to extract information on the volatility structure of interest rates; the procedure, however, is necessarily more involved.

5.11 CMS Swaps, Caps and Floors

As the market in plain vanilla swaps is both deep and very active, market quotes of corresponding swap rates can be used as “indexes”, i.e. market variables that can themselves be used in defining payoffs of other securities. The demand for such products is often driven by particular segments of fixed income markets. For example, mortgage lenders are primarily concerned with hedging interest rate risk arising from holding residential loans, some of which may have maturities as long as thirty years. Because of potential prepayments, the interest rate risk of a pool of such mortgages is often assumed to be closely connected to movements in the 10 year swap rate; hence, mortgage lenders are natural consumers of interest rate securities linked to the 10 year swap rate.

A constant-maturity swap (CMS) rate is defined as a break-even swap rate (see (5.5)) on a standard swap of a fixed maturity, e.g. 10 years or 30 years. A *CMS swap* works just like a standard fixed-floating (Libor) swap, except for the fact that floating leg payments are based on CMS, rather than Libor, rates. Formally, let $S_{n,m}(\cdot)$ be the m -period swap rate with the first fixing date T_n , as defined by (5.14). Then an m -period (payer) CMS swap’s value is given by

¹²Nevertheless, this practice has been widespread until recently, and may still be in use in some institutions.

$$V_{\text{cmsswap}}(t) = \beta(t) \sum_{n=0}^{N-1} \tau_n E_t \left(\beta(T_{n+1})^{-1} (S_{n,m}(T_n) - k) \right)$$

or, using the T_{n+1} -forward measure for each period,

$$V_{\text{cmsswap}}(t) = \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) E_t^{T_{n+1}} ((S_{n,m}(T_n) - k)).$$

While standard swaps can be valued solely from knowledge of the term structure of interest rates, CMS swaps require an interest rate model for valuation; we return to a complete discussion in Chapter 16.

CMS caps and floors are defined as strips of European options on CMS rates, just like regular caps and floors are strips of European options on Libor rates:

$$\begin{aligned} V_{\text{cmscap}}(t) &= \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) E_t^{T_{n+1}} ((S_{n,m}(T_n) - k)^+), \\ V_{\text{cmsfloor}}(t) &= \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) E_t^{T_{n+1}} ((k - S_{n,m}(T_n))^+). \end{aligned}$$

CMS caplets are related to European swaptions, as both are European-style options on swap rates. The connection between the two types of securities is, however, subtle, as we shall discuss later in this book.

5.12 Bermudan Swaptions

A Bermudan swaption is an option to enter into a fixed-floating swap on any (or any from a given subset) of its fixing dates. For a given tenor structure (5.3), the holder of a standard Bermudan swaption has the right to exercise it on any of the dates $\{T_n\}_{n=0}^{N-1}$. Once exercised on date T_n , say, the option goes away, and the holder enters the swap with the first fixing date T_n and the final payment date T_N . The period up to $T_0 > 0$ is known as the *lockout* or *no-call* period. In common jargon, a Bermudan swaption on, say, a 10 year swap with a 2 year lockout period (at inception) is known as a “10 no-call 2”, or “10nc2”, Bermudan swaption.

Formally, at time T_n , the value of a payer¹³, if exercised, is therefore

$$\begin{aligned} U_n(T_n) &= \beta(t) \sum_{i=n}^{N-1} \tau_i E_{T_n} \left(\beta(T_{i+1})^{-1} (L_i(T_i) - k) \right) \\ &= \sum_{i=n}^{N-1} \tau_i P(T_n, T_{i+1}) (L_i(T_n) - k). \end{aligned}$$

¹³Upon exercise, the holder of a payer (receiver) Bermudan swaption will pay the fixed (floating) leg of the swap.

Here, $U_n(T_n)$ here denotes the *exercise* value of the Bermudan swaption; loosely speaking, a Bermudan swaption contract is an option to chose between $U_n(T_n)$ for different $n = 0, \dots, N-1$. More succinctly, we recall from Section 1.10 that the Bermudan option value at time T_n will be the maximum of $U_n(T_n)$ and the *hold value* $H_n(T_n)$, the latter defined as the value of a Bermudan swaption with the exercise dates $\{T_i\}_{i=n+1}^{N-1}$ only (compare to Sections 1.10 and 3.5).

Demand for Bermudan swaptions comes from different segments of fixed income markets. Mortgage companies use them to hedge pools of mortgages, with the flexibility of Bermudan exercise convenient in matching the uncertain timing of prepayments in mortgage pools. Investors seeking higher current income sell Bermudan-style options on swaps to increase the coupons they receive, as explained later in the context of callable Libor exotics. Bermudan swaptions are also used as hedges for callable coupon bonds.

While it may be tempting to think of Bermudan swaptions as straightforward generalizations of European swaptions, they are substantially more difficult to model and price. Indeed, it is fair to say that many valuation methods and techniques covered in this book were developed in response to the need to value and risk manage Bermudan swaptions. Bermudan swaptions are, by far, the most liquid exotic fixed income securities, with all interest rate dealers holding large inventories.

5.13 Exotic Swaps and Structured Notes

With market sophistication ever on the rise, clients demand increasingly complicated payouts, often in a familiar swap or bond format (although the appetite has waned somewhat post-crisis). In an *exotic swap*, a regular floating Libor leg is swapped against structured coupons that are allowed to be arbitrary functions of observed interest rates (such as Libor or CMS rates). A standard fixed-floating vanilla swap is an obvious and trivial example where the structured coupon simply is a fixed rate. A cap (or a floor) can be seen as another, less trivial, example. In particular, note that

$$\begin{aligned}(k - L_n(T_n))^+ &= \left((k - L_n(T_n))^+ + L_n(T_n) \right) - L_n(T_n) \\ &= \max(k, L_n(T_n)) - L_n(T_n),\end{aligned}$$

which demonstrates that a floor can be represented as an exotic swap in which a Libor rate is exchanged for a *floored payoff* $\max(k, L_n(T_n))$.

Exotic swaps often start their life as bonds, or notes, sold by banks to investors. In a structured note, the investor pays an up-front *principal amount* (e.g., \$10,000,000) to the issuer of the note, who in turn pays the investor a structured coupon, and repays the principal at the maturity of the note. The principal amount is invested by the issuer (or the trading

desk to which the issuer passes the note for risk management), and pays the Libor rate plus or minus a spread. From the perspective of the issuer (or the trading desk), the net cash flows of the note are those of an exotic swap.

In terms of valuation, if C_n is the structured coupon for the n -th period, the value of the exotic swap is equal to (from the perspective of structured leg buyer)

$$V_{\text{exotic}}(t) = \beta(t) \sum_{n=0}^{N-1} \tau_n E_t \left(\beta(T_{n+1})^{-1} (C_n - L_n(T_n)) \right),$$

where we for brevity have assumed that both legs of the swap pay at the end date of each coupon period (see Remark 5.7.1). As discussed earlier, in this valuation equation, the coupon C_n can be a complicated function of interest rates, structured to reflect investors' views on the market, or to take advantage of current interest rate market conditions. For example, a floored payoff can be offered to an investor who believes that interest rates are poised for a fall in the future.

There is no universally agreed "taxonomy" for exotic swaps, but for our purposes we can distinguish between exotic swaps that are i) Libor-based, ii) CMS-based, iii) multi-rate, iv) range accruals, and v) generally path dependent. We proceed to described each type of swap in more details.

5.13.1 Libor-Based Exotic Swaps

In a Libor-based exotic swap, the structured coupon is a function of a Libor rate:

$$C_n = C_n(L_n(T_n)).$$

A large variety of structured coupons $C_n(\cdot)$ can be used. For example:

- A standard swap,

$$C_n(x) = k.$$

- Capped and floored floaters. For strike s , gearing g , cap c and floor f ,

$$C_n(x) = \max(\min(g \times x - s, c), f). \quad (5.16)$$

- Capped and floored inverse floaters. For spread s , gearing g , cap c and floor f ,

$$C_n(x) = \max(\min(s - g \times x, c), f). \quad (5.17)$$

- Digitals. For strike s and coupon k ,

$$C_n(x) = k \times 1_{\{x > s\}}$$

or

$$C_n(x) = k \times 1_{\{x < s\}}.$$

- “Flip-flops” or “tip-tops”. For strike s and two coupons, k_1 and k_2 ,

$$C_n(x) = \begin{cases} k_1, & x \leq s, \\ k_2, & x > s. \end{cases}$$

Different coupon types can be combined together to create new types of structured coupons.

A Libor-based exotic swap can usually be decomposed¹⁴ into a sum of simpler instruments such as ordinary swap floating legs, fixed legs, caps and floors, and digital caps and floors. Therefore, if the prices of these simple contracts are available in the market (as is typically the case), Libor-based exotic swaps can be perfectly replicated by a one-time transaction in market-available instruments, a strategy referred to as *static replication*. Hence, by themselves, these instruments rarely present major valuation challenges. They do, however, serve as building blocks for more complicated securities.

5.13.2 CMS-Based Exotic Swaps

The payoffs from the previous section can be applied to CMS, rather than Libor, rates. Structured coupons are then deterministic functions of CMS rates. If an m -period rate is used, then a structured coupon for period n can be defined by

$$C_n = C_n(S_{n,m}(T_n)),$$

with $C_n(x)$ as defined in the previous section.

CMS-based exotic swaps can be decomposed into linear combinations of CMS swaps and CMS caps/floors and rarely present any extra modeling difficulties beyond those already present in CMS swaps and caps.

5.13.3 Multi-Rate Exotic Swaps

Multi-rate exotic swaps differ from the structures in Sections 5.13.1 and 5.13.2 by referencing multiple market rates (Libor or CMS) for the calculation of structured coupons. The most common example is a *CMS spread coupon*. To describe this contract, let $S_{n,a}(\cdot)$ and $S_{n,b}(\cdot)$ be two collections of CMS rates, fixing on T_n , $n = 0, \dots, N - 1$, and covering a and b periods, respectively. A CMS spread coupon with gearing g , spread s , cap c and floor f is then defined by

$$C_n = \max(\min(g \times (S_{n,a}(T_n) - S_{n,b}(T_n)) + s, c), f).$$

A typical example would be a 10 year/2 year (often abbreviated as 10y2y) CMS call spread option where a is 40 (40 quarterly periods to cover 10 years) and b is 8, with the quarterly coupon given by

¹⁴Indeed, this is the case for all the payouts listed above, a fact that we invite the reader to verify.

$$C_n = \max(S_{10y}(T_n) - S_{2y}(T_n), 0)$$

(using somewhat loose notation). A relatively liquid broker market exists for spread options on Euro and US dollar CMS rates.

A more general example is obtained by using one of the payoff functions $C_n(x)$ defined in Section 5.13.1, applied to the spread $x = S_{n,a}(T_n) - S_{n,b}(T_n)$. In particular, digital and flip-flop CMS spread swaps are quite popular.

Multi-rate exotic swaps typically cannot be decomposed into “standard” instruments (such as vanilla swaps, caps, etc.). Therefore, they, as a rule, cannot be valued by replication arguments, and a valuation model is required. Such a model, however, does not always need to be a full-blown term structure model: we shall show later that some types of spread-linked payoffs can be efficiently valued and risk managed by vanilla models (see footnote 9).

It should be noted that more than two rates can be used in the definition of a coupon. For example, in the so-called *curve cap* one takes a standard capped and floored payoff on a Libor or CMS (or CMS spread!) rate — see (5.16) — and makes the cap c and the floor f functions of, potentially different, CMS spreads:

$$\begin{aligned} C_n(x) &= \max(\min(g \times x - s, c), f), \\ c &= \max(\min(g_1 \times (S_{n,a_1}(T_n) - S_{n,b_1}(T_n)) + s_1, c_1), f_1), \\ f &= \max(\min(g_2 \times (S_{n,a_2}(T_n) - S_{n,b_2}(T_n)) + s_2, c_2), f_2). \end{aligned} \quad (5.18)$$

5.13.4 Range Accruals

A *range accrual* structured coupon is defined as a given rate — fixed in the simplest case, but potentially a Libor, CMS or a CMS spread rate — that only “accrues” when a different *reference* rate is inside (or, sometimes, outside) a given range. So, let $R_n(t)$ be the payment rate and $X_n(t)$ be the reference rate, and let l be the low bound and u be the upper bound. A range accrual coupon then pays

$$C_n = R_n(T_n) \times \frac{\#\{t \in [T_n, T_{n+1}] : X_n(t) \in [l, u]\}}{\#\{t \in [T_n, T_{n+1}]\}}, \quad (5.19)$$

where $\#\{\cdot\}$ is used to denote the number of days that a given criteria is satisfied.

The most common choice of the payment rate $R_n(t)$ is either a constant or Libor, but a CMS rate or any other structured coupon rate are also occasionally used. The reference rate $X_n(t)$ can be any market-observable rate such as a Libor rate fixing at t , a CMS rate fixing at t , or even a CMS spread rate.

We note that a range accrual coupon can always be decomposed into simpler digital payoffs, because

$$\#\{t \in [T_n, T_{n+1}] : X_n(t) \in [l, u]\} = \sum_{t \in [T_n, T_{n+1}]} 1_{\{X_n(t) \in [l, u]\}}, \quad (5.20)$$

where the sum on the right-hand side is over all business days in the period. This decomposition is particularly useful for fixed rate ($R_n(T_n) \equiv k$) range accruals, as simple digital options can be priced directly from the market information on European options (see Section 7.1.2). For floating, or more complicated, range accruals the decomposition is useful but requires further work to turn it into valuation formulas — see Section 17.5 for further details.

The basic payout (5.19) can be extended to include more than one range condition. In a *dual range accrual*, the position of two different reference rates relative to the range are monitored, and (5.19) is generalized to

$$C_n = R_n(T_n) \frac{\#\{\{t \in I_n : X_{n,1}(t) \in [l_1, u_1]\} \diamond \{t \in I_n : X_{n,2}(t) \in [l_2, u_2]\}\}}{\#\{t \in I_n\}}$$

with $I_n = [T_n, T_{n+1}]$ and \diamond denoting either intersection \cap or union \cup . In the former case, one counts the number of days when *both* reference rates $X_{n,1}$ and $X_{n,2}$ are within their ranges; in the latter case, one counts the number of days when *either* of the two reference rates are within their ranges.

In a *curve cap range accrual*, the lower and upper bounds become functions of CMS spreads themselves, similar to (5.18).

A *product-of-ranges* range accrual multiplies up all range accrual factors to date to define the multiplier that is used for the current coupon, e.g.

$$C_n = R_n(T_n) \times Y_n,$$

where

$$Y_n = Y_{n-1} \times \frac{\#\{t \in [T_n, T_{n+1}] : X_n(t) \in [l, u]\}}{\#\{t \in [T_n, T_{n+1}]\}}, \quad (5.21)$$

$$Y_{-1} = 1.$$

5.13.5 Path-Dependent Swaps

The payoff (5.21) is an example of *path-dependence* in the payoff, where a coupon depends on rate observations from previous coupon periods. More commonly, path-dependence in exotic swaps is introduced by linking a structured coupon not only to interest rates observed during the coupon period, but to previous coupon(s) as well. This is often referred to as a *snowball* feature. The “original” snowball structure involved a coupon of an inverse floating type, with the n -th coupon C_n given by

$$C_n = (C_{n-1} + s_n - g_n \times L_n(T_n))^+. \quad (5.22)$$

Here $\{s_n\}$ and $\{g_n\}$ are contractually specified deterministic sequences of spreads and gearings. This type of a swap is sometimes also called a *ratchet* or a *ladder* swap.

The term snowball originates with the tendency of high initial coupons to spill into subsequent coupons, in a compounding or “snowballing” fashion”. Indeed, a little reflection reveals that a snowball has a highly leveraged exposure to its first few coupons, a feature that makes it more attractive to some investors, but also quite difficult to risk manage.

A large number of snowball-like payoffs have been created, often with “snow”-themed — and rather nonsensical — names, such as “snowrange”, “snowbear”, and “snowstorm”. For example, a “snowrange” combines a range-accrual feature and a snowball feature, in the following way

$$C_n = C_{n-1} \times Y_n + s_n + g_n \times X_{n,1}(T_n), \quad (5.23)$$

where $X_{n,1}(T_n)$ is some reference rate, and Y_n is a range-accrual factor depending on a second rate $X_{n,2}$,

$$Y_n = \frac{\#\{t \in [T_n, T_{n+1}] : X_{n,2}(t) \in [l, u]\}}{\#\{t \in [T_n, T_{n+1}]\}}.$$

The range accrual factor Y_n may be a product-of-ranges accrual factor, as in (5.21). Also, additional caps and/or floors are often added to the coupon (5.23).

Path-dependent swaps typically require a term structure model for valuation; for obvious reasons, Monte Carlo methods are often mandatory.

5.14 Callable Libor Exotics

5.14.1 Definitions

As described in Section 5.12, Bermudan swaptions are Bermudan-style options to enter a regular fixed-floating swap. If we alter the swap underlying the Bermudan swaption from a regular swap to an exotic swap (see previous section), then a so-called *callable Libor exotic* (CLE) is created. CLEs most often emerge as part of callable structured notes in which an issuer receives the principal from an investor and pays a structured coupon in return. In addition, the issuer has the right to cancel — or *call* — the note on a schedule of dates; typically, this call schedule will coincide with coupon fixing dates, after some initial lock-out (or *no-call*) period. Should a note be called by its issuer¹⁵, the principal is returned to the investor and no future coupons are paid.

A callable structured note is typically passed through by the issuer to an exotics trading desk (which could, but does not have to, be internal to the issuing bank) to deal with its risk management. Also, the principal is

¹⁵The call decision is most often made by the issuer’s swap counterparty who is actually managing the risk, see next paragraph.

invested and pays a Libor rate, plus or minus a spread depending on the cost of financing. From a trading desk perspective what is left is an exotic swap paying structured coupons and receiving Libor, plus a Bermudan-style right to cancel the swap. For clarity, Figures 5.1–5.3 list the cash flow diagrams of a callable structured note¹⁶.

Fig. 5.1. Callable Note: Flows at Inception

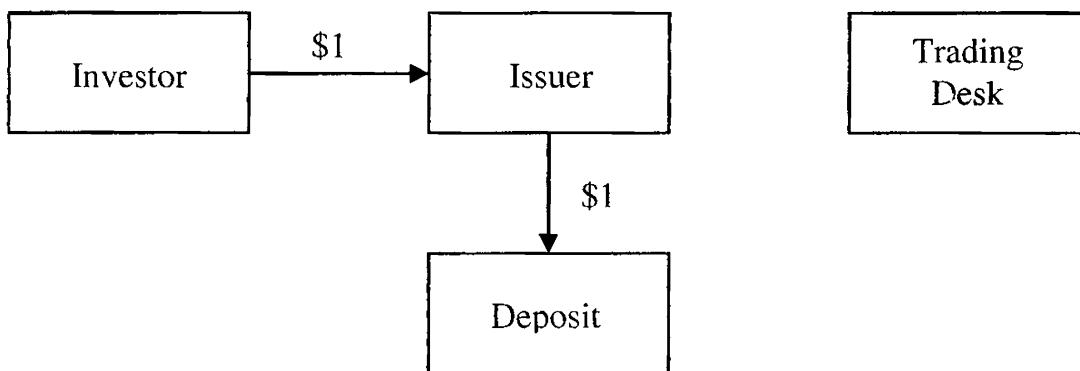
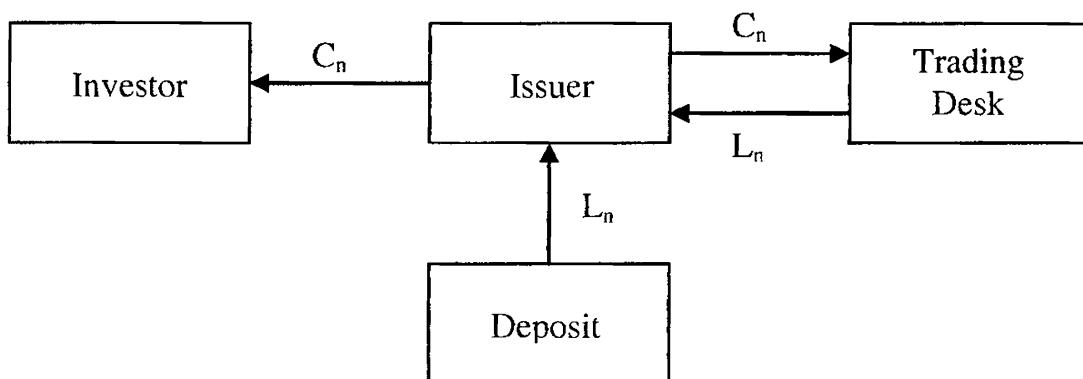


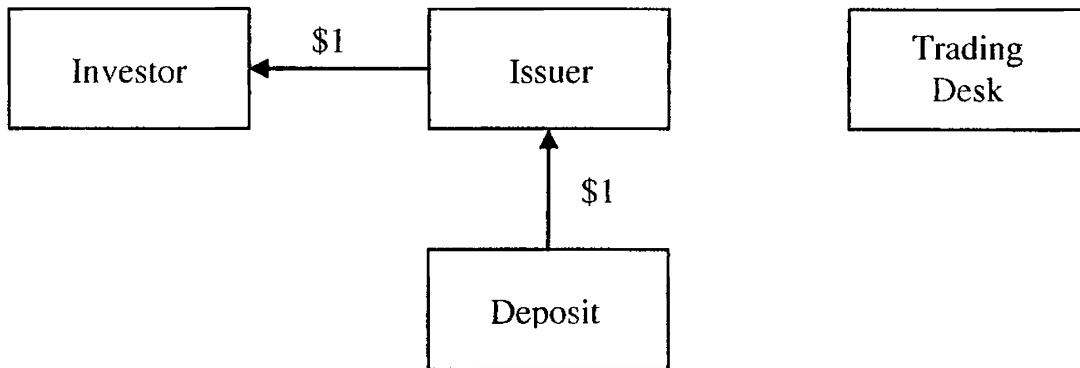
Fig. 5.2. Callable Note: Flows at Payment Times



Sometimes it is convenient to represent a cancelable exotic swap as a straight exotic swap, plus a Bermudan-style option to enter a reverse swap, i.e. a swap where legs are reversed relative to the original one. Beyond providing a break-down that is convenient for valuation purposes, this representation emphasizes the fact that the cancelability feature of a CLE benefits the party that owns it (typically a structured note issuer). Indeed, the feature is

¹⁶While, conceptually, the principal is deposited into a Libor-paying account, in practice it is used as part of cash management activities by the issuer. A structured note issuance program often provides cheaper funding to a bank than would be attainable by other means.

Fig. 5.3. Callable Note: Flows at Termination (Maturity or early Cancellation)



often added to a structured note as a way to offer a more attractive coupon to the investor, in return for the Bermudan-style option. Often, the coupons inside the non-call period are fixed rate coupons, and a typical way for the issuer to “pay” for the Bermudan option is to make these first coupons high, often much higher than the return available elsewhere. This “optical illusion” of high rate of return on investment is, at least in part, what drives investor interest in structured callable notes.

Consider a CLE on an exotic swap with structured coupons $\{C_n\}_{n=0}^{N-1}$. As for regular Bermudan swaptions, we denote the value of the exotic swap that one can exercise on date T_n by

$$U_n(T_n) = \beta(T_n) \sum_{i=n}^{N-1} \tau_i E_{T_n} \left(\beta(T_{i+1})^{-1} \times (C_i - L_i(T_i)) \right). \quad (5.24)$$

Here, we recall that from a trading desk prospective, the cancelability feature of CLE involves an option to enter a reverse swap, with receipt of structure coupons and payment of Libor. Hold values are also defined analogously to the Bermudan swaption case: the n -th hold value $H_n(T_n)$ is defined as the time T_n value of the CLE on the same exotic swap, but with exercise dates $\{T_i\}_{i=n+1}^{N-1}$ only. That is, $H_n(T_n)$ is the time T_n value of the CLE provided it has not been exercised on or before T_n .

5.14.2 Pricing Callable Libor Exotics

A significant part of this book is dedicated to efficient methods for pricing and risk-managing callable Libor exotics. To provide a brief preview of the difficulties involved, we notice that the call feature embedded in CLEs may suggest application of PDE methods, using backward induction arguments outlined in previous chapters. However, often a CLE has explicit path dependence that makes the application of PDE methods impractical. In other cases, as we shall describe in greater depth later in the book, models that admit an efficient PDE representation are often too inflexible for

application to anything but the simplest of CLEs. Hence, CLEs more often than not must be valued using Monte Carlo methods. We have seen a preview of how optimal exercise can be handled in Monte Carlo in Section 3.5; many more details are provided in Chapter 18.

5.14.3 Types of Callable Libor Exotics

Any exotic swap can be used as an underlying for a callable Libor exotic. For our purposes, the taxonomy of callable Libor exotics can follow closely that of exotic swaps, see Section 5.13. We can thereby distinguish various types of CLEs, e.g. Libor-based, CMS-based, multi-rate, callable range accruals, callable snowballs, and so on. Many variations on the basic CLE design exist, most of which are driven by a desire to increase the value of the option to cancel that the investor sells the trading desk, in order for a higher coupon to be paid. It is difficult to classify all the features that have been invented: we content ourselves with merely listing some of the more popular ones.

5.14.4 Callable Snowballs

A callable snowball is a CLE with a snowball (or snowrange, etc.) underlying. From a modeling prospective, they are notable for being one of the first widely popular instruments that combine both strong path-dependence and optimal exercise. It is possible to incorporate snowball-type path-dependence into a PDE framework by introducing auxiliary variables, following the principles of Section 2.7.5; Section 18.4.5 discusses details specific to snowballs. Alternatively — and often preferably — optimal exercise can be incorporated into the Monte Carlo method, as discussed in Section 3.5 and later on in Chapter 18.

5.14.5 CLEs Accreting at Coupon Rate

Typically the notional of the underlying swap of a CLE is fixed throughout the life of the deal, but it does not have to be. For instance, it is not uncommon for the notional to vary deterministically, e.g. increase or decrease by non-random additive or multiplicative amounts each coupon period. Such deterministic accretion/amortization rarely adds extra complications from a modeling prospective. Occasionally, however, a contract specifies that the notional of the swap accretes at the structured coupon rate, in which case the accretion rate will be random. For such CLEs, the exercise value in (5.24) must be amended. Specifically, if q_i is the notional in place for the period $[T_i, T_{i+1}]$, then q_i is obtained from the notional over the previous period q_{i-1} by multiplication with the structured coupon over the previous period. Formally,

$$U_n(t) = \beta(t) \sum_{i=n}^{N-1} \tau_i E_t \left(\beta(T_{i+1})^{-1} \times q_i \times (C_i - L_i(T_i)) \right),$$

$$q_i = q_{i-1} \times (1 + \tau_{i-1} C_{i-1}),$$

where the initial notional q_0 is contractually specified.

Interestingly, the random accretion feature above can be incorporated into a PDE-based scheme without any extra cost, see Section 18.4.5.

5.14.6 Multi-Tranches

The more optionality an investor can sell to the issuer, the better coupon she can receive. As described earlier, the option to call the note is already present in a callable structured note. Another option that is sometimes embedded is a right for the issuer to increase the size of the note, i.e. to put more of the same note to the investor, whether she wants it or not. The name of this feature, a “multi-tranche” callable structured note, originates with the fact that these possible notional increases are formalized as tranches¹⁷ of the same note that the issuer has the right to put to the investor. The times when the issuer has the right to increase the notional of the note typically come before the times when the note can be canceled altogether. Callability usually applies jointly to all tranches of the note.

By itself, the multi-tranche feature rarely presents modeling issues, although one must be mindful of it and plan for a pricing infrastructure that is flexible enough to handle it.

5.15 TARNs and Other Trade-Level Features

While sometimes the precise split is a little arbitrary, it is often helpful to think of a Libor exotic as being defined by

- A definition of its coupon, i.e. a formula that converts rates observed during a coupon period (and, sometimes, previous coupons) into the amount of money paid to the investor.
- A collection of trade-level features, i.e. features that cannot conveniently be expressed as coupon definitions, but instead “act” on the whole trade.

We have already seen examples of both features. For instance, a callable snowball CLE has a coupon definition given by the formula (5.22) on top of which callability has been added as a trade-level feature. In the next few sections we review some other trade-level features.

¹⁷“Slices” in French; here meaning “similar securities offered as part of the same transaction”.

5.15.1 Knock-out Swaps

A knock-out swap is just an exotic swap that disappears on the first fixing date on which some reference rate is above (or below) a given barrier. If the knock-out rate for the period n is denoted by $X_n(t)$, the coupon by C_n , the Libor rate by L_n , and the knock-out barrier by R , the value of a down-and-out¹⁸ knock-out swap is given by

$$V_{\text{KO}}(t) = \beta(t) E_t \left(\sum_{n=0}^{N-1} \tau_n \beta(T_{n+1})^{-1} \times (C_n - L_n(T_n)) \times \prod_{i=0}^n 1_{\{X_i(T_i) > R\}} \right).$$

5.15.2 TARNs

Callable structured notes have proved to be popular with investors, but suffer from the drawback that investors rarely know when the issuer will call the note — indeed, the decision to exercise a Bermudan-style option is driven by a model rarely accessible to the average investor. A relatively recent innovation, the *Targeted Redemption Note* (TARN), presents one possible solution to the problem. In a TARN (see Piterbarg [2004c]), the total investor return, defined as the sum of all structured coupons paid to date, is recorded over time. When the total return exceeds a pre-specified target level— hence the name of the structure — the note is terminated and the principal is returned to the investor.

As with callable notes, issuers do not keep TARN structures on their books, but swap them out with a trading desk. Since the principal payment from investors is invested at the Libor rate, to a trading desk a TARN looks like an exotic swap that knocks out on the total sum of structured coupons. Formally, let the structured coupon¹⁹ for the period $[T_n, T_{n+1}]$ be C_n . The coupon over the period $[T_n, T_{n+1}]$ is only paid if the sum of structured coupons up to (and not including) time T_n is below a total return R . Thus, the value of the TARN at time 0 from the investor's viewpoint is given by

$$V_{\text{tarn}}(t) = \beta(t) E_t \left(\sum_{n=1}^{N-1} \tau_n \beta(T_{n+1})^{-1} \times (C_n - L_n(T_n)) \times 1_{\{Q_n < R\}} \right), \quad (5.25)$$

$$Q_n = \sum_{i=1}^{n-1} \tau_i C_i, \quad Q_1 = 0.$$

¹⁸I.e. disappearing upon some variable breaching a barrier from above; compare to up-and-out options discussed in Chapter 2.

¹⁹In the original TARN product, an inverse floating coupon (5.17) was used, but any structured coupon can be employed.

We note that a TARN typically pays some fixed coupons to an investor before the knock-out feature starts, mirroring the non-call structure of CLEs, see Section 5.14. We omit these from the contract description as their present value can be computed statically off an interest rate curve separately, as the payments are known in advance.

Various features can be added to the description of the TARN we just described. For example, the last coupon, i.e. the coupon that pushes the total return over the target R , can be paid only partially to make the total return exactly R and not more. This feature is known as a *cap at trigger* or a *lifetime cap*. Also, if the total return over the life of the TARN never reaches the target R , a TARN equipped with a so-called *lifetime floor* will make a *make whole* payment at the TARN maturity to ensure that the total return exactly equals R , and not less. These features do not alter the general modeling framework for TARNs that we develop in later chapters, and we shall generally ignore them.

While it could be argued that a TARN is really just a swap with a different coupon definition (namely, $C_n \times 1_{\{Q_n < R\}}$)²⁰ we prefer to classify a TARN feature as trade-level, reflecting its historic importance and its relationship with the callability feature.

5.15.3 Global Cap

As discussed, a TARN can have a feature that restricts total return to an investor to be exactly the trigger level R . This feature can be decoupled from the TARN definition and used by itself, often called a *global cap*. Specifically, an exotic swap with a global cap R pays a structured coupon C_n to an investor until the sum of the coupons has reached the level R . Note that a swap typically does not terminate at this point, i.e. the trading desk will continue to receive Libor until the maturity of the trade or some other agreed termination event.

5.15.4 Global Floor

A *global floor* guarantees the investor a minimum cumulative return of the note. Specifically, if the sum of structured coupons paid to the investor does not reach a global floor value of F by the maturity of the deal, the issuer will pay a make whole amount to the client, equal to F minus the sum of actual coupons paid. The payment is made at the maturity of the swap or some other termination event (such as when the swap is canceled, in case of a callable global floor note).

²⁰The same could be said about knock-out swaps previously defined.

5.15.5 Pricing and Trade Representation Challenges

Trade-level features are often combined with each other, and with various coupon formulas. For example, one can be asked to price a callable, TARN’ed snowball with a global cap. As we shall argue in later chapters, the only modeling solution that is sufficiently scalable to accommodate such arbitrary combinations of various trade features involves a generic, flexible model that is calibrated to a broad collection of market volatility information. In such a framework, adding extra features to a trade is ultimately not much of a modeling problem, but could be a significant trade representation challenge. While outside the scope of this book, such a challenge should be addressed by a software framework that is flexible enough to represent any, current or future, trade-level features, and incorporate them into a pricing engine without significant extra effort. A successful implementation of such a trade representation framework requires careful planning and considerable investment. One fairly common route is to use a domain-specific programming language for trade representation, see for example Jones et al. [2000] for a commercially available version or Frankau et al. [2009] for an example of an in-house one.

5.16 Volatility Derivatives

In a nutshell, a *volatility derivative* is a contingent claim whose underlying is the volatility of a financial observable, rather than a financial observable itself. The simplest example of such a derivative is the *variance swap* (see Carr and Lee [2009b]), a structure that first emerged in equity and foreign exchange trading. In the last few years, similar ideas and structures have entered the fixed income derivative arena.

The demand for volatility derivatives in interest rates is driven by the same factors as in other asset classes; a common motivation is the desire of some market participants — hedge funds in particular — to have direct exposure to interest rate volatility, but not to the outright level of rates, say. In other cases, the product development follows the usual path of creating structured notes with appealing payoff profiles.

Different interest rates constitute different financial observables for defining a volatility, hence one needs to be rather specific when defining a volatility-linked payoff.

5.16.1 Volatility Swaps

A *volatility swap* in interest rates is a contract that measures realized volatility (or a related quantity) of a given rate, although it is structured somewhat differently from volatility swaps in equity or FX markets. Let

$X_n(t)$ be the rate used for period n ; then the most common coupon of a volatility swap is given by

$$C_n = |X_{n+1}(T_{n+1}) - X_n(T_n)|,$$

or a capped version

$$C_n = \min(|X_{n+1}(T_{n+1}) - X_n(T_n)|, c).$$

The value of the (structured) leg of a volatility swap measures realized variation of the rate $X_n(\cdot)$,

$$\begin{aligned} V_{\text{volswap}}(t) &= \beta(t) E_t \left(\sum_{n=0}^{N-1} \tau_n \beta(T_{n+1})^{-1} \times |X_{n+1}(T_{n+1}) - X_n(T_n)| \right) \\ &\quad + V_{\text{floatleg}}(t), \end{aligned} \tag{5.26}$$

where

$$V_{\text{floatleg}}(t) \triangleq \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}) L_n(t) = 1 - P(0, T_N).$$

There are two common choices for the rate X_n . One choice, a *fixed-tenor* volatility swap, involves a swap rate of the same tenor on each of the fixing dates. Technically speaking, *different* swap rates are therefore used for different periods,

$$X_n(t) = S_{n,m}(t)$$

with a fixed value of m , the number of periods in the swap rate (see the definition (5.14)). For example, a rolling 10 year CMS rate could be used. The other choice, a *fixed-expiry* volatility swap, specifies the swap rate to have a fixed expiry and tenor, i.e.

$$X_n(t) = S_{K,m}(t).$$

With this definition, the volatility swap pays the absolute variation of a rate with the fixing date T_K and spanning m periods of the tenor structure $\{T_n\}_{n=0}^{N+m}$. Often $K = N$, so that the variability of the rate $S_{N,m}$ is measured over the whole of its life.

Recently, volatility swaps on CMS spread rates have appeared. As the name implies, they measure the variation of the spread of two rates, e.g. $X_n = S_{n,m_1} - S_{n,m_2}$.

5.16.2 Volatility Swaps with a Shout

Sometimes, the investor in a volatility swap is given an option to *shout*, that is to choose when the fixing of the rate occurs for the purposes of calculating the coupon payoff. In particular, the payoff of the n -th coupon is then

$$C_n = |X_{n+1}(\eta_n) - X_n(T_n)|,$$

where the random stopping time $\eta_n \in [T_n, T_{n+1}]$ is chosen by the investor coupon-by-coupon. For the uncapped version of this payoff, it is intuitively clear that it is always optimal to postpone the shout until the end of the period, i.e. $\eta_n = T_{n+1}$. So, the option given to the investor is actually worthless, while designed to appear to have some value. Interestingly, for the capped version, it is optimal²¹ to shout at the lesser of T_{n+1} and the first time that $|X_{n+1}(\eta_n) - X_n(T_n)| \geq c$, where c is the cap level. As a consequence, a capped volatility swap with a shout option is equivalent to a volatility swap with a barrier on each coupon:

$$\begin{aligned} C_n &= c \times 1_{\left\{\max_{t \in [T_n, T_{n+1}]} |X_{n+1}(t) - X_n(T_n)| \geq c\right\}} \\ &\quad + |X_{n+1}(T_{n+1}) - X_n(T_n)| \times 1_{\left\{\max_{t \in [T_n, T_{n+1}]} |X_{n+1}(t) - X_n(T_n)| < c\right\}}. \end{aligned}$$

This decomposition follows from results in Broadie and Detemple [1995] and is discussed in more detail in Chapter 20. The fact that one can replace an optimal exercise feature with a known static barrier is quite convenient and allows for easy Monte Carlo valuation.

5.16.3 Min-Max Volatility Swaps

The structured coupon for a min-max volatility swap is given by

$$C_n = M_n - m_n,$$

where

$$\begin{aligned} M_n &= \max_{t \in [T_n, T_{n+1}]} X_n(t), \\ m_n &= \min_{t \in [T_n, T_{n+1}]} X_n(t). \end{aligned}$$

The coupon represents the spread between the maximum and the minimum values that a given rate achieves during a coupon period.

While the two products appear quite different at a first glance, there is an interesting connection between min-max and standard volatility swaps. We shall explore this further in Chapter 20.

5.16.4 Forward Starting Options and Other Forward Volatility Contracts

The value of a standard European swaption and a standard fixed-expiry volatility swap are both linked to the volatility of the swap rate over its

²¹Ignoring some small convexity effects and the difference between discrete vs. continuous shout option rights.

entire life, i.e. from the valuation date to the fixing date of the swap rate. Some clients, however, prefer securities that are linked to the volatility of a swap rate as measured over only a sub-period of this time; in effect, the clients want exposure to what is often known as *forward* volatility. Precise definitions and the importance of forward volatility are left for future chapters; for now we content ourselves by listing a few relevant varieties of forward volatility derivatives.

Midcurve swaptions are swaptions whose expiry T^e is strictly before the fixing date T_0 of the underlying swap rate. Their value depend on the volatility of the swap rate over the period $[t, T^e] \subset [t, T_0]$.

Given a swap rate $S(\cdot)$ with the fixing date T_0 , and a date $T^s < T_0$, a *forward-starting swaption straddle*²² is given by the payoff

$$A(T_0) \times |S(T_0) - S(T^s)|$$

paid at T_0 . Essentially, this contract is a combination of a receiver and a payer swaption, both of which will have their strikes fixed at time T^s to the then-prevailing level of the underlying swap rate. That is, the contract pays the value of the at-the-money straddle on the rate $S(\cdot)$ at time T^s with expiry T_0 . The value of a forward-starting swaption straddle is driven by the volatility of the swap rate over the period $[T^s, T_0]$.

Recall that European swaptions are typically quoted in terms of their implied volatilities. Due to this convention, some clients find the forward starting straddle too indirect and instead want to receive implied volatility itself. Particularly popular are contracts that pay implied Normal²³ volatility as defined on p. 204. Fortunately, the at-the-money Normal volatility has a direct relationship to the swaption price, and the payoff of an *implied Normal volatility contract* is

$$\sqrt{\frac{\pi}{2(T_0 - T^s)}} \times |S(T_0) - S(T^s)|$$

paid at T_0 . Apart from the factor $A(T_0)$ and a non-consequential scaling factor, the payoffs of a forward-starting straddle and the implied Normal volatility contract are identical. The differences in their prices are just a matter of a minor convexity correction, a topic we return to in Chapter 20.

²²The term *straddle* is used to denote the sum of a put and a call option with identical strikes.

²³Contracts paying implied log-normal, or Black, volatility are possible but less common, due to the common perception that interest rates are more Gaussian than log-normal, i.e. the implied Gaussian volatility is less sensitive to the changes in the level of interest rates than the implied Black volatility.

5.A Appendix: Day Counting Rules and Other Trivia

In this appendix, we very briefly cover some of the finer details of how schedules are constructed and how interest rate payments accrue under market conventions. We generally ignore these details in the main body of the book, and our treatment here only scratches the surface. For a full account, see Mayle [1993] or Stigum and Robinson [1996].

5.A.1 Libor Rate Definitions

Consider the 6 month Libor rate L fixing at time T . According to (4.2), we would compute this rate as simply

$$L(T) = L(T, T, T + 1/2) = (P(T, T + 1/2)^{-1} - 1) / \tau, \quad \tau = 1/2. \quad (5.27)$$

In reality, this computation ignores a number of quoting conventions. First, a 6 month USD Libor rate that fixes at time T , does not truly cover an *accrual period* of $[T, T + 1/2]$. Instead, the start date T^s of the accrual is set to be $T^s = T + \delta^s$, where δ^s is a delay of two business days²⁴. In other words, the quoted spot Libor rate is in actuality based on a *forward starting* CD that is entered into with time lag of δ^s after the quotation date T . As for the end date T^e of the accrual period, it is normally determined by counting 6 months ahead starting from T^s , adjusting the resulting date to ensure that it is a valid business day. The precise mechanism used to make such a business day adjustments of T^s is determined from a *date rolling convention*. For USD Libor, one always uses the so-called “Modified-Following” convention where weekend or holiday dates are rolled forward to the next business day, unless doing so would cause T^e to lie in the next calendar month, in which case the payment date is rolled to the previous business day. Other rolling conventions are discussed in Mayle [1993] and Stigum and Robinson [1996].

Once the correct accrual period $[T^s, T^e]$ has been determined, to compute rate accrual it remains to compute the proper *year fraction* (or *accrual factor* or *day count fraction*) τ representing how many whole years are spanned by $[T^s, T^e]$. For the purposes of our book, we normally write simply $\tau = T^e - T^s$, but a little thought shows that expressions like this are ambiguous when T^e and T^s are thought of as actual (discrete) calendar dates, rather than as arbitrary numbers on the real line. For instance, given the existence of leap years, how many days are there in a standard year? For quant purposes, it is common to assume that a calendar year has 365.25 years, such that $T^e - T^s$ is obtained by simply counting the number of days between T^s and T^e and then dividing this number by 365.25. This “convention” is sometimes known as Actual/365.25 (or sometimes just A365.25), and is rarely, if ever, used for actual market quoting purposes. Instead, for quotation of Libor

²⁴Libor rates in other currencies may have different delays. For instance, GBP Libor has zero business day delay.

rates the standard is to use an Actual/360 (A360) convention, where the number of days between T^s and T^e are converted to a year fraction by dividing by 360. As a consequence, the true value of τ used for 6 month Libor quotation purposes is typically slightly larger than 1/2. For additional year-count conventions (of which there are many), see Mayle [1993] and Stigum and Robinson [1996].

Due to the quoting standards used in real Libor markets, the relationship between discount bonds and quoted Libor rates is more complicated than (5.27). Specifically, if $L_{\text{mkt}}(T)$ represents the true quoted 6 month Libor rate at time T , we instead have

$$L_{\text{mkt}}(T) = \left((P(T, T^e)/P(T, T^s))^{-1} - 1 \right) \frac{360}{D(T^s, T^e)}, \quad (5.28)$$

where by $D(T^s, T^e)$ we denoted the number of days between T^s and T^e according to the convention used. Notice in particular how the formula now involves a forward starting zero-coupon bond $P(T, T^s, T^e) = P(T, T^e)/P(T, T^s)$, as a reflection of the settlement delay δ^s . Using existing “idealized” Libor rate notation (see (4.2)), we may write this expression as

$$L_{\text{mkt}}(T) = L(T, T^s, T^e) \times \frac{360}{365.25} \approx L(T, T^s, T^e) \times 0.986.$$

The difference between $L_{\text{mkt}}(T)$ and $L(T)$ is small enough for us to ignore it in most of this book, but any real system implementation obviously should use precise day counting rules when computing Libor fixings.

5.A.2 Swap Payments

The payments on swaps (and other instruments, such as CDs and FRAs) are subject to similar conventions as the Libor rate. Consider for instance a standard fixed-for-floating interest rate swap issued at time t (today). First, a²⁵ schedule $\{T_i\}_{i=0}^N$ for interest rate accrual must be constructed, starting from a given base frequency of the swap (e.g.: semi-annual). As was the case above, the schedule normally starts one or two business days after time t , i.e. $T_0 = t + \delta_0$ where δ_0 is some contractually specified delay. Date T_0 is known as the *effective date* of the swap. Given T_0 , the remaining T_i , $i = 1, \dots, N$, are computed by first laying out “unadjusted” dates according to the swap base frequency, and then applying a date rolling convention (typically Modified-Following) to each of the dates. As part of the swap contract, associated with each accrual period $[T_i, T_{i+1}]$ are then:

- A fixing date T_i^f : the date on which the floating leg index (Libor, most often) is observed. Typically T_i^f is two business days before time T_i .

²⁵We assume that the fixed and floating legs pay interest on the same schedule, but in reality, this may not be the case. For instance, in USD, the standard frequency for the fixed leg is six months, and three months for the floating leg.

- A payment date T_i^p : the date on which the swap payments are made. Normally $T_i^p = T_{i+1}$, but it is not uncommon to have payment delays of 1 or 2 business days after T_{i+1} .
- A fixed leg year fraction τ_i^{fix} : the year fraction used to determine the payment at time T_i^p on the fixed leg. In the US, the most common convention for the fixed leg is²⁶ 30/360.
- A floating leg year fraction τ_i^{flt} : the year fraction used to determine the payment at time T_i^p on the floating leg²⁷. In the US, the most common convention for the floating leg is Actual/360.

At time t , the value of a payer swap paying a fixed coupon c against Libor is therefore

$$V_{\text{swap}}(t) = \sum_{i=1}^N P(t, T_i^p) E^{T_i^p} \left(\tau_i^{flt} L_{\text{mkt}}(T_i^f) - \tau_i^{fix} c \right),$$

where we have used the T_i^p -forward measure to state the valuation, and where L_{mkt} is defined in (5.28). In the book, we normally simplify this to

$$V_{\text{swap}}(t) = \sum_{i=1}^N P(t, T_i) \tau_i E^{T_i} (L(T_{i-1}) - c), \quad \tau_i = T_i - T_{i-1}.$$

²⁶When counting days in the 30/360 convention, each month is assumed to have 30 days. The number of days used to determine interest rate accrual will therefore differ from the *actual* number of days (which distinguishes 30/360 from Actual/360).

²⁷Another, often ignored, complication with the floating leg is that the periods that define the (payment) year fractions τ_i^{flt} are sometimes slightly different from those that define forward Libor rates, due to certain conventions surrounding date adjustments. Again, see details in Mayle [1993] or Stigum and Robinson [1996].

Yield Curve Construction and Risk Management

In a nutshell, the job of an interest rate model is to describe the random movement of a curve of discount bond prices through time, starting from a known initial condition. In reality, however, only a few short-dated discount bonds are directly quoted in the market at any given time, a long stretch from the assumption of many models that an initial curve of discount bond prices is observable for a continuum of maturities out to 20–30 years or more. Fortunately, a number of liquid securities depend in relatively straightforward fashion on discount bonds, opening up the possibility of uncovering discount bond prices from prices of such securities. Still, as only a finite set of securities are quoted in the market, constructing a continuous curve of discount bond prices will inevitably require us to complement market observations with an interpolation rule, based perhaps on direct assumptions about functional form or perhaps on a regularity norm to be optimized on. A somewhat specialized area of research, discount curve construction relies on techniques from a number of fields, including statistics and computer graphics. While we cannot possibly do the subject full justice, discount curve construction is a fundamental step in the modeling exercise, and no book on fixed income models is complete without a discussion of basic techniques.

As mentioned in the Preface to this book, the crisis of 2007–2009 have lead to changes in the foundations of interest rate modeling, not least in the area of yield curve construction and risk management. Pre-crisis, it was often sufficient to construct only a single (Libor) discount curve, but nowadays the task is more complicated as a whole collection of inter-related curves is required. Nevertheless, the traditional techniques used for single-curve construction are by no means obsolete, and their mastery is required before more ambitious curve algorithms can be attempted. Accordingly, we have split this chapter into three parts. In the first, and most significant, part, we introduce notations and cover a number of curve construction techniques, moving from simply bootstrapped C^0 curves through “local spline” C^1 curves to full C^2 smoothing splines with and without tension.

Perturbation locality is discussed, as are methods to control behavior under perturbations. In the second part we discuss the management of interest rate curve risk, covering both basic approaches as well as more advanced methods based on Jacobian techniques. In the last part, we discuss a number of specialized issues and contemporaneous extensions, most notably turn-of-year adjustments and techniques to construct separate discount and forward curves. The need for such a separation has long been recognized (albeit neglected in the literature) as a requirement to avoid arbitrages in markets for foreign exchange forwards and for floating-floating cross-currency swaps. More recently, similar issues have appeared in purely domestic markets where the Libor rate is no longer considered a good proxy for the risk-free discount rate, and where a significant *tenor basis* has developed in floating-floating single-currency swaps. Accordingly, we conclude the chapter with a description of techniques for building a *multi-index curve group*, a self-consistent arbitrage-free collection of discount and forward curves suitable for valuation of different types of swaps and other interest rate derivatives.

6.1 Notations and Problem Definition

6.1.1 Discount Curves

Throughout this chapter, we use the abbreviated notation $P(T) = P(0, T)$ where $P : [0, \mathcal{T}] \rightarrow (0, 1]$ is a continuous, monotonically decreasing *discount curve*. \mathcal{T} denotes the maximum maturity considered, typically given as the longest maturity in the set of securities the curve is built from. Let there be N such securities — the *benchmark set* — with observable prices V_1, \dots, V_N . We assume that the time 0 price $V_i = V_i(0)$ of security i can be written as a linear combination of discount bond prices at different maturities,

$$V_i = \sum_{j=1}^M c_{i,j} P(t_j), \quad i = 1, \dots, N, \quad (6.1)$$

where $0 < t_1 < t_2 < \dots < t_M \leq \mathcal{T}$ is a given finite set of dates, in practice obtained by merging together the cash flow dates of each of the N benchmark securities. Let T_1, T_2, \dots, T_N denote the final maturities of the N benchmark securities, in which case we necessarily must have

$$c_{i,j} = 0, \quad t_j > T_i.$$

Securities that can be represented by pricing expressions of the form (6.1) obviously include coupon and discount bonds, but also FRAs and fixed-floating interest rate swaps. For instance, consider a newly issued unit-notional fixed-floating swap, paying a coupon of $c\tau$ at times $\tau, 2\tau, 3\tau, \dots, n\tau$. If no spread is paid on the floating rate, the time 0 total swap value to the fixed payer is

$$V_{\text{swap}} = 1 - P(t_n) - \sum_{j=1}^n c\tau P(j\tau),$$

as already discussed in Chapter 5. We can rewrite this as

$$1 - V_{\text{swap}} = P(t_n) + \sum_{j=1}^n c\tau P(j\tau), \quad (6.2)$$

which is in the form¹ (6.1) once we interpret $V_i = 1 - V_{\text{swap}}$. In practice, swaps used for discount curve construction are nearly always newly issued and par-valued, in the sense that the coupon c is set to make $V_{\text{swap}} = 0$. To give another example, consider an FRA on the $[T, T + \tau]$ Libor rate, for which formula (5.2) in Chapter 5 gives, at $t = 0$,

$$V_{\text{FRA}} = \tau (L(0, T, T + \tau) - k) P(T + \tau), \quad (6.3)$$

where k is the quoted FRA rate. From the definition of $L(0, T, T + \tau)$ this is just

$$V_{\text{FRA}} = P(T) - P(T + \tau) - k\tau P(T + \tau) = P(T) - (1 + k\tau), P(T + \tau)$$

which is in the form (6.1). As for swaps, FRAs used for curve construction are newly issued and typically have k set such that $V_{\text{FRA}} = 0$.

The choice of the securities to be included in the benchmark set depends on the market under consideration. For instance, to construct a Treasury bond curve, it is natural to choose a set of Treasury bonds and T-Bills. On the other hand, if we are interested in constructing a discount curve applicable for bonds issued by a particular firm, we would naturally use bonds and loans issued by the firm in question. For our purposes, the most important discount curve is the *Libor curve*, constructed out of market quotes for Libor deposits, swaps and Eurodollar futures. In the construction of this curve, most firms would use a few certificates of deposit for the first 3 months of the curve, followed by a strip of Eurodollar futures² (with maturities staggered 3 months apart) out to 3 or 4 years. Par swaps are then used for the rest of the curve, with typical maturities being 5, 7, 10, 12, 15, 20, 25, and 30 years.

¹For swaps where payment schedules do not coincide perfectly with the accrual periods of the Libor rates, the expression (6.2) is only an approximation, albeit a very good one. In practice we can construct the yield curve assuming that (6.2) holds, and then perform a small post-processing clean-up iteration, along the lines of the algorithm in Section 6.5.2.4.

²We note that Eurodollar futures contracts do not allow for a pricing expression of the form (6.1), so a pre-processing step is normally employed to convert the futures rate quote to a forward rate (FRA) quote. See Proposition 4.5.3 or Chapter 16 for more on this.

6.1.2 Matrix Formulation

Define the M -dimensional discount bond vector³

$$\mathbf{P} = (P(t_1), \dots, P(t_M))^T,$$

and let $\mathbf{V} = (V_1, \dots, V_N)^T$ be the vector of observable security prices. Also let $\mathbf{c} = \{c_{i,j}\}$ be an $(N \times M)$ -dimensional matrix containing all the cash flows produced by the chosen set of securities. The matrix \mathbf{c} would typically be quite sparse, with many rows containing only a few non-zero entries. A typical, albeit simplified, form of the matrix \mathbf{c} might be (\times marks a non-zero element)

$$\mathbf{c} = \begin{pmatrix} \times & & & & & & & \\ \times & & & & & & & \\ & \times \times & & & & & & \\ & & \times \times & & & & & \\ & & & \times \times & & & & \\ & & & & \times \times & & & \\ & & & & & \times \times & & \\ \times \times \times \times \times \times \times \times & & & & & & & \\ \times \times \times \times \times \times \times \times \times & & & & & & & \\ \times & & & & & & & \end{pmatrix},$$

corresponding to two certificates of deposit (first two rows); four FRAs or Eurodollar futures (next four rows); and three swaps (last three rows).

In a consistent, friction-free market without arbitrage opportunities, the fundamental relation

$$\mathbf{V} = \mathbf{c}\mathbf{P} \quad (6.4)$$

must be satisfied, giving us a starting point to determine \mathbf{P} .

6.1.3 Construction Principles and Yield Curves

In practice, we normally have more cash flow dates than benchmark security prices, i.e. $M > N$, in which case (6.4) is insufficient to uniquely determine \mathbf{P} . The problem of curve construction essentially boils down to supplementing (6.4) with enough additional assumptions to allow us to extract \mathbf{P} and to determine $P(T)$ for values of T not in the cash flow timing set $\{t_j\}_{j=1}^M$.

As it is normally easier to devise an interpolation scheme on a curve that is reasonably flat (rather than exponentially decaying), it is common to perform the curve fitting exercise on *discount yields*, rather than directly on bond prices⁴. Specifically, we introduce a continuous yield function $y : [0, T] \rightarrow \mathbb{R}_+$ given by

³For extra clarity, throughout this chapter we use boldface type for vectors and matrices.

⁴See e.g. Shea [1984] for a discussion of the pitfalls associated with curve interpolators that work directly on the discount function $P(T)$.

$$e^{-y(T)T} = P(T) \quad \Rightarrow \quad y(T) = -T^{-1} \ln P(T), \quad (6.5)$$

such that in (6.4)

$$\mathbf{P} = \left(e^{-y(t_1)t_1}, \dots, e^{-y(t_M)t_M} \right)^\top.$$

The mapping $T \mapsto y(T)$ is known as the *yield curve*; it is related to the discount curve by the simple transformation (6.5). Of related interest is also the *instantaneous forward curve* $f(T)$, given by

$$P(T) = e^{-\int_0^T f(u)du}. \quad (6.6)$$

Notice that

$$f(T) = y(T) + \frac{dy(T)}{dT} T. \quad (6.7)$$

For alternative transformations, and a discussion of their relative merits, see Andersen [2005]. Unless explicitly stated, in the remainder of this chapter we shall work with yields, i.e. we treat $y(T)$ as the fundamental curve to be estimated.

Whatever space we elect to work in, we have at least three options for solving (6.4).

1. We can introduce new and unspanned securities such that $N = M$ and (6.4) allows for exactly one solution.
2. We can use a parameterization of the yield curve with precisely N parameters, using the N equations in (6.4) to recover these parameters.
3. We can search the space of all solutions to (6.4) and choose the one that is “optimal” according to a given criterion.

Let us provide some comments to these three ideas. First, in option 1 above, introduction of new securities might not truly be possible — such securities may simply not exist — but sometimes interpolation rules applied to the given benchmark set may allow us to provide reasonable values for an additional set of “fictitious” securities. Although it can occasionally be useful in pre-processing to pad an overly sparse benchmark set, this idea will often require some quite ad-hoc decisions about the specifics of the fictitious securities, and excessive use may ultimately lead to odd-looking curves and suboptimal hedge reports. When an interpolation rule is to be used, it is typically better to apply it directly on more fundamental quantities such as zero-coupon yields or forward rates, thereby maintaining a higher degree of control over the resulting discount curve.

In option 2 above, parametric functional forms (e.g. Nelson and Siegel [1987]) are sometimes used, but it is far more common to work with a spline representation with N user-selected knots (typically at the maturity dates of the benchmark securities), with the level of the yield curve at these knots constituting the N unknowns to be solved for. We discuss the details of

this approach in Section 6.2, using a number of different spline types. Some required elements of basic spline theory can be found in Appendix 6.A of this chapter.

Option 3 constitutes the most sophisticated approach and can often be stated in completely non-parametric terms, with the yield curve emerging naturally as the solution to an optimization problem. If carefully stated, this approach can easily be modified to handle the situation where the system of equations (6.4) is (near-) singular, in the sense that either no solutions exist or all solutions are irregular and non-smooth⁵. Technically, we handle this by working with *smoothing splines*, in the process replacing (6.4) with a penalized least-squares optimization problem. We discuss elements of this idea in Section 6.3 below.

6.2 Yield Curve Fitting with N -Knot Splines

In this section we discuss a number of well-known yield curve algorithms based on polynomial and exponential (tension) splines of various degrees of differentiability. Throughout, we assume that we can select and arrange our benchmark set of securities to guarantee that the maturities of the benchmark securities satisfy

$$T_i > T_{i-1}, \quad i = 2, 3, \dots, N, \quad (6.8)$$

where the inequality is strict. Equation (6.8) constitutes a “spanning” condition and allows us to select the N maturities as distinct knots in our yield curve splines.

6.2.1 C^0 Yield Curves: Bootstrapping

If continuity of the yield curve is all that we require, we can work with a common iterative procedure known as *bootstrapping*. The basic idea is encapsulated in the following iteration:

1. Let $P(t_j)$ be known for $t_j \leq T_{i-1}$, such that prices for benchmark securities $1, \dots, i-1$ are matched.
2. Make a guess for $P(T_i)$.
3. Use an interpolation rule to fill in $P(t_j)$, $T_{i-1} < t_j < T_i$.
4. Compute V_i from the now-known values of $P(t_j)$, $t_j \leq T_i$.
5. If V_i equals the value observed in the market, stop. Otherwise return to Step 2.
6. If $i < N$, set $i = i + 1$ and repeat.

⁵Intuitively, this situation can arise if, say, two or more securities in the benchmark set have near-identical cash flows, yet have significantly different present values.

The updating of guesses when iterating over Steps 2 through 5 can be handled by any standard one-dimensional root-search algorithm (e.g., the Newton-Raphson or secant methods).

There are strong limitations on what kind of interpolation rule can be applied in Step 3. For instance, one might consider using a representation in terms of instantaneous forwards $f(T)$ (see (6.6)), with the assumption that $f(T)$ is a continuous piecewise linear function on the maturity grid $\{T_i\}_{i=1}^N$. While based on seemingly natural assumptions, this interpolation rule can, however, be shown to be numerically unstable and prone to oscillations. Some stable, and standard, choices for interpolation rules are covered in the next two sections; common for both is that the resulting yield curve is continuous, but non-differentiable. This, in turn, implies that the instantaneous forward curve is discontinuous (see (6.7)).

6.2.1.1 Piecewise Linear Yields

The most common discount curve bootstrap algorithm assumes that the continuously compounded yield $y(T)$ in (6.5) is a continuous piecewise linear function on $\{T_i\}_{i=1}^N$. Formally, the interpolation rule in Step 3 of the algorithm in Section 6.2.1 writes $P(T) = e^{-y(T)T}$, where

$$y(T) = y(T_i) \frac{T_{i+1} - T}{T_{i+1} - T_i} + y(T_{i+1}) \frac{T - T_i}{T_{i+1} - T_i}, \quad T \in [T_i, T_{i+1}]. \quad (6.9)$$

To initiate the iterative bootstrap algorithm, we note that the interpolation rule (6.9) may require us to provide an equation for $y(t)$, $t < T_1$. There are a number of ways to do this; one common choice is to simply set $y(t) = y(T_1)$, $t < T_1$.

To give a feel for the types of yield curves produced by linear yield bootstrapping, let us consider a simple example with a benchmark set of $N = 10$ swaps, with maturities and quoted par swap rates as given in Table 6.1⁶.

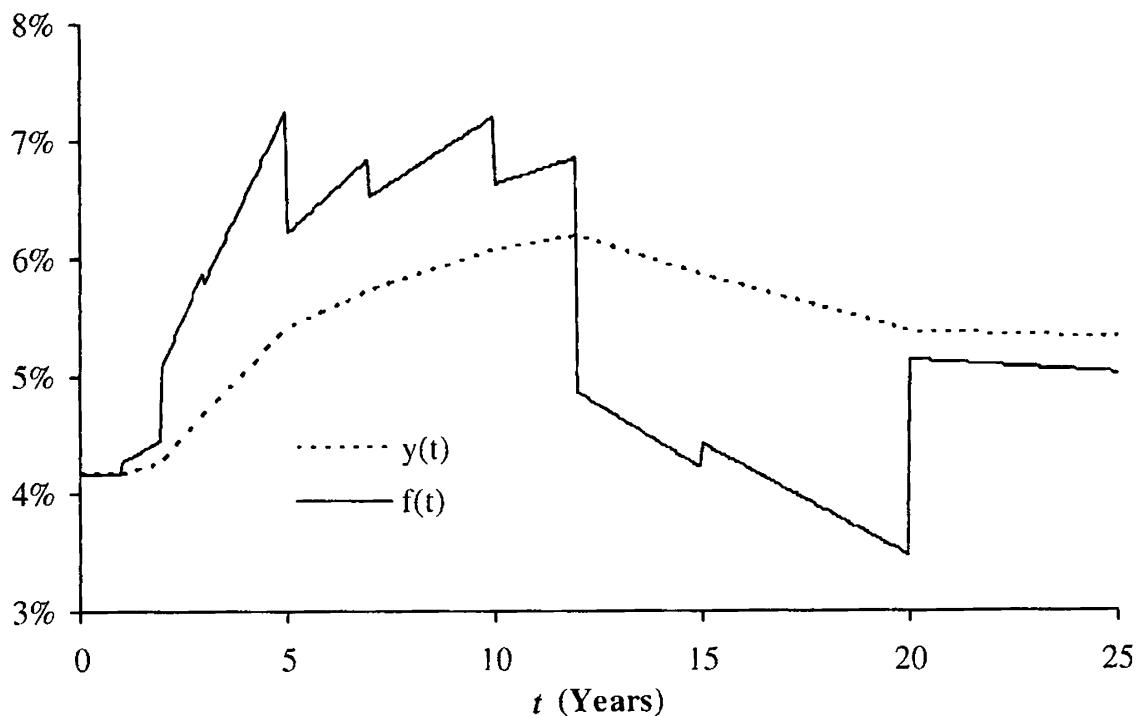
The swaps are assumed to pay on a semi-annual basis,

$$t_j = j \cdot 0.5, \quad j = 1, 2, \dots, 50.$$

Setting $y(t) = y(1)$, $t < 1$, and then running the bootstrap procedure on the swap price expression (6.2) results in the yield shown in Figure 6.1. The same figure also shows the continuously compounded forward curve, as computed by equation (6.7). The discontinuous “saw-tooth” shape of the forward curve is characteristic for bootstrapped yield curves with piecewise linear yield.

⁶In actual markets, swap yields are most often increasing functions of the swap maturity, rather than humped as in Table 6.1. The data in Table 6.1 was picked to stress the curve construction algorithms, in order to emphasize their strengths and weaknesses.

Maturity (Years)	Swap Par Rate
1	4.20%
2	4.30%
3	4.70%
5	5.40%
7	5.70%
10	6.00%
12	6.10%
15	5.90%
20	5.60%
25	5.55%

Table 6.1. Swap Benchmark Set for Numerical Tests**Fig. 6.1.** Yield and Forward Curve

Notes: Yield curve is constructed by bootstrapping, assuming piecewise linear yields. Swap data is in Table 6.1.

6.2.1.2 Piecewise Flat Forward Rates

Assume now that the instantaneous forward curve is piecewise flat, switching to a new level at each point in $\{T_i\}$, i.e.

$$f(T) = f(T_i), \quad T \in [T_i, T_{i+1}), \quad (6.10)$$

with $T_0 \triangleq 0$. This corresponds to an interpolation rule where $\ln P(T)$ is linear in T , or

$$P(T) = P(T_i) e^{-f(T_i)(T-T_i)}, \quad T \in [T_i, T_{i+1}),$$

where a bootstrap algorithm can be used to establish the values of the N unknown constants $f(T_0), f(T_1), \dots, f(T_{N-1})$. From the equation

$$y(T)T = \int_0^T f(u) du,$$

we see that the assumption of piecewise flat forwards gives, for $T \in [T_i, T_{i+1})$,

$$y(T) = \frac{y(T_i)T_i + f(T_i)(T - T_i)}{T} = f(T_i) + \frac{(y(T_i) - f(T_i))T_i}{T},$$

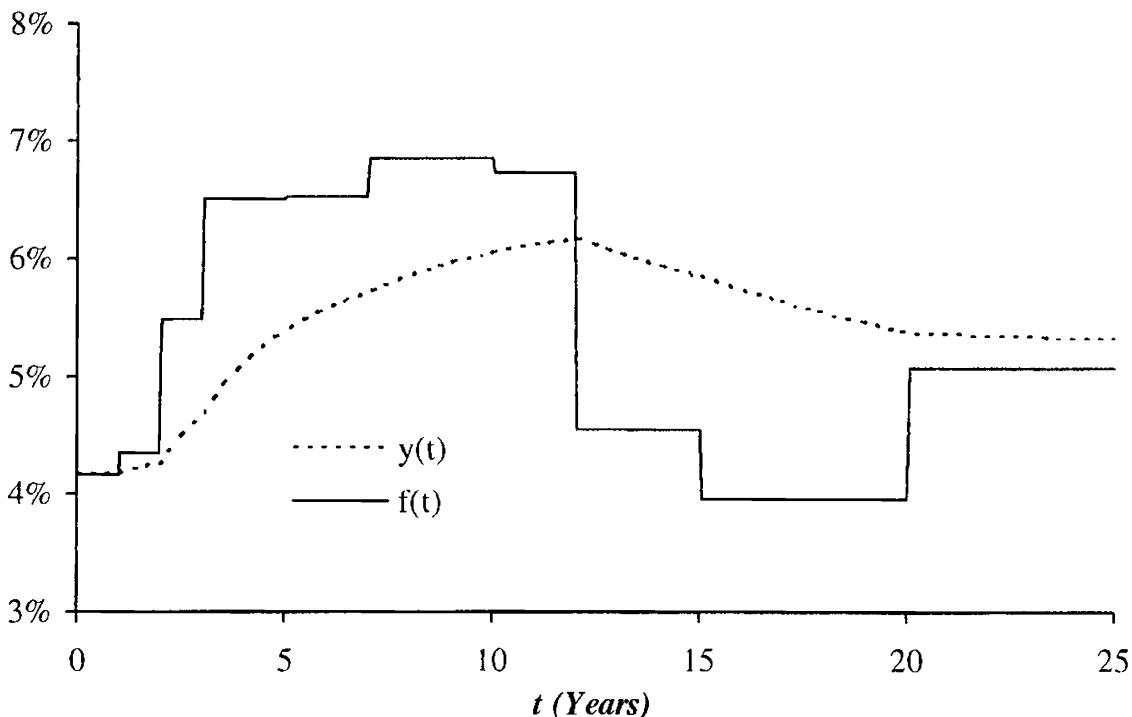
or

$$y(T) = \frac{1}{T} \left(T_i y(T_i) \frac{T_{i+1} - T}{T_{i+1} - T_i} + T_{i+1} y(T_{i+1}) \frac{T - T_i}{T_{i+1} - T_i} \right).$$

The yield curve will remain continuous.

Figure 6.2 below shows the results of applying (6.10) to the swap data in Table 6.1. Notice the non-linear behavior of yields between maturity dates and the staircase shape of the forward curve.

Fig. 6.2. Yield and Forward Curve



Notes: Yield curve is constructed by bootstrapping, assuming piecewise flat forward rates. Swap data is in Table 6.1.

6.2.2 C^1 Yield Curves: Hermite Splines

As we have seen, simply bootstrapped curves generally result in a discontinuous forward curve. From an empirical/economic perspective, such discontinuities are often unrealistic, and may also result in distortions of derivative prices⁷ and technical difficulties in dynamic yield curve models. In this section, we consider a simple scheme to extend the bootstrapping technique to produce a once-differentiable yield curve and a continuous forward curve. Our scheme relies on *Hermite cubic splines*, where we write

$$y(T) = a_{3,i}(T - T_i)^3 + a_{2,i}(T - T_i)^2 + a_{1,i}(T - T_i) + a_{0,i}, \quad T \in [T_i, T_{i+1}], \quad (6.11)$$

for a series of constants $a_{3,i}$, $a_{2,i}$, $a_{1,i}$, $a_{0,i}$ to be determined from given values of $y(T_i)$, $y(T_{i+1})$, $y'(T_i)$, and $y'(T_{i+1})$. Appendix 6.A.1 contains a review of Hermite spline theory.

A particularly popular choice among Hermite splines is the *Catmull-Rom spline*, where derivatives $y'(T_i)$, $i = 1, \dots, N$, are constructed by finite differences, relieving the user from directly specifying them. As shown in Appendix 6.A.1, for the Catmull-Rom spline we can organize (6.11) in a vector/matrix form as

$$y(T) = \mathbf{D}_i(T)^\top \mathbf{A}_i \begin{pmatrix} y_{i-1} \\ y_i \\ y_{i+1} \\ y_{i+2} \end{pmatrix}, \quad T \in [T_i, T_{i+1}], \quad i = 1, \dots, N-1, \quad (6.12)$$

where, adapting as necessary the notation from the previous section,

$$\mathbf{D}_i(T) = \begin{pmatrix} d_i^3 \\ d_i^2 \\ d_i \\ 1 \end{pmatrix}, \quad d_i = \frac{T - T_i}{h_i}, \quad y_i = y(T_i), \quad h_i = T_{i+1} - T_i,$$

and the matrix \mathbf{A}_i is as given in (6.54)–(6.56) in Appendix 6.A.1. While nominally (6.12) involves the values y_{N+1} and y_0 , the matrices \mathbf{A}_{N-1} and \mathbf{A}_1 are such that these values are irrelevant.

The Catmull-Rom spline prescription (6.12) completely specifies the yield curve on the interval $[T_1, T_N]$, given the N constants y_1, \dots, y_N . To extend the yield curve to cover the interval $[0, T_1]$, we need to supply additional extrapolation assumptions. As in bootstrapping, possible choices for this additional equation is $y_0 = y(0) = y_1$, or perhaps the slope condition

$$\frac{y_1 - y_0}{h_0} = \frac{y_2 - y_1}{h_1}. \quad (6.13)$$

⁷For instance, as deal maturity crosses a point of discontinuity on the forward curve, the price of an FRA or a caplet on a short-tenor rate will jump.

Away from the boundaries, we notice that the price of security i depends only on y_1, \dots, y_{i+1} , as the pricing equations take the diagonal form

$$\begin{aligned} V_1 &= F_1(y_1, y_2, y_3), \\ V_2 &= F_2(y_1, y_2, y_3), \\ V_3 &= F_3(y_1, y_2, y_3, y_4) \\ &\vdots \\ V_{N-1} &= F_{N-1}(y_1, \dots, y_N), \\ V_N &= F_N(y_1, \dots, y_N), \end{aligned}$$

for non-linear functions F_i . Here F_i is typically only mildly sensitive to y_{i+1} , so the system of equations is nearly, but not quite, in bootstrap form. This makes solving for the y_i 's an easy fare for a standard non-linear root-search algorithm (see Press et al. [1992] for several algorithms). We can also consider an iteration on a series of bootstrap procedures. To describe this idea, let $y_i^{(k)}$ be the value for y_i found in the k -th iteration, and consider then the following algorithm:

1. Let $y_j^{(k)}$, $j = 1, \dots, i-1$, and $y_{i+1}^{(k-1)}$ all be known.
2. Make a guess for $y_i^{(k)}$.
3. Compute $V_i = F_i(y_1^{(k)}, \dots, y_i^{(k)}, y_{i+1}^{(k-1)})$.
4. If V_i equals the market value stop. Otherwise return to Step 2.
5. If $i < N$, set $i = i + 1$ and repeat.

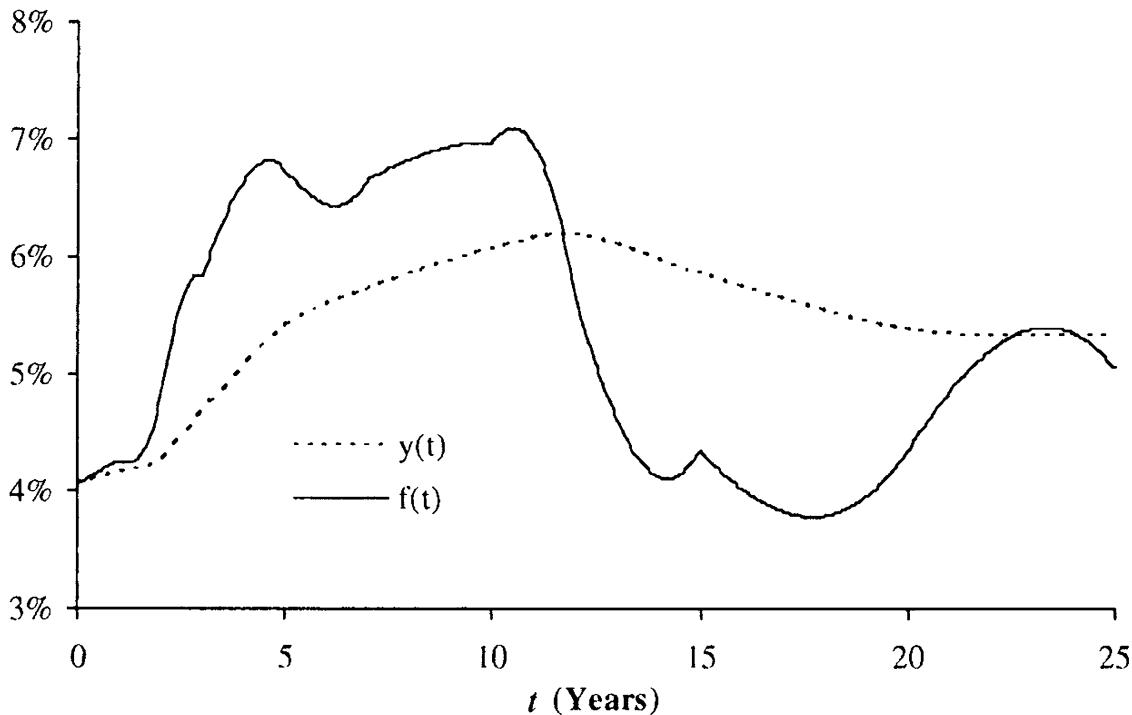
We emphasize that the iteration over Steps 2–4 is still only one-dimensional, as in the bootstrapping algorithm of Section 6.2.1. Upon completion, the algorithm above yields $y_1^{(k)}, \dots, y_N^{(k)}$. Iterating over k , we repeat the algorithm until the differences between the yields found at the k -th and $(k+1)$ -th iteration are sufficiently small, say when

$$N^{-1} \sum_{i=1}^N \left(y_i^{(k+1)} - y_i^{(k)} \right)^2 < \varepsilon^2,$$

where ε is a given tolerance. To initialize the iteration over k , we need a starting guess $y_1^{(0)}, \dots, y_N^{(0)}$; a good choice is the yield curve constructed by regular bootstrapping.

In Figure 6.3, we show the results of applying the algorithm above (using the boundary choice (6.13)) to the numerical example of Sections 6.2.1.1 and 6.2.1.2. We see that, as desired, the yield curve is smooth and the instantaneous forward curve is continuous. As the yield curve by construction is only once differentiable, equation (6.7) shows that the forward curve is not differentiable; this is obvious from the figure.

We can easily extend the procedure above beyond Catmull-Rom splines to more complicated C^1 cubic splines in the Hermite class, using results

Fig. 6.3. Yield and Forward Curve

Notes: Yield curve is assumed to be a Catmull-Rom cubic spline. Swap data is in Table 6.1.

from Appendix 6.A. For instance, it is relatively straightforward to add *tension* to the Catmull-Rom spline. We cover twice-differentiable tension splines later in this chapter.

6.2.3 C^2 Yield Curves: Twice Differentiable Cubic Splines

While the spline method introduced in the previous section often produces acceptable yield curves, the method is heuristic in nature and ultimately does not produce a smooth forward curve. To improve on the latter, one alternative is to remain in the realm of cubic splines, but now insist that the curve is twice differentiable everywhere on $[T_1, T_N]$. We then write (see Appendix 6.A.2)

$$\begin{aligned} y(T) = & \frac{(T_{i+1} - T)^3}{6h_i} y''_i + \frac{(T - T_i)^3}{6h_i} y''_{i+1} + (T_{i+1} - T) \left(\frac{y_i}{h_i} - \frac{h_i}{6} y''_i \right) \\ & + (T - T_i) \left(\frac{y_{i+1}}{h_i} - \frac{h_i}{6} y''_{i+1} \right), \quad T \in [T_i, T_{i+1}], \end{aligned} \quad (6.14)$$

where $y''_i = d^2y(T_i)/dT^2$, $y_i = y(T_i)$, and $h_i = T_{i+1} - T_i$. The appendix demonstrates that continuity of the second derivative across the $\{T_i\}$ knots requires that the y''_i and y_i are connected through a tri-diagonal linear system of equations, see equation (6.62). To state the expressions explicitly in matrix format, let $\mathbf{y}'' = (y''_2, y''_3, \dots, y''_{N-2}, y''_{N-1})^\top$ and

$\mathbf{y} = (y_2, y_3, \dots, y_{N-2}, y_{N-1})^\top$ such that

$$\mathbf{B}\mathbf{y}'' = \mathbf{C}\mathbf{y} + \mathbf{M}(y_1, y_N, y_1'', y_N''), \quad (6.15)$$

where the matrices \mathbf{B} and \mathbf{C} are both $(N - 2) \times (N - 2)$ tri-diagonal, with elements given by

$$B_{i,i} = \frac{h_i + h_{i+1}}{3}, \quad B_{i,i+1} = \frac{h_{i+1}}{6}, \quad B_{i,i-1} = \frac{h_i}{6},$$

and

$$C_{i,i} = -\left(\frac{1}{h_i} + \frac{1}{h_{i+1}}\right), \quad C_{i,i+1} = \frac{1}{h_{i+1}}, \quad C_{i,i-1} = \frac{1}{h_i}.$$

The $(N - 2)$ -dimensional vector $\mathbf{M}(y_1, y_N, y_1'', y_N'')$ captures boundary terms at T_1 and T_N . The most important — and, as discussed later, in a sense *best* — boundary specification is that of the *natural spline*, where we set $y_1'' = y_N'' = 0$. In this case, we have

$$\mathbf{M}(y_1, y_N, y_1'', y_N'') = \mathbf{M}(y_1, y_N) = \left(\frac{y_1}{h_1}, 0, 0, \dots, 0, 0, \frac{y_N}{h_{N-1}}\right)^\top.$$

Notice that application of a natural boundary condition at time T_1 allows us to recover yields inside the time period $[0, T_1]$ by linear interpolation, using the gradient $y'(T_1)$ at time T_1 (which can easily be found by differentiating (6.14)).

We notice that (6.14) combined with (6.15) allows us to turn any guess of y_1, y_2, \dots, y_N into a guess for the vector \mathbf{P} in (6.4). Specifically, we perform the following steps:

1. Compute the right-hand side of (6.15).
2. Use a standard tri-diagonal LU solver (see Press et al. [1992]) to invert (6.15) and recover \mathbf{y}'' .
3. Apply (6.14) to determine⁸ all values of $y(t_j)$, $j = 1, \dots, M$, extrapolating as necessary when $t_j < T_1$.
4. Use (6.5) to establish \mathbf{P} .

The computational effort of Steps 1 through 4 are $O(N)$, $O(N - 2)$, $O(M)$, and $O(M)$, respectively.

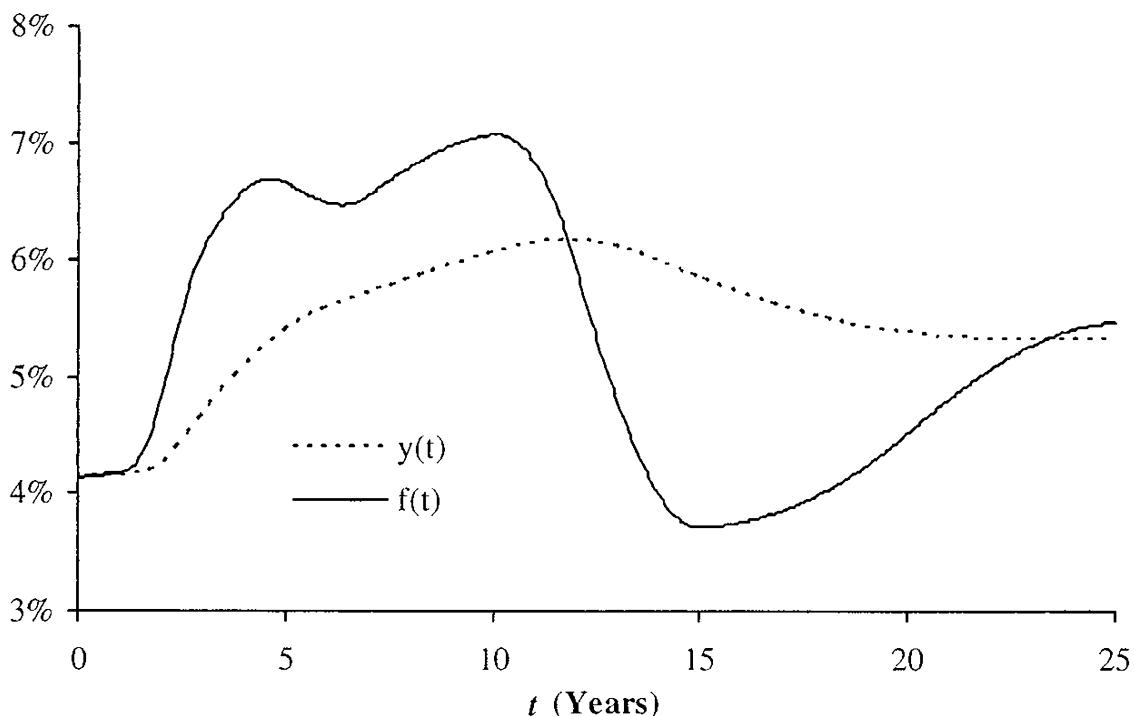
To solve for the correct values of y_1, y_2, \dots, y_N , we iterate on Steps 1–4 using a non-linear root-search algorithm, terminating when (6.4) is satisfied to within acceptable tolerances. The fitting problem is typically good-natured, and virtually all standard root-search packages (see Press et al. [1992]) can tackle it successfully. Tanggaard [1997], for instance, uses

⁸For computational reasons, the terms multiplying the various y and y'' in (6.14) should be pre-cached, to avoid wasting effort when we ultimately perform an iteration.

a simple Gauss-Newton scheme with good results. Whatever root-search algorithm is selected, a good first guess can always be found by simple bootstrapping.

In Figure 6.4, we show the results of applying the algorithm above to a natural cubic spline representation of the yield curve example used in earlier sections. The yield curve is smooth and, unlike the Hermite spline case in Figure 6.3, the instantaneous forward curve is now differentiable, as desired.

Fig. 6.4. Yield and Forward Curve



Notes: Yield curve is assumed to be a C^2 natural cubic spline. Swap data is in Table 6.1.

While the C^2 cubic spline discussed here has attractive smoothness, it is not necessarily an ideal representation of the yield curve. As discussed in Andersen [2005] and Hagan and West [2004], among others, twice differentiable cubic spline yield curves are often subject to oscillatory behavior, spurious inflection points, poor extrapolatory behavior, and non-local behavior when prices in the benchmark set are perturbed. We shall return to the concept of non-local perturbation effects in Section 6.4 below, but for now just note that perturbation of a single benchmark price can cause a slow-decaying “ringing” effect on the C^2 cubic yield curve, with the effect of the perturbation of the benchmark instrument price spilling into the entire yield curve. This behavior is not surprising, given that the spline is constructed through a full $(N - 2) \times (N - 2)$ matrix system, where interpolation behavior on the interval $[T_i, T_{i+1}]$ depends on *all* values y_j , $j = 1, \dots, N$. In contrast, the simple linear-yield bootstrapping method in Section 6.2.1 interpolation on

the interval $[T_i, T_{i+1}]$ involves only the two points y_i and y_{i+1} , and the Hermite spline approach involves only the four points $y_{i-1}, y_i, y_{i+1}, y_{i+2}$.

6.2.4 C^2 Yield Curves: Twice Differentiable Tension Splines

C^1 Hermite cubic splines are less prone to non-local behavior than C^2 cubic splines, but accomplish this in a somewhat ad-hoc fashion by giving up one degree of differentiability. Rather than taking such a draconian step, one wonders whether there may be a way to retain the C^2 feature of the cubic spline in Section 6.2.3, yet still allow control of curve locality and “stiffness”. As it turns out, an attractive remedy to the shortcomings of the pure C^2 cubic spline is to insert some tension in the spline, that is, to apply a tensile force to the end-points of the spline. Appendix 6.A.3 lists the necessary details of this approach, using the classical *exponential tension spline* construction⁹ in Schweikert [1966]. When applied to the yield-curve setting, the construction involves a modification of the cubic equation (6.14) for $y(T)$, $T \in [T_i, T_{i+1}]$, to

$$\begin{aligned} y(T) = & \left(\frac{\sinh(\sigma(T_{i+1} - T))}{\sinh(\sigma h_i)} - \frac{T_{i+1} - T}{h_i} \right) \frac{y''_i}{\sigma^2} \\ & + \left(\frac{\sinh(\sigma(T - T_i))}{\sinh(\sigma h_i)} - \frac{T - T_i}{h_i} \right) \frac{y''_{i+1}}{\sigma^2} \\ & + y_i \frac{T_{i+1} - T}{h_i} + y_{i+1} \frac{T - T_i}{h_i}, \quad (6.16) \end{aligned}$$

where $\sigma \geq 0$ is the *tension factor*, and where we recall the definition $h_i = T_{i+1} - T_i$.

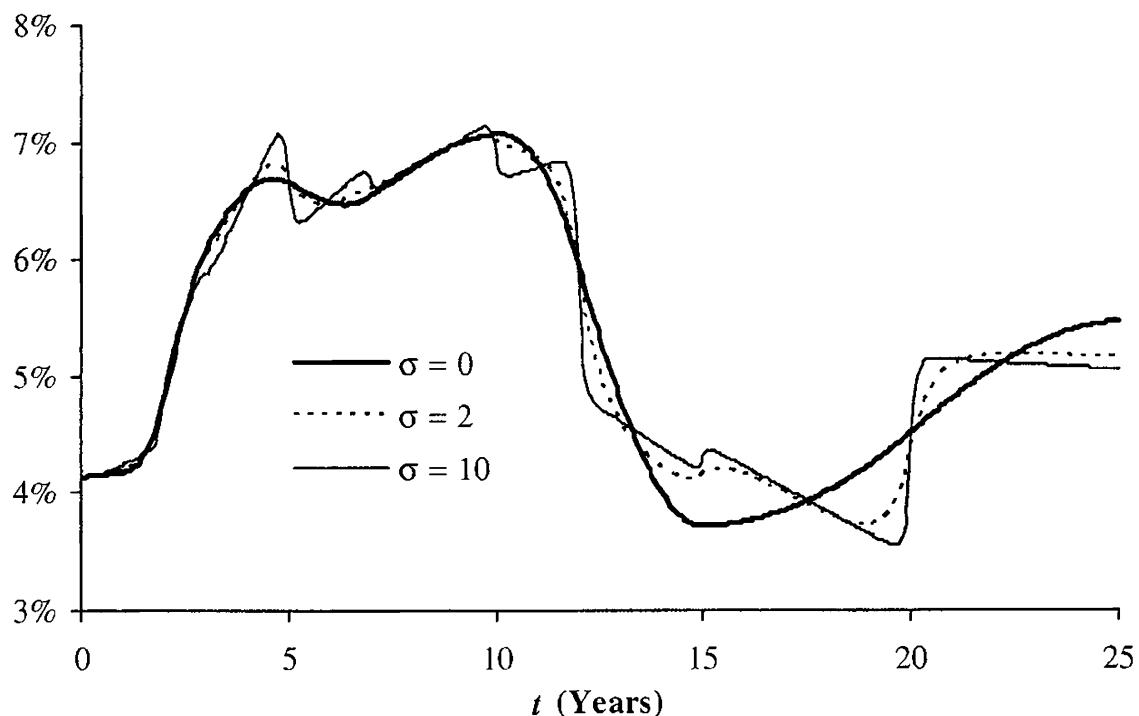
Appendix 6.A.3 discusses a number of properties of tension splines, the most important perhaps being the fact that setting $\sigma = 0$ will recover the ordinary C^2 cubic spline, whereas letting $\sigma \rightarrow \infty$ will make the tension spline uniformly approach a linear spline (i.e. the spline we used in Section 6.2.1.1). Loosely, we can thus think of a tension spline as a twice differentiable hybrid between a cubic spline and a linear spline. Equally loosely: as we increase σ , spurious inflections and ringing in the cubic spline are gradually “stretched” out of the curve, accompanied by rising (absolute values of) second derivatives at the knot points. More details on tension splines can be found in Andersen [2005], who also discusses application of computationally efficient local spline bases and the usage of a T -dependent tension factor to gain further control of the curve.

We observe that (6.16) is structurally similar to (6.14), and allows for a matrix representation of the same form as (6.62), albeit with suitably

⁹The exponential tension spline is not the only class of twice differentiable tension splines, but is probably the most common. Other classes are discussed in Kvasov [2000] and Andersen [2005].

modified definitions of the vector \mathbf{M} and the matrices \mathbf{B} and \mathbf{C} ; we leave these modifications as an exercise to the reader. Suffice to say that once a value of σ has been decided upon, the numerical search for the unknown levels y_i , $i = 1, \dots, N$, can proceed along the same principles as in Section 6.2.3 above. Figure 6.5 below shows an example; notice how increasing the tension parameter gradually moves us from cubic spline behavior to bootstrap behavior.

Fig. 6.5. Forward Curve



Notes: The yield curve is constructed as a C^2 natural tension spline, with tension parameters as given in the graph (only the forward curve $f(t)$ is shown). Swap data is in Table 6.1.

Remark 6.2.1. If the tension spline is applied not to yields, but to the logarithm of discount factors $\ln P(t)$, the limit of $\sigma \rightarrow \infty$ will produce a piecewise flat forward curve, as in Figure 6.2.

The reader may at this point wonder whether there are any firm rules as to what σ should be. We have no definitive answers to this question, and we do not try to determine σ automatically (although such routines do exist, see Renka [1987]). Instead, we normally treat σ as an “extra knob” that allow users to balance curve smoothness, shape preservation, and perturbation locality to their particular tastes. Inevitably some element of experimentation is required here.

6.3 Non-Parametric Optimal Yield Curve Fitting

The techniques we have outlined so far generally suffice for constructing a discount curve from a “clean” set of non-duplicate benchmark securities, including the carefully selected set of liquid staggered-maturity deposits, futures, and swaps most banks assemble for the purpose of constructing a Libor yield curve. In some settings, however, the benchmark set may be significantly less well-structured, involving illiquid securities with little order in their cash flow timing and considerable noise in their prices. This situation may, say, arise when one attempts to construct a yield curve from corporate bonds. While construction of a Libor curve is the most important task for the purposes of this book, we nevertheless wish to say a few words about techniques suitable for less cooperative benchmark security sets. These techniques can also be applied to Libor curve construction, of course, and are particularly relevant for applications where we are willing to sacrifice some precision in the fit to benchmark prices in return for a smoother yield curve.

6.3.1 Norm Specification and Optimization

When the input benchmark set is noisy, a direct solution of (6.4) may be erratic or may not exist. To overcome this, and to reflect that noise in the input data may make us content to solve (6.4) only to within certain error bounds, we now proceed to replace this equation with a problem of minimization of a penalized least-squares norm. Specifically, define the space $\mathcal{A} = C^2([t_1, t_M])$ of all functions $[t_1, t_M] \rightarrow \mathbb{R}$ that are twice differentiable with continuous second derivative, and introduce the M -dimensional discount vector

$$\mathbf{P}(y) = \left(e^{-y(t_1)t_1}, \dots, e^{-y(t_M)t_M} \right)^\top.$$

Also, let \mathbf{W} be a diagonal $N \times N$ weighting matrix. Then, as our best estimate \hat{y} of the yield curve we will here use

$$\hat{y} = \underset{y \in \mathcal{A}}{\operatorname{argmin}} \mathcal{I}(y), \quad (6.17)$$

with

$$\begin{aligned} \mathcal{I}(y) \triangleq \frac{1}{N} (\mathbf{V} - \mathbf{c}\mathbf{P}(y))^\top \mathbf{W}^2 (\mathbf{V} - \mathbf{c}\mathbf{P}(y)) \\ + \lambda \left(\int_{t_1}^{t_M} [y''(t)^2 + \sigma^2 y'(t)^2] dt \right), \end{aligned} \quad (6.18)$$

where λ and σ^2 are positive constants. The norm $\mathcal{I}(y)$ consists of three separate terms:

- A least-squares penalty term

$$\begin{aligned} \frac{1}{N} (\mathbf{V} - \mathbf{cP}(y))^\top \mathbf{W}^2 (\mathbf{V} - \mathbf{cP}(y)) \\ = \frac{1}{N} \sum_{i=1}^N W_i^2 \left(V_i - \sum_{j=1}^M c_{i,j} e^{-y(t_j) t_j} \right)^2, \end{aligned}$$

where W_i is the i -th diagonal element of \mathbf{W} . This term is an outright precision-of-fit norm and measures the degree to which the constructed discount curve can replicate input security prices. The weights W_i can be used to assign different importance to the various securities in the benchmark set, and/or to translate the precision of the fit from raw dollar amounts into more intuitive quantities, such as security-specific quoted yields¹⁰. Clearly, if (6.4) can be satisfied, then the least-squares penalty term will attain its minimum (of zero) for all yield curves that satisfy (6.4).

- A weighted smoothness term $\lambda \int_{t_1}^{t_M} y''(t)^2 dt$, penalizing high second-order derivatives of y to avoid kinks and discontinuities.
- A weighted curve-length term $\lambda \sigma^2 \int_{t_1}^{t_M} y'(t)^2 dt$, penalizing oscillations and excess convexity/concavity.

Our choice of calibration norm is, we believe, an attractive one, but other choices obviously are available as well. For instance, in Adams and van Deventer [1994] the norm contains no curve-length term and the smoothing norm is expressed on the forward curve, rather than on the yield curve. Due to the lack of the curve-length penalty term, the resulting curve will tend to behave like the C^2 cubic spline in Section 6.2.3; see Hagan and West [2004] for some numerical tests.

The following result is shown by variational methods in Andersen [2005]:

Proposition 6.3.1. *The curve \hat{y} that satisfies (6.17) is a natural exponential tension spline with tension factor σ and knots at all cash flow dates t_1, t_2, \dots, t_M .*

Proposition 6.3.1 establishes that the curve we are looking for is a tension spline with tension factor σ , but does not in itself allow us to identify the optimal spline directly, beyond the fact that i) it is a natural spline with boundary conditions $y''(t_1) = y''(t_M) = 0$; and ii) it has knots at all t_j , $j = 1, \dots, M$. Identification of the correct spline involves solving for unknown

¹⁰Most fixed-income securities are quoted through some type of yield, e.g. $V_i = g_i(r_i)$ where r_i is the quoted yield and g_i is a function that encapsulates the quoting convention. The quantity $D_i = -dg_i/dr_i$ is known as the *duration* of V_i . Setting $W_i = 1/D_i$ in the least-squares norm will turn price deviations into yield deviations.

levels¹¹ $y(t_1), y(t_2), \dots, y(t_M)$ to optimize directly (6.18). In this exercise, the following lemma is useful.

Lemma 6.3.2. *For a natural tension spline interpolating the values $y(t_1), y(t_2), \dots, y(t_M)$, we have*

$$\begin{aligned} y'(t_j) &= \left(-\sigma \frac{\cosh(\sigma(t_{j+1} - t_j))}{\sinh(\sigma(t_{j+1} - t_j))} + \frac{1}{t_{j+1} - t_j} \right) \frac{y''(t_j)}{\sigma^2} \\ &+ \left(\frac{\sigma}{\sinh(\sigma(t_{j+1} - t_j))} - \frac{1}{t_{j+1} - t_j} \right) \frac{y''(t_{j+1})}{\sigma^2} + \frac{y(t_{j+1})}{t_{j+1} - t_j} - \frac{y(t_j)}{t_{j+1} - t_j}, \end{aligned}$$

and

$$\lambda \left(\int_{t_1}^{t_M} [y''(t)^2 + \sigma^2 y'(t)^2] dt \right) = -\lambda \sum_{j=1}^{M-1} d_j (y(t_{j+1}) - y(t_j)), \quad (6.19)$$

where $y''(t_1) = y''(t_M) = 0$, and

$$d_j \triangleq \frac{y''(t_{j+1}) - \sigma^2 y(t_{j+1})}{t_{j+1} - t_j} - \frac{y''(t_j) - \sigma^2 y(t_j)}{t_{j+1} - t_j}. \quad (6.20)$$

Proof. The result for $y'(t_j)$ follows from direct differentiation of the basic equations for a tension spline (see (6.16) above, applied to the knot grid $\{t_j\}$). To show (6.19), consider the interval $[t_j, t_{j+1}]$ and the integral

$$\int_{t_j}^{t_{j+1}} (y''(t)^2 + \sigma^2 y'(t)^2) dt = \int_{t_j}^{t_{j+1}} (y''(t) \cdot y''(t) + \sigma^2 y'(t) \cdot y'(t)) dt.$$

Integration by parts yields

$$\begin{aligned} &\int_{t_j}^{t_{j+1}} (y''(t)^2 + \sigma^2 y'(t)^2) dt \\ &= [y''(t)y'(t)]_{t_j}^{t_{j+1}} - \int_{t_j}^{t_{j+1}} (y^{(3)}(t) - \sigma^2 y'(t)) y'(t) dt \\ &= y''(t_{j+1})y'(t_{j+1}) - y''(t_j)y'(t_j) - d_j (y(t_{j+1}) - y(t_j)), \end{aligned} \quad (6.21)$$

where d_j is given in (6.20). Here, we have used that, by definition, hyperbolic tension splines have $y^{(3)}(t) - \sigma^2 y'(t)$ piecewise constant and equal to d_j on each interval $[t_j, t_{j+1}]$ (see equation (6.63) in Appendix 6.A). The result (6.20) follows by addition of the terms (6.21) and using the condition $y''(t_1) = y''(t_M) = 0$. \square

¹¹In Andersen [2005], the search for yield levels has been replaced by the more contemporary idea of searching for weights in a local basis representation of the spline.

Lemma 6.3.2 shows us that we can compute the value of the integral penalty term in (6.18) directly from knowledge of yield levels $y(t_1), \dots, y(t_M)$ and second derivatives $y''(t_2), \dots, y''(t_{M-1})$. For each guess for the M unknown levels $y(t_1), y(t_2), \dots, y(t_M)$ we can proceed as follows.

1. Compute the least-squares penalty term $\frac{1}{N}(\mathbf{V} - \mathbf{cP}(y))^\top \mathbf{W}^2 (\mathbf{V} - \mathbf{cP}(y))$ directly from the definition of $\mathbf{P}(y)$.
2. Use the results in Section 6.2.4 to solve for $y''(t_2), \dots, y''(t_{M-1})$ by solving a tri-diagonal set of equations.
3. Use Lemma 6.3.2 to compute $\lambda(\int_{t_1}^{t_M} [y''(t)^2 + \sigma^2 y'(t)^2] dt)$, thereby completing the computation of the norm $\mathcal{I}(y)$.

Embedding Steps 1–3 above in a multi-variate numerical optimizer ultimately allows us to determine the optimal solution \hat{y} . A good generic routine for this optimization step would be the Levenberg-Marquardt algorithm; see Press et al. [1992]. The optimization problem at hand is generally good-natured, and one can also use a simpler Gauss-Newton method, as discussed in Andersen [2005]. If possible, it is often useful to use a simpler method (e.g. bootstrapping) to establish a good guess for the yield curve levels $y(t_1), y(t_2), \dots, y(t_M)$. A proper implementation of the algorithm should typically construct a yield curve in less than one-tenth of a second on a standard PC.

Remark 6.3.3. If we let $\sigma = 0$, the solution to the optimization problem becomes a *cubic smoothing spline*; see Tanggaard [1997] for more details on this case.

Remark 6.3.4. If we let $\lambda \downarrow 0$, the resulting spline will often end up hitting all benchmark prices exactly, i.e. will satisfy (6.4) in the limit. The resulting spline is then the optimal *interpolating* curve, in the sense that of all twice differentiable yield curves that match the benchmark prices, the spline is the minimizer of the regularity term $\int_{t_1}^{t_M} [y''(t)^2 + \sigma^2 y'(t)^2] dt$. If, for $\lambda \downarrow 0$, we do not satisfy (6.4), then the resulting spline can be considered a *least-squares regression spline* solution.

6.3.2 Choice of λ

So far, we have assumed that the parameter λ has been specified exogenously by the user. In practice, however, a good magnitude of λ may sometimes be hard to ascertain by inspection, and a procedure to estimate λ directly from the data is often useful. One possibility is to use a cross-validation approach, either outright or through the more efficient Generalized Cross-Validation (GCV) criterion by Craven and Wahba [1979]. Some results along these lines can be found in Tanggaard [1997] and Andersen [2005], but are outside the scope of our treatment here. A more pragmatic approach is to replace the optimization problem (6.17) with the constrained optimization problem

$$\hat{y} = \operatorname{argmin}_{y \in \mathcal{A}} \int_{t_1}^{t_M} [y''(t)^2 + \sigma^2 y'(t)^2] dt, \quad (6.22)$$

$$\text{subject to } \frac{1}{N} (\mathbf{V} - \mathbf{cP}(y))^\top \mathbf{W}^2 (\mathbf{V} - \mathbf{cP}(y)) = \gamma^2, \quad (6.23)$$

where γ is an exogenously specified constant. Note that γ is just the allowed weighted root-mean-square (RMS) error in the fit to benchmark securities, an intuitive quantity that most users should have no problem specifying directly based on, say, observed bid-offer spreads. The Lagrangian for the above problem becomes

$$\begin{aligned} \hat{y} = \operatorname{argmin}_{y \in \mathcal{A}} & \left(\int_{t_1}^{t_M} [y''(t)^2 + \sigma^2 y'(t)^2] dt \right. \\ & \left. + \rho \left[\frac{1}{N} (\mathbf{V} - \mathbf{cP}(y))^\top \mathbf{W}^2 (\mathbf{V} - \mathbf{cP}(y)) - \gamma^2 \right] \right), \end{aligned} \quad (6.24)$$

where the Lagrange multiplier ρ must be determined such that the constraint (6.23) is satisfied at the optimum of (6.24). Apart from a constant scale, (6.24) is identical to (6.17), so we solve the constrained optimization problem (6.22)–(6.23) through the following iteration over λ :

1. Given a guess for λ , find the optimum value of $y(t_1), y(t_2), \dots, y(t_M)$, as a solution of (6.17).
2. Compute $\mathcal{S} = \frac{1}{N} (\mathbf{V} - \mathbf{cP}(y))^\top \mathbf{W}^2 (\mathbf{V} - \mathbf{cP}(y))$.
3. If $\mathcal{S} = \gamma^2$, stop; otherwise update λ and go to Step 1.

In Step 1, we can proceed as discussed in Section 6.3.1 above. In general, the precision norm $\mathcal{S} = \mathcal{S}(\lambda)$ will be a declining function in λ and, provided that a root to $\mathcal{S}(\lambda) = \gamma^2$ exists¹², the updating in Step 3 can be done by any standard root search algorithm.

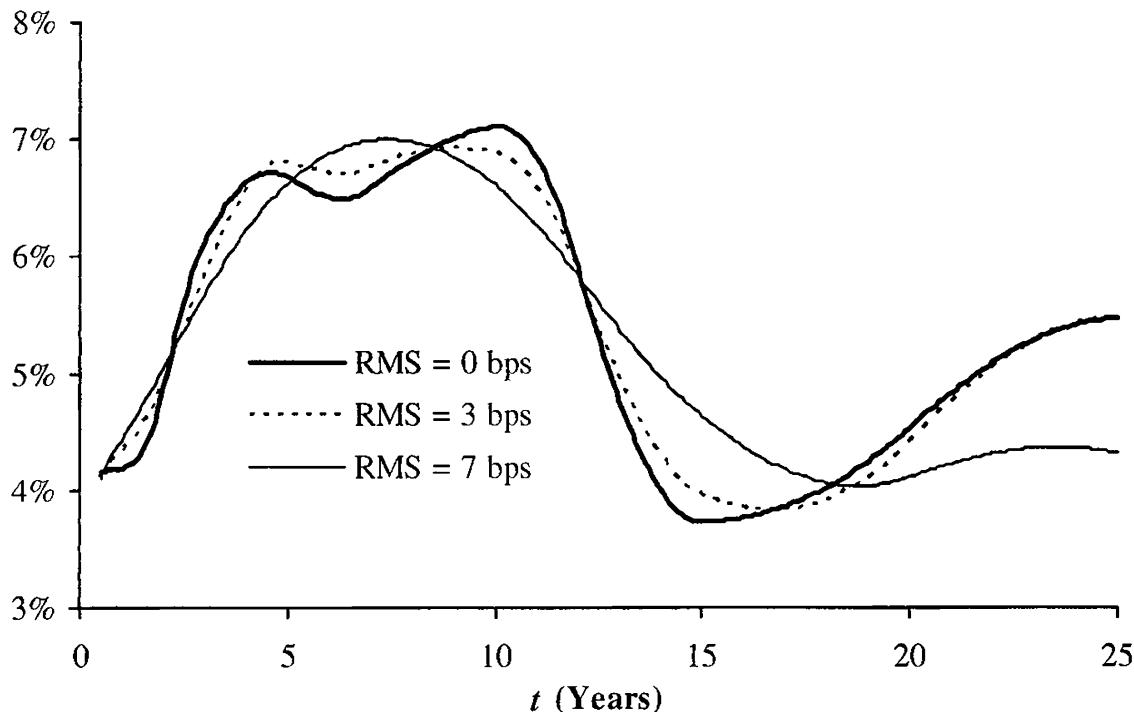
6.3.3 Example

To illustrate the effect of λ , we now apply the algorithm in Section 6.3.2 to the test data in Table 6.1 above. In doing so, we use the matrix \mathbf{W} to normalize (see footnote 10) all price errors to yield-to-maturity errors, allowing us to consider γ in (6.23) as the root-mean-square (RMS) yield error. Setting $\sigma = 0$, the forward curves for various choices of γ are shown in Figure 6.6. As one would expect, the higher we allow γ to be, the smoother the forward (and yield) curves become.

For our test case, the zero-RMS optimal (M -knot) forward curve in Figure 6.6 is virtually identical to the N -knot cubic spline solution in Figure 6.4. In general, the N -knot interpolating curve can be interpreted as a

¹²There may be instances where $\mathcal{S}(0) > \gamma^2$. If the desired precision is unattainable, we can either increase γ^2 or perhaps prune the benchmark security set.

Fig. 6.6. Forward Curve



Notes: The yield curve is constructed as an optimal C^2 natural tension spline, with an RMS yield error constraint as listed in the graph (only the forward curve $f(t)$ is shown). The tension parameter is set to $\sigma = 0$ for all curves. Swap data is in Table 6.1.

constrained solution to (6.17) with $\lambda = 0$, with the constraint requiring that knots be placed only at benchmark maturities $\{T_i\}_{i=1}^N$, rather than at all cash flow dates $\{t_j\}_{j=1}^M$. The effect of enforcing this additional constraint is often rather small, at least for the purposes of constructing a Libor curve.

6.4 Managing Yield Curve Risk

Consider a portfolio of securities with value V_0 , where V_0 is a function of the yield curve $y(t)$. The securities in V_0 would typically not be in the benchmark set and could contain, say, interest rate options, seasoned swaps, and so forth. As the yield curve is a function of the benchmark set values $\mathbf{V} = (V_1, \dots, V_N)^\top$, we may write

$$V_0 = V_0(V_1, \dots, V_N; \theta),$$

where the vector θ contains model parameters (e.g. volatilities) and where the function $V_0(\cdot)$ is determined both from the valuation model of the security in question, and from the curve construction algorithm employed. Clearly, then

$$dV_0 = \sum_{i=1}^N \frac{\partial V_0}{\partial V_i} dV_i + \sum_i \frac{\partial V_0}{\partial \theta_i} d\theta_i,$$

or, for non-infinitesimal moves,

$$\Delta V_0 \approx \sum_{i=1}^N \frac{\partial V_0}{\partial V_i} \Delta V_i + \sum_i \frac{\partial V_0}{\partial \theta_i} \Delta \theta_i. \quad (6.25)$$

For the purpose of managing first-order risk exposure to moves in the yield curve, (6.25) suggests that the collection of derivatives $\partial V_0 / \partial V_i$, $i = 1, \dots, N$ — often called *(bucketed) interest rate deltas* — forms a natural metric for portfolio risk. In particular, if all these derivatives are zero, our portfolio would, to first order, be immunized against any move in the yield curve that is consistent with the chosen curve construction algorithm. On the other hand, if some or all of the derivatives are non-zero, we could manage our risk by setting up a hedge portfolio of benchmark securities, with notional $-\partial V_0 / \partial V_i$ on the i -th security. We emphasize that the resulting hedge would typically *not* be model-consistent: most interest models assume that yield curve risk originates from only a few stochastic yield curve factors that tend to move the curve smoothly¹³, in a predominantly parallel fashion. Theoretically, a bucket-by-bucket immunization against all terms ΔV_i may then be considered an overkill — we typically hedge against far too many risk factors (N) — but is nevertheless standard industry practice and has proven to be robust. Notice that bucket hedging along these lines would, for instance, correctly reject the notion that we could perfectly hedge a 20 year swap with a 1 month FRA, something that a one-factor interest rate model (see Chapter 4 and Chapter 10) would happily accept. We pick up this subject again in Chapter 22.

6.4.1 Par-Point Approach

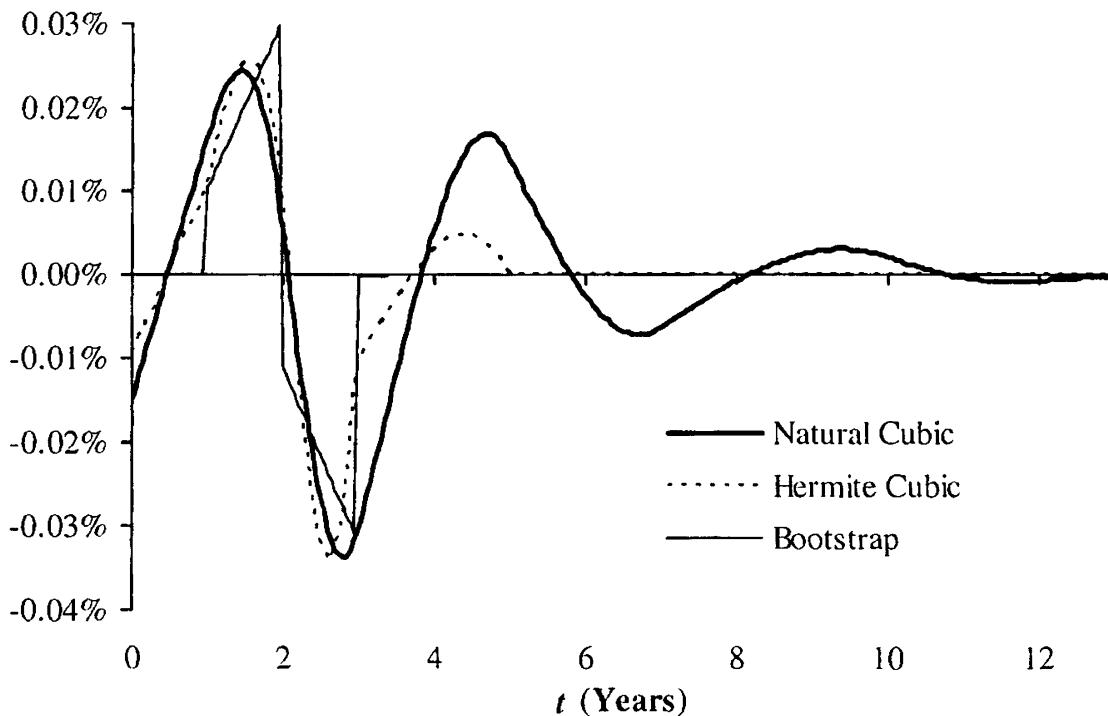
The simplest approach to computation of the delta $\partial V_0 / \partial V_i$ involves a manual bump¹⁴ to V_i , followed by a reconstruction of the yield curve, and a subsequent repricing of the portfolio V_0 . This procedure is sometimes known as the *par-point approach*, and resulting derivatives *par-point deltas*. For the approach to work properly, it is important that the yield curve construction algorithm is fast and produces clean, local perturbations of the yield curve when benchmark prices are shifted. For instance, perturbing a short-dated FRA price should not cause noticeable movements in long-term yields, lest we reach the erroneous conclusion (again) that we can perfectly hedge a 20 year swap with a 1 month FRA. As we have discussed earlier, Hermite splines and bootstrapped yield curves both exhibit good perturbation locality, but cubic C^2 splines often do not. To illustrate this, Figure 6.7 considers the

¹³See the principal components analysis in Chapter 14 for more on this.

¹⁴In practice, rather than bumping the price V_i outright, one may instead bump the yield of the i -th benchmark security (typically by 1 basis point). See also footnote 10.

effect on the forward curves in Figures 6.1, 6.3, and 6.4 from a 1 basis point up-move in the par yield of the 2 year swap in Table 6.1. As we can see, the move causes a noisy, ringing perturbation in the C^2 cubic spline solution, spreading into short- and long-dated parts of the forward curve.

Fig. 6.7. Forward Curve Move



Notes: Change in instantaneous forward curve, from a 1 basis point shift in the 2 year swap yield in Table 6.1. The curve construction methods tested are: bootstrapping with piecewise linear yields ("Bootstrap"), Hermite C^1 cubic spline ("Hermite"), and C^2 natural cubic spline ("Natural Cubic"). Swap data is in Table 6.1.

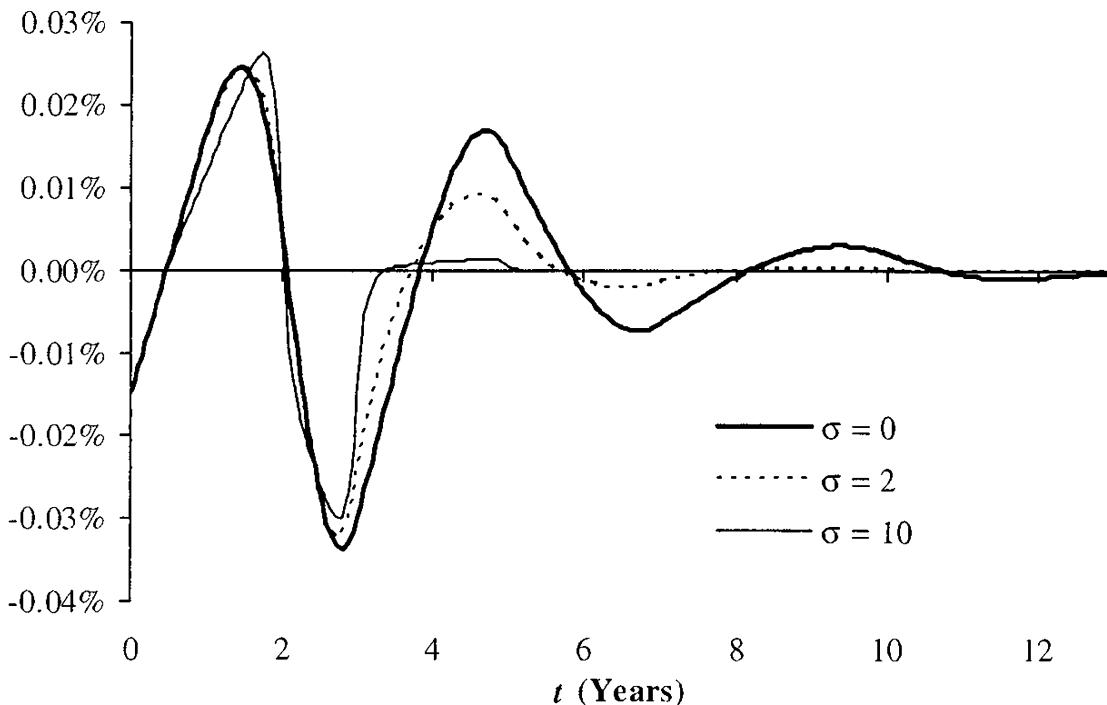
In Figure 6.8, we have followed the recommendations of Section 6.2.4 and added tension to the C^2 spline, causing a dampening of the perturbation noise. Clearly, the usage of a tension factor can have a beneficial impact on risk reports produced by the par-point approach.

6.4.2 Forward Rate Approach

As an alternative to direct perturbation of benchmark security prices, we can consider applying perturbations directly to the discount curve, thereby mostly avoiding the introduction of artifacts specific to the curve construction algorithm. In practice, this technique typically focuses on the forward curve¹⁵ $f(t)$, to which we apply certain functional shifts $\mu_k(t)$, $k = 1, \dots, K$. Writing

¹⁵Perturbations may also be performed on discretely, rather than continuously, compounded forward rates.

Fig. 6.8. Forward Curve Move



Notes: Change in instantaneous forward curve, from a 1 basis point shift in the 2 year swap yield in Table 6.1. The yield curve was constructed as a tension spline, with tension factors as given in the graph. Swap data is in Table 6.1.

(loosely) $V_0 = V_0(f)$ to highlight the dependence of V_0 on the forward curve, we then compute functional (Gâteaux) derivatives¹⁶ for V_0 :

$$\partial_k V_0 = \left. \frac{dV_0(f(t) + \varepsilon \mu_k(t))}{d\varepsilon} \right|_{\varepsilon=0}, \quad k = 1, \dots, K. \quad (6.26)$$

Standard choices for $\mu_k(t)$ are

$$\begin{aligned} \text{Piecewise Triangular: } \mu_k(t) &= \frac{t - t_{k-1}}{t_k - t_{k-1}} 1_{\{t \in [t_{k-1}, t_k)\}} \\ &\quad + \frac{t_{k+1} - t}{t_{k+1} - t_k} 1_{\{t \in [t_k, t_{k+1})\}}, \end{aligned} \quad (6.27)$$

$$\text{Piecewise Flat: } \mu_k(t) = 1_{\{t \in [t_k, t_{k+1})\}}, \quad (6.28)$$

where $\{t_k\}$ is a user-specified discretization grid. The resulting sensitivities are often called *forward rate deltas*.

It is common practice to use $\{t_k\}$ grids spaced three months apart, with dates on Eurodollar futures maturities. The number of deltas K is thus typically a rather large number, and the K derivatives $\partial_k V_0$ give a detailed picture of where the portfolio risk is concentrated on the forward curve. As forward rate contracts and Eurodollar futures cease to be liquid beyond 4

¹⁶For a proper definition of the Gâteaux derivative, see Gâteaux [1913].

or 5 year maturities, the forward rate deltas do not directly suggest hedging instruments for the medium and long end of the yield curve exposure; however it is not difficult to translate forward rate deltas into a hedging portfolio (see the next section). The choice of par point versus forward rate deltas is largely a matter of personal preference, and it is not uncommon for traders to use both at the same time.

6.4.3 From Risks to Hedging: The Jacobian Approach

A collection of forward rate shifts $\mu_k(t)$, $k = 1, \dots, K$, defines a certain view on the (first-order) risk on the portfolio $V_0(f)$ via the functional derivatives (6.26). In order to be useful, this risk view ultimately needs to be translated into a portfolio of hedging instruments that offsets the risks of V_0 . While fixed income traders normally are quite adept at mentally translating forward rate risk into actual hedge transactions, some linear algebra can help out in this exercise, as we now show.

Suppose a set of L hedging instruments is available, with values $\mathbf{H} = (H_1, \dots, H_L)^\top$. This set may or may not coincide with the benchmark set used for curve construction; for example, one may want to exclude some benchmark securities from the hedging set due to poor liquidity, or one may want to add instruments to the benchmark set to fine-tune hedging. Using (6.26), we denote the sensitivities of hedging instruments to the shifts $\mu_k(t)$ by $\partial_k H_l$, $l = 1, \dots, L$, $k = 1, \dots, K$. If the l -th hedging instrument is included in the hedging portfolio with notional weight p_l , and $\mathbf{p} = (p_1, \dots, p_L)^\top$, then the sensitivity of the hedge portfolio value

$$H_0(\mathbf{p}) = \mathbf{p}^\top \mathbf{H}$$

to the k -th perturbation is given by

$$\partial_k H_0(\mathbf{p}) = \mathbf{p}^\top \partial_k \mathbf{H},$$

where we have denoted

$$\partial_k \mathbf{H} = (\partial_k H_1, \dots, \partial_k H_L)^\top.$$

In most cases¹⁷ we would like to choose the weights \mathbf{p} in such a way that $\partial_k H_0(\mathbf{p})$ offsets as much of $\partial_k V_0$ as possible, for all $k = 1, \dots, K$. Let W_k be the relative importance of offsetting the k -th derivative, and U_l a relative “reluctance” to using the l -th hedging instrument (a function of the bid-ask spread, for example). Then, the optimal hedging weights $\hat{\mathbf{p}}$ can be defined by the condition

¹⁷Sometimes traders deliberately wish to keep some risk on their books, as a way to speculate on interest rate movements. A non-zero target risk profile is easily accommodated by a change of the optimization target in (6.29).

$$\hat{\mathbf{p}} = \underset{\mathbf{p}}{\operatorname{argmin}} \left(\sum_{k=1}^K W_k^2 (\partial_k H_0(\mathbf{p}) - \partial_k V_0)^2 + \sum_{l=1}^L U_l^2 p_l^2 \right). \quad (6.29)$$

Define the matrix $\partial\mathbf{H}$ to have columns $\partial_1\mathbf{H}, \dots, \partial_K\mathbf{H}$, the vector $\partial\mathbf{V}_0$ by

$$\partial\mathbf{V}_0 = (\partial_1 V_0, \dots, \partial_K V_0)^\top,$$

the matrix \mathbf{W} to be diagonal with W_k 's on the diagonal, and the matrix \mathbf{U} to be diagonal with U_l 's on the diagonal. With this notation (6.29) can be recast as a least-squares problem,

$$(\partial\mathbf{H}^\top \mathbf{p} - \partial\mathbf{V}_0)^\top \mathbf{W}^2 (\partial\mathbf{H}^\top \mathbf{p} - \partial\mathbf{V}_0) + \mathbf{p}^\top \mathbf{U}^2 \mathbf{p} \rightarrow \min. \quad (6.30)$$

The problem (6.30) can be solved by standard methods; a formal solution is given by the linear system

$$(\partial\mathbf{H} \mathbf{W}^2 \partial\mathbf{H}^\top + \mathbf{U}^2) \hat{\mathbf{p}} = \partial\mathbf{H} \mathbf{W}^2 \partial\mathbf{V}_0. \quad (6.31)$$

We note that the addition of the \mathbf{U} term to the optimization problem (6.31) is sometimes called *Tikhonov regularization*, a technique that we shall return to in Chapter 18.

When solving (6.30), one should carefully consider the relative dimensions of the matrices involved. First, if there are fewer hedging instruments than shifts to be immunized ($L < K$), then, in general, not all risks can be offset. In this case, the weights \mathbf{W} gain in importance as they allow the user to focus hedging on risk buckets deemed more important, at the expense of other, less critical ones. Also, when $L < K$ the weights \mathbf{U} are less important, and in most cases can safely be set to zero. Second, if there are more hedging instruments than risk buckets to immunize against ($L > K$), then there are typically multiple hedging portfolios that perfectly offset all risks. In this case, the weights \mathbf{W} can normally be ignored (all set to 1), but the weight matrix \mathbf{U} becomes more critical as it allows one to choose which of the possible hedging portfolios one “likes” best (e.g., the least costly). Finally, if $L = K$, then normally there exists exactly one portfolio that hedges all risks. Both \mathbf{W} and \mathbf{U} are then often of little consequence, although one might still want to specify non-zero weights \mathbf{U} to avoid oscillatory or unstable solutions if the linear equations are ill-posed. We note that in the simple case of $L = K$, $\mathbf{W} = \mathbf{1}$, $\mathbf{U} = \mathbf{0}$ and $\partial\mathbf{H}$ invertible, the solution to the optimization problem is given by

$$\mathbf{p} = (\partial\mathbf{H}^\top)^{-1} \partial\mathbf{V}_0. \quad (6.32)$$

The method of constructing a hedge portfolio from derivatives to arbitrary shocks of the forward curve via the optimization problem (6.30) is known as the *Jacobian method* for interest rate deltas; the name originates from the fact that the matrix $\partial\mathbf{H}$ is the Jacobian matrix of a hedge set with respect to the forward curve shocks. Combined with the forward rate deltas,

the Jacobian method helps aggregate fine-grained risks into various sets of hedges. The approach has considerable generality as the risk basis functions μ_k and the hedge portfolio can be chosen freely by the user — note, for instance, that even the par-point approach can be seen as a special type of the Jacobian method where we effectively choose the hedging set to coincide perfectly with the benchmark set and where the μ_k 's are set to be the shifts of the forward curve that correspond to the bumps of benchmark securities. In this special case, the Jacobian $\partial \mathbf{H}$ is then a unit matrix and (from (6.32)) the original par-point deltas are recovered.

The Jacobian method serves to decouple risk calculations from curve construction. This, potentially, allows for combining smooth curves with localized risk, a feat that is difficult to achieve by other methods. The Jacobian is also useful in applications where curves need to be rebuild over and over, to address the fact that Libor and Treasury benchmark security prices (or yields) change very quickly, often quicker than a sophisticated curve construction algorithm can rebuild the curves. With the aid of a Jacobian, changes in benchmark prices can be quickly translated into changes of the forward curve via a matrix multiplication. A full curve rebuild needs only be triggered when the benchmark prices have moved sufficiently far from their initial values.

6.4.4 Cumulative Shifts and other Common Tricks

As evident in Figure 6.7 (the bootstrap case), a shift to a single swap rate (while keeping other swap rates fixed) typically results in a strong “see-saw” impact on the forward rate curve. Let us attempt to gain some intuition about the magnitude of the forward rate shock. For a back-of-the-envelope calculation, we can assume that a swap rate is a linear combination of forward Libor rates (see (4.11)),

$$S_n \approx \sum_{i=1}^n w_{i,n} L_i,$$

where S_n denotes a swap rate for a swap covering n periods (for simplicity assume that each period is 1 year), L_i denotes a forward Libor rate for i -th period (from year $i-1$ to year i), and $w_{i,n} \approx 1/n$. Inverting this relationship yields

$$L_n \approx n \left(S_n - \frac{n-1}{n} S_{n-1} \right). \quad (6.33)$$

As part of a par-point report, assume now that S_n is shifted by the amount δ , but S_{n-1} and S_{n+1} remain unchanged. According to (6.33) L_n will then shift by approximately $n\delta$, and L_{n+1} by $-n\delta$. For instance, if a 30 year swap yield is shifted by 1 basis point, while 29 year and 31 year are kept unchanged, then evidently the forward Libor rate L_{30} will move by 30

basis points, and the rate L_{31} will move by -30 basis points. If the portfolio whose deltas we are computing happens to contain, say, a spread option on the difference between L_{30} and L_{31} , the underlying rate of this option would be shifted by 60 basis points (!). And clearly, a shift of 60 basis points (or 30 basis points, for that matter) is not small, and may be inappropriate for calculating a first-order derivative. We emphasize that what appears to be a benign 1 basis point rate shift translates into a much larger forward curve move that can potentially affect underlying instruments in unexpected ways.

The example above highlights the importance of applying shifts to the forward curve that are consistent with real moves of interest rates. Obviously, it is highly unlikely that a forward curve would move in such a way that a 30 year swap rate has changed but the 29 and 31 year rates have not.

One tweak to the standard par-point approach that goes some way towards the goal of realistic curve shifts is the so-called *cumulative par-point approach* (also known as a *waterfall par-point approach*). The idea is simple: the shift to the i -th benchmark security is retained while calculating the derivative to the $(i+1)$ -th (and subsequent) securities. In other words, the two curves for the $(i+1)$ -th derivative are constructed from the prices

$$(V_1 + \Delta V_1, \dots, V_i + \Delta V_i, V_{i+1}, V_{i+2}, \dots, V_N)$$

(base) and

$$(V_1 + \Delta V_1, \dots, V_i + \Delta V_i, V_{i+1} + \Delta V_{i+1}, V_{i+2}, \dots, V_N)$$

(perturbed). The standard deltas are then computed as differences of two consecutive (cumulative) derivatives. While the resulting deltas should coincide with the standard par-point deltas in the limit of $\Delta V \rightarrow 0$, they differ for non-vanishing perturbations.

The forward curve shifts implied by the cumulative par-point method are less extreme than those of the ordinary par-point method, making the cumulative par-point method quite attractive in practice. Another practical advantage of the method is the fact that the sum of deltas computed by the method is always (by definition) exactly equal to the *parallel delta*, i.e. the delta that is obtained by shifting all benchmark yields by the same amount at the same time. Because of the second-order effects, the same is only true for the standard par-point method in the limit of vanishing shifts, not for the non-infinitesimal perturbations used in practice.

The cumulative par-point approach is easy to mimic (and even improve) in the Jacobian framework of Section 6.4.3. Clearly, from (6.33), the i -th cumulative shift roughly corresponds to a piecewise flat move of the forward curve between the maturities of $(i-1)$ -th and i -th benchmark. Hence, we can define

$$\mu_i(t) = 1_{\{t \in [T_{i-1}, T_i)\}}, \quad i = 1, \dots, N, \quad (6.34)$$

with $T_0 \triangleq 0$. Note that this specification involves benchmark maturities $\{T_i\}$, in contrast to (6.28) which is typically set on a 3 month grid; in

particular, (6.34) involves as many shocks as there are benchmark securities. Application of the Jacobian method to (6.34) yields an attractive variation of the cumulative par-point method where all forward curve shocks are similarly scaled, in contrast to the basic cumulative par-point where the size of forward curve shocks grows linearly with maturity, as implied by (6.33).

We should note that to improve accuracy, one may compute deltas as the average of deltas computed using first positive shocks, then negative shocks. This idea applies to par-point, forward-rate, Jacobian, cumulative-par-point, or any other delta calculation method. For the simple par-point method, this boils down to using two-sided finite difference approximations versus one-sided for approximating derivatives, a standard trick. For other methods the relationship is not as straightforward but the end result is the same: improved accuracy and stability of deltas. Using averaged deltas is typically particularly useful for security prices that depend on yields in a strongly non-linear fashion.

Finally, let us mention another popular trick. We have spent a good part of Section 6.2.4 describing ways to build smooth yield curves that exhibit good locality under perturbations. A more simplistic approach to tackle the same problem is based on building two different curves. One — smooth — is then used for pricing and the other — bootstrapped and with good locality — used for risk computations. While certainly helpful in a pinch, this approach tends to suffer from poor P&L predict, in the sense that changes in valuations of a portfolio between two dates are not well explained by first-order sensitivities (because values and sensitivities are calculated using different curves). We spend more time on P&L predict in Chapter 22.

6.5 Various Topics in Discount Curve Construction

6.5.1 Curve Overlays and Turn-of-Year Effects

Many of the curve construction algorithms so far have been designed around the implicit idea that the forward curve should ideally be *smooth*. While this is, indeed, generally a sound principle, exceptions do exist. For instance, it may be reasonable to expect instantaneous forwards to jump on or around meetings of monetary authorities, such as the Federal Reserve in the US. In addition, other “special” situations may exist that might warrant introduction of discontinuities into the forward curve. A well-known example is the turn-of-year (TOY) effect where short-dated loan premiums spike for loans between the last business day of the year and the first business day of the following calendar year.

One common way of incorporating TOY-type effects is to exogenously specify an *overlay curve* $\varepsilon_f(t)$ on the instantaneous forward curve. Specifically, the forward curve $f(t) = f(0, t)$ is written as

$$f(t) = \varepsilon_f(t) + f^*(t), \quad (6.35)$$

where $\varepsilon_f(t)$ is user-specified — and most likely contains discontinuities around special event dates — and $f^*(t)$ is unknown. The yield curve algorithm is then subsequently applied to the construction of $f^*(t)$. That is, rather than solving $\mathbf{cP} = \mathbf{V}$ (see equation (6.4)), we instead write

$$P(T) = e^{-\int_0^T \varepsilon_f(t) dt} e^{-\int_0^T f^*(t) dt} \triangleq P_\varepsilon(T) P^*(T) \quad (6.36)$$

and solve

$$\mathbf{c}_\varepsilon \mathbf{P}^* = \mathbf{V}, \quad (6.37)$$

where $\mathbf{P}^* = (P^*(t_1), \dots, P^*(t_M))^\top$, and \mathbf{c}_ε is a modified $N \times M$ coupon matrix, with elements

$$(\mathbf{c}_\varepsilon)_{i,j} = c_{i,j} P_\varepsilon(t_j).$$

Construction of \mathbf{c}_ε can be done as a pre-processing step, after which any of the algorithms discussed earlier in this chapter can be applied to attack (6.37). Once the curve $P^*(t)$ (or, equivalently, the yield curve $y^*(t) = -t^{-1} \ln P^*(t)$) has been constructed, any subsequent use of the curve for cash flow discounting requires, according to (6.36), a multiplicative adjustment of time t discount factors by the quantity $P_\varepsilon(t)$.

6.5.2 Cross-Currency Curve Construction

In this section we consider the issues involved in constructing yield curves simultaneously in multiple currencies. As it turns out, the market for foreign exchange (FX) forwards and cross-currency basis swaps imposes certain arbitrage constraints that must be considered in the curve construction exercise.

6.5.2.1 Basic Problem

To provide some motivation, consider a US dollar (USD) based firm receiving \$1 for certain at some future time T . Assuming that we have available a risk-free discount curve $P(\cdot)$ for USD-denominated cash flows, we compute the value of this security simply as $P_\$(T)$. Suppose now that the firm enters into a (costless) FX forward where it commits to pay \$1 at time T against receipt of a Japanese yen (JPY) amount of ¥ $Y(T)$; $Y(T)$ thereby represents the time 0 forward JPY/USD exchange rate for delivery at time T . By transacting the FX forward, the firm has effectively turned the receipt of \$1 into receipt of ¥ $Y(T)$, the USD PV at time 0 of which is

$$X(0) P_\$(T) Y(T),$$

where $P_\$(T)$ is a JPY discount factor and $X(0)$ is the time 0 foreign exchange rate in \$/¥ terms. To avoid an arbitrage, we evidently need

$$P_{\$}(T) = X(0)P_{\text{¥}}(T)Y(T) \quad \Rightarrow \quad P_{\text{¥}}(T) = \frac{P_{\$}(T)}{Y(T)X(0)}. \quad (6.38)$$

Suppose, say, that we have blindly estimated discount curves $P(\cdot)$ and $P_{\text{¥}}(\cdot)$ from the market for USD- and JPY-denominated interest rate swaps, respectively, without paying any attention to FX markets. The discount curves $P_{\$}(\cdot)$ and $P_{\text{¥}}(\cdot)$ estimated in this fashion will very likely *not* satisfy (6.38), implying the existence of cross-currency arbitrages. The degree to which (6.38) is typically violated is often small, but any such violation can be highly problematic for a firm engaging in trading of significant amounts of both USD- and JPY-denominated assets.

6.5.2.2 Separation of Discount and Forward Rate Curves

It may appear that there is no way out of this conundrum: after all, our curve construction algorithms imply a unique discount curve out of given swap prices and have few, if any, means of incorporating additional requirements such as (6.38). However, built into our assumptions about how to price a swap (see (6.2)) was an implicit assumption that Libor itself is the proper discount rate for flows based on Libor fixings. As Libor rates represent lending rates between banks, they contain a certain amount of credit risk¹⁸ and it is ex-ante unclear that they are suitable proxies for a “risk-free” rate (or, at least, are suitable for discounting of swap cash flows). More details about this can be found in Collin-Dufresne and Goldstein [2001] and Duffie and Huang [1996]. For our purposes, it suffices to introduce the notion that when computing a swap value we may need two curves: i) the Libor “pseudo-discount” curve $P^{(L)}(t) = P^{(L)}(0, t)$, used to project the Libor-based floating cash flows on the floating leg of the swap; and ii) a real discount curve $P(t) = P(0, t)$, used to discount all cash flows. For, say, a regular JPY-based fixed-floating swap paying a coupon c on a schedule $\{t_i\}_{i=1}^n$, the swap valuation equation thus becomes

$$V_{\text{¥}}(0) = \underbrace{\sum_{i=0}^{n-1} c\tau_i P_{\text{¥}}(0, t_{i+1})}_{\text{Fixed Leg}} - \underbrace{\sum_{i=0}^{n-1} L_{\text{¥}}(0, t_i, t_{i+1})\tau_i P_{\text{¥}}(0, t_{i+1})}_{\text{Floating Leg}}, \quad (6.39)$$

where $\tau_i = t_{i+1} - t_i$, $t_0 = 0$, and where we compute $L_{\text{¥}}(t, t_i, t_{i+1})$ as (compare to (4.2))

¹⁸Reflecting the average bank credit rating, it is common to think of the Libor curve as a proxy for a AA-rated funding curve. In reality, however, this is not quite accurate, as banks with deteriorating credit are eliminated from the consortium of banks polled when determining the Libor rates. As such, the medium- and long-term forwards of the Libor curve contain *less* credit risk adjustment than similar forwards for a curve used to discount obligations to a single AA-rated firm. For more on this, see Collin-Dufresne and Goldstein [2001].

$$L_{\mathbb{Y}}(t, t_i, t_{i+1}) = \frac{1}{\tau_i} \left(P_{\mathbb{Y}}^{(L)}(t, t_i) / P_{\mathbb{Y}}^{(L)}(t, t_{i+1}) - 1 \right).$$

A similar construction can be done for any USD swap, by means of introducing curves $P_{\$}^{(L)}(t)$ and $P_{\$}(t)$. Technically speaking, the Libor forwards $L_{\mathbb{Y}}(t, t_i, t_{i+1})$ in (6.39) represent expectations in the t_{i+1} -forward measure — i.e. the martingale measure associated with the numeraire price $P_{\mathbb{Y}}(t, t_{i+1})$ (*not* $P_{\mathbb{Y}}^{(L)}(t, t_{i+1})$) — of quoted *spot* Libor rates,

$$L_{\mathbb{Y}}(t, t_i, t_{i+1}) \triangleq E_t^{t_{i+1}, \mathbb{Y}} (L_{\mathbb{Y}}(t_i, t_i, t_{i+1})). \quad (6.40)$$

In this view, the quoted Libor rate is effectively reduced to an observable index that may have little, if any, relationship to a true discount rate. For this reason, the time 0 pseudo-discount curves $P_{\$}^{(L)}(t)$ and $P_{\mathbb{Y}}^{(L)}(t)$ are often referred to as *index* curves.

It should be clear that the introduction of the pseudo-discount curves $P_{\$}^{(L)}(t)$ and $P_{\mathbb{Y}}^{(L)}(t)$ equips us with enough degrees of freedom to fit both USD-denominated swaps, JPY-denominated swaps, and the market for FX forward contracts. In fact, we have *too many* degrees of freedom: four curves, but only three separate markets to calibrate to. One way of handling this issue is to impose additional assumptions about the relationship between the curves $P(t)$ and $P^{(L)}(t)$ in one chosen currency. Before the 2007–2009 crisis, the following assumption was common.

Assumption 6.5.1. *In USD, the Libor pseudo-discount curve coincides with the real discount curve, i.e. $P_{\$}^{(L)}(t, T) = P_{\$}(t, T)$ for all t and T , $T \geq t$.*

Assumption 6.5.1 amounts to a convention where the liquidity and credit basis of non-USD Libor rates should all be measured relative to a neutral “bed-rock” established by USD Libor rates. Embedded into Assumption 6.5.1 also is the notion that most firms world-wide can fund themselves by borrowing in USD at levels close to USD-Libor; in the past this was often not a bad assumption. As we discuss in Section 6.5.3 below, post-crisis a non-trivial basis between index and discounting curves has emerged in the US. For simplicity of exposition we proceed in this section with Assumption 6.5.1, but the index-discounting basis in the US could be easily incorporated into the algorithm. The problem of accounting for this basis in *single-currency* curve construction is postponed until Section 6.5.3.

It is common to measure the difference between $P^{(L)}(t)$ and $P(t)$ in yield space, writing

$$P^{(L)}(t) = P(t)e^{-s(t)t},$$

where $s(t)$ is a yield spread often known as the *cross-currency (CRX) yield spread*. By Assumption 6.5.1, $s(t)$ is zero for USD, but will rarely be zero for any other currency. As indicated earlier, $s(t)$ will generally be quite small, often in the magnitude of a few basis points or less. Occasionally, however,

the CRX yield spread may blow out, particularly if banks in a particular country are perceived as having below-average credit quality. For instance, in the late 1990's, the CRX yield spread reached somewhere around -40 basis points in JPY as Japanese banks were perceived as being in economic trouble. During that period of time, foreign banks could generally fund themselves in USD at USD Libor, but in JPY at rates significantly below JPY Libor (due to their superior credit relative to Japanese banks). Had FX forward rates traded without any large CRX basis, foreign banks could have borrowed in JPY and used the FX forward markets to turn their obligations into USD-denominated ones at a borrowing cost below USD Libor, which would have indicated the existence of an inconsistency and an arbitrage. Conversely, in early 2008 the CRX basis spread became significantly positive (up to $+60$ basis points) as the hedging demands of long-dated FX books increased rapidly on the back of significant strengthening of the Yen versus the US Dollar. During the 2007–2009 crisis, many other currencies (including EUR) have experienced similar dramatic moves in the CRX basis spreads against USD.

6.5.2.3 Cross-Currency Basis Swaps

The discussion so far has assumed the existence of a liquid market for FX forwards, as means to observe and tie down the CRX basis between rates in two separate currencies. In reality, the interbank FX forward market is rarely liquid beyond maturities of one year, a far cry from the 30+ year horizons to which we often want to build yield curves. Rather than relying on FX forwards, instead we can turn to the market for *floating-floating cross-currency (CRX) basis swaps*. Briefly speaking, CRX basis swaps are contracts where floating Libor payments in one currency are exchanged for floating Libor payments in another currency, plus or minus a spread. The swaps involve an exchange of notional amounts at trade inception and at maturity; the ratio between the two notional amounts is normally set to equal the spot FX exchange rate prevailing at trade inception. CRX basis swaps are closely related to FX forward contracts — indeed a one-period CRX basis swap is identical to an FX forward contract.

As was the case with FX forward contracts, failing to fit to the market for CRX basis swaps can lead to arbitrable inconsistencies. For instance, consider the pricing of a stream of fixed USD cash flows. One way to determine the JPY price of these cash flows would be through simple discounting by the USD discount curve, followed by a conversion to JPY at the spot exchange rate. Alternatively, the following zero-cost scheme could be implemented to turn the stream of USD cash flows into fixed JPY cash flows:

1. Swap the fixed cash flows to streams of USD Libor plus some spread x , in a regular USD interest rate swap.

2. Swap USD Libor + x against JPY Libor + $e + x$ in a CRX basis swap, e being a market-quoted CRX basis swap spread.
3. Swap JPY Libor + $e + x$ against a stream of fixed JPY coupons, in a regular JPY interest rate swap.

The USD value of the cash flows in Step 3 can be determined through discounting with the JPY discount curve, and subsequent conversion to USD at the spot USD/JPY exchange rate. If the JPY discount curve is inconsistent with the basis-swap market, the value computed this way may not equal the value computed by discounting the original USD cash flows at the USD discount curve. Since the swap transactions 1–3 above are costless, this discrepancy will indicate an arbitrage.

We can use the pricing formalism developed in Section 6.5.2.2 to provide an explicit expression for the value of a CRX basis swap. For concreteness, we again turn to the USD/JPY market and consider a CRX basis swap, where a USD-based corporation receives USD Libor flat in exchange for payments of JPY Libor plus a fixed spread, $e_{\text{¥}}$. With payment dates $\{t_i\}_{i=1}^n$, the USD price $V_{\text{basisswap},\$}$ of the basis swap is (assuming a \$1 notional)

$$V_{\text{basisswap},\$}(0) \quad (6.41)$$

$$\begin{aligned} &= \sum_{i=0}^{n-1} L_{\$}(0, t_i, t_{i+1}) \tau_i P_{\$}(0, t_{i+1}) + P_{\$}(0, t_n) \\ &\quad - X(0) \left(\sum_{i=0}^{n-1} (L_{\text{¥}}(0, t_i, t_{i+1}) + e_{\text{¥}}) \tau_i P_{\text{¥}}(0, t_{i+1}) + P_{\text{¥}}(0, t_n) \right) \\ &= 1 - X(0) \\ &\quad \times \left(\sum_{i=0}^{n-1} \left(\frac{P_{\text{¥}}^{(L)}(0, t_i)}{P_{\text{¥}}^{(L)}(0, t_{i+1})} - 1 + e_{\text{¥}} \tau_i \right) P_{\text{¥}}(0, t_{i+1}) + P_{\text{¥}}(0, t_n) \right), \end{aligned} \quad (6.42)$$

where we have used the fact that $P_{\$}$ and $P_{\$}^{(L)}$ are identical (by Assumption 6.5.1), in order to reduce the time 0 price of the USD floating leg to \$1. The market quotes par values $e_{\text{¥}}^{\text{par}}$ — that is, the value of $e_{\text{¥}}$ that will make $V_{\text{basisswap},\$}(0) = 0$ — in a wide range of maturities extending out to 30 years or more.

6.5.2.4 Modified Curve Construction Algorithm

By Assumption 6.5.1, construction of the USD discount and Libor curves can be accomplished through direct application of the routines in Sections 6.2 or 6.3 on benchmark securities consisting of deposits, FRAs, and swaps. For non-USD currencies, however, matters are more complicated as we must now simultaneously estimate both curves $P(t)$ and $P^{(L)}(t)$, $t > 0$,

in a manner ensuring that i) Libor benchmark securities are correctly priced; and ii) par-valued cross-currency swaps against USD are correctly priced. In performing this exercise, we apply (6.42) and adjust valuation expressions for the benchmark securities according to the principles of the swap-pricing¹⁹ equation (6.39). We can make the curve construction problem quite complicated if we insist on $P(t)$ and $P^{(L)}(t)$ both being smooth functions of t ; instead, here we will show a simpler idea that applies earlier algorithms in this chapter in iterative fashion.

Working as before with JPY as the foreign currency, we start by assuming that we have somehow managed to construct the correct Libor curve $P_{¥}^{(L)}(t)$. Were we — erroneously — to pretend that $P_{¥}^{(L)}(t)$ were a proper discount curve, we would get, for our N benchmark securities, a vector of values $\mathbf{V}^{(L)}$ that would generally *not* equal the correct JPY market prices \mathbf{V} :

$$\mathbf{cP}_{¥}^{(L)} = \mathbf{V}^{(L)}, \quad \mathbf{V}^{(L)} \neq \mathbf{V}. \quad (6.43)$$

As the $P_{¥}^{(L)}(t)$ discount curve will be used to project forward rates, the yields and forward rates implied by $P_{¥}^{(L)}(t)$ should ideally be smooth. The smoothness requirements of the discount curve $P_{¥}(t)$, however, are significantly lower, as we shall never need to (in effect) differentiate this curve to produce forward Libor rates. Assuming that we have CRX basis swaps maturing on the benchmark set maturity dates $\{T_i\}_{i=1}^N$, it is thus, for instance, not unreasonable to write

$$P_{¥}(t) = P_{¥}(T_i) \frac{P_{¥}^{(L)}(t)e^{-\varepsilon_i \cdot (t-T_i)}}{P_{¥}^{(L)}(T_i)}, \quad t \in [T_i, T_{i+1}), \quad (6.44)$$

which assumes that the instantaneous forward rates generated by $P_{¥}(t)$ are given by those computed from $P_{¥}^{(L)}(t)$ plus a piecewise flat function:

$$f_{¥}(t) = f_{¥}^{(L)}(t) + \varepsilon(t), \quad \varepsilon(t) = \sum_{i=0}^{N-1} \varepsilon_i 1_{\{t \in [T_i, T_{i+1})\}}, \quad (6.45)$$

where $T_0 = 0$ as before.

In a cross-currency setting, the expression (6.4) no longer holds and must be replaced with something like

$$\mathbf{f}(\mathbf{P}_{¥}^{(L)}, \mathbf{P}_{¥}) = \mathbf{V} \quad (6.46)$$

for a non-linear vector-valued function \mathbf{f} . Indeed, according to (6.39), many of the coupon payments (\mathbf{c}) become a non-linear function of points on $\mathbf{P}_{¥}^{(L)}$ and cannot be considered constants. To salvage the methodologies discussed

¹⁹Pricing of short-term deposits only involves the discount curve P , whereas FRAs can be treated as one-period swaps. We leave details to the reader.

in Sections 6.2 or 6.3, we avoid working with (6.46) directly, and instead use an iteration based on equations (6.42), (6.43), and (6.45). The iteration attempts to estimate the unknown quantity $\mathbf{V}^{(L)}$ in (6.43) and works as follows:

1. Let $\mathbf{V}^{(L)}(j)$ be the j -th iteration for $\mathbf{V}^{(L)}$. Use $\mathbf{V}^{(L)}(j)$ along with (6.43) to estimate the curve $P_{\mathbb{Y}}^{(L)}(t)$, using any of the curve construction methods discussed in earlier sections of this chapter.
2. Given knowledge of $P_{\mathbb{Y}}^{(L)}(t)$, use (6.44)–(6.45) combined with (6.42) to imply the N constants $\varepsilon_0, \varepsilon_1, \dots, \varepsilon_{N-1}$, by calibration to the N par-valued CRX basis swaps maturing at time T_1, \dots, T_N . This calibration exercise can be done by simple bootstrapping, and establishes the current guess for $P_{\mathbb{Y}}(t)$.
3. Given estimates for both $P_{\mathbb{Y}}^{(L)}(t)$ and $P_{\mathbb{Y}}(t)$, compute guessed prices $\mathbf{V}(j)$ of all benchmark securities, i.e. evaluate the left-hand side of (6.46). If $\mathbf{V}(j)$ and \mathbf{V} are within a given tolerance, we are done.
4. Update the guess for $\mathbf{V}^{(L)}$ according to $\mathbf{V}^{(L)}(j+1) = \mathbf{V}^{(L)}(j) - (\mathbf{V}(j) - \mathbf{V})$, and proceed to Step 1.

The iteration is initiated at $j = 0$ with the estimate $\mathbf{V}^{(L)}(0) = \mathbf{V}$ and runs until the termination criterion in Step 3 is satisfied. As the approximation $\mathbf{V}^{(L)} \approx \mathbf{V}$ is normally very accurate, only a few iterations are needed to reach acceptable precision.

In this book we shall mostly ignore the existence of a non-zero CRX basis spread. In construction of a model for the evolution of the Libor curve, the reader should, however, keep in mind that it may be necessary to adjust the curve slightly before using it to discount any cash flows. In a dynamic setting, it is quite common to perform this adjustment by simply assuming that $\varepsilon(t)$ in (6.45) is deterministic. A discussion of how to incorporate both stochastic and deterministic spreads in a dynamic model for interest rate evolution can be found in Section 15.5. For now, we note that using deterministic spreads is generally safe, unless pricing securities with strong convexity in $\varepsilon(t)$ — e.g. an option on a CRX basis swap — in which case a separate stochastic model for $\varepsilon(t)$ may be needed.

6.5.3 Tenor Basis and Multi-Index Curve Group Construction

Section 6.5.2 relied extensively on the notion of separating the discount curves used for Libor projection and for outright discounting. This idea is quite powerful and has applications in other settings, including some where only a single currency is involved. For instance, for swaps that pay a non-Libor index — e.g. the Bond Mark Association (BMA) index in the US — it is natural to introduce a basis spread that measures the difference between forward rates of the non-Libor index curve and the Libor curve itself.

More recently, a similar technique has become important even for curves used for pricing standard Libor-based contracts. We have already mentioned (Section 5.1) that the Fed funds rate, the overnight rate used for balances of bank deposits with the Federal reserve, is often considered the closest proxy to the risk-free rate in the US (with Eonia and Sonia rates, see Section 5.1, fulfilling the same function for Euro and GBP). One argument for this choice is that most inter-dealer transactions are collateralized under the *International Swaps and Derivatives Association* (ISDA) Master Agreement, with the rate paid on collateral being the Fed funds rate (for USD; Eonia and Sonia for Euro and GBP), see Piterbarg [2010] (and also the discussion in Section 5.1). While the spread between the Fed funds rate and 3 month Libor rate used to be very small — in the order of a few basis points — after September 2007 it went up to as much as 275 basis points, and it is now generally accepted now that the Libor rate is no longer a good proxy for a discounting rate on collateralized trades. Uncollateralized derivative contracts are subject to credit risk, and a fully consistent pricing approach needs to incorporate the cost of hedging this risk (the co-called *credit valuation adjustment* or *CVA*). These computations are outside the scope of this book and can get very complex, in part because collateral rules can be complicated and are normally enforced on entire counterparty portfolios and not just on individual trades. See Gregory [2009] for further details.

As we discussed in Section 6.5.2, if we make an assumption on the index-discounting basis in one currency (say, USD), we can translate it into the index-discounting basis in any other currency through the market quotes for forward FX contracts and cross-currency basis swaps. However, to estimate this basis in USD (say), we need to rely on domestic markets only; doing otherwise will introduce a circularity into our arguments. Fortunately, the market in the appropriate instruments, the OIS (overnight index swap, a swap of payments based on a compounded Fed funds rate versus fixed rate, see (5.7)–(5.8)) and the Fed funds/Libor basis swaps (see Section 5.7) — has developed in the US with a range of maturities actively traded. Hence, using techniques that we already discussed, we can construct a pair of curves — a curve for discounting and a curve for projecting 3 month (say) forward Libor rates — in a self-consistent way from the market quotes on deposits, FRAs, swaps with 3 month frequency, and overnight index swaps.

Currently there are no countries where both the OIS market and the cross-currency basis swap (vs. USD) market are very liquid, and we can always use one or the other to find the index-discounting basis. As the markets evolve, there may come a time when there will be two liquid sources of discounting curve information. It turns out that potential conflict between the two can be resolved by carefully analyzing the collateral mechanisms used in the two markets and the implications for yield curve construction. This discussion is outside of scope of the current edition of our book, but the interested reader could consult Fujii et al. [2010] for details.

The challenges of curve construction do not end with building separate discount and forecasting curves to take into account the index-discounting basis. We also need to account for the *tenor basis* between vanilla single-currency swaps trading at different frequencies, e.g. 1 month, 3 months, and 6 months. Before proceeding, let us explain in more detail what we mean by tenor basis.

Suppose we construct Libor and discount curves based on, say, vanilla swaps (and for non-USD currencies also CRX basis swaps) paying 3 month Libor on a quarterly schedule. If the resulting index and discount curves are subsequently used to price a vanilla swap paying 1 month Libor on a monthly schedule, the resulting price is typically different from actual market quotes. In other words, there is a basis between the 3 month and 1 month Libor index curves, a basis arising partly from credit considerations and partly liquidity considerations (banks have a natural desire to have longer-term deposits to better match their loan commitments). Historically this basis has also been low; for example, the difference between 1 month and 3 month Libor rates was in the order of one basis point up until September 2007, but since then has been as wide as 50 basis points.

When various basis levels were small, the small discrepancies between different Libor-tenor swaps were often accounted for by building a unique discount curve for the subset of swaps referencing the Libor rate of a particular tenor; this curve would, in addition to generating the floating leg forward rates, then be used to discount floating and fixed cash flows of swaps of that frequency. In a swap pricing framework, this can create an arbitrage since it implies that fixed flows (from the fixed leg) will be discounted at different discount curves, depending on which Libor tenor the fixed flows happen to be paid against. Moreover, it is not clear how to deal with swaps that involve multiple Libor tenors²⁰, or how to aggregate risks coming from unrelated, individually constructed curves. Again, when the differences were small, these issues were largely ignored.

More recently, the naive approach above has evolved into the idea of using a *multi-index curve group*, a collection consisting of a single discount curve and multiple index curves, one for each Libor tenor covered by the multi-index curve group. The index curves are used in a tenor-specific manner to project Libor forward rates, and the universal discount curve serves to discount all floating and fixed cash flows, irrespective of tenor. The index curves are built sequentially as spreads off previously-built curves, which provides linkage between index curves and also a convenient risk parameterization. This relatively recent idea is discussed in Traven [2008] in considerable details, from where most of the material of this section is derived. Another good reference here is Fujii et al. [2010].

²⁰The most common example of this is a swap with a short front stub, i.e. a swap where at inception the first payment period is shorter than subsequent ones.

To discuss multi-index curve groups in detail, let us introduce a superscript k , $k = 1, \dots, K$, to distinguish quantities related to different tenors. In particular, let t_i^k be the i -th date in the tenor structure for tenor k ; $\tau_i^k = t_{i+1}^k - t_i^k$ the corresponding tenor offset; $L^k(t_i^k, t_i^k, t_{i+1}^k)$ the spot Libor rate of tenor k for the i -th period; and so on. If we denote the expected value of $L^k(t_i^k, t_i^k, t_{i+1}^k)$ under the t_{i+1}^k -forward measure by

$$L^k(t, t_i^k, t_{i+1}^k) \triangleq E_t^{t_{i+1}^k} (L^k(t_i^k, t_i^k, t_{i+1}^k))$$

(compare to (6.40)), then the value of a fixed-floating k -Libor-tenor swap with n periods and rate c is given by

$$V^k(0) = \underbrace{\sum_{i=0}^{n-1} c \tau_i^k P(t_{i+1}^k)}_{\text{Fixed Leg}} - \underbrace{\sum_{i=0}^{n-1} L^k(0, t_i^k, t_{i+1}^k) \tau_i^k P(t_{i+1}^k)}_{\text{Floating Leg}}. \quad (6.47)$$

Here $P(t)$ is the universal discounting curve. The time 0 index curve for tenor k , $P^k(t)$, is defined by the condition that forward Libor rates (of tenor k) be defined by the familiar formula

$$L(0, t_i^k, t_{i+1}^k) = \frac{1}{\tau_i^k} (P^k(t_i^k) / P^k(t_{i+1}^k) - 1). \quad (6.48)$$

A multi-index curve group is defined as a collection $\{P(\cdot), P^1(\cdot), \dots, P^K(\cdot)\}$ of the universal discounting curve and one index curve per tenor, with swaps priced via (6.47) (and equivalent formulas for other linear instruments) and (6.48).

Let us outline how to calibrate a multi-index curve group to market instruments referencing rates of different tenors. For each market, fixed-floating swaps referencing a Libor rate of one particular tenor L^k are usually the most liquid. Assume that this curve is the first index curve in the group, i.e. $k = 1$. The method from Section 6.5.2 can be used to construct a discounting curve, and a *base* index curve $P^1(\cdot)$ from

- Funding instruments such as deposits, forward FX contracts, OIS, overnight rate basis swaps (i.e. floating-floating swaps of an overnight rate versus L^1 , see Section 5.7 and a discussion of more general floating-floating single-currency basis swaps in the next paragraph), cross-currency basis swaps, and the like.
- Vanilla instruments referencing L^1 such as FRAs on L^1 and fixed-floating swaps on L^1 versus a fixed rate.

To construct $P^2(\cdot)$, we assume that prices of *floating-floating single-currency basis swaps* are available in the market. A floating-floating basis swap is a swap of payments linked to a Libor rate of a particular tenor — such as L^1 — versus payments based on a Libor rate of different tenor —

such as L^2 . Each leg pays on its own schedule of a corresponding tenor. A fixed spread is typically added to one of the legs to make the swap value at inception equal to zero. If a floating-floating basis swap is not traded or not liquid, it can be synthesized from two fixed-floating swaps referencing L^1 and L^2 .

If a floating-floating basis swap is traded at par in the market, the values of both legs should be the same at time 0:

$$\begin{aligned} & \sum_{i=0}^{n^2(T)-1} L^2(0, t_i^2, t_{i+1}^2) \tau_i^2 P(t_{i+1}^2) \\ &= \sum_{i=0}^{n^1(T)-1} (L^1(0, t_i^1, t_{i+1}^1) + e^{1,2}(T)) \tau_i^1 P(t_{i+1}^1). \quad (6.49) \end{aligned}$$

Here $n^k(T)$ is the number of periods in the tenor structure to date T for tenor k , and $e^{1,2}(T)$ is the quoted floating-floating basis spread for exchanging L^1 for L^2 to maturity T , quoted on the L^1 leg. It could be positive or negative, depending on perceived desirability of payments linked to L^1 versus L^2 .

Similar to (6.44)–(6.45), we represent $P^2(\cdot)$ as a multiplicative spread to $P^1(\cdot)$:

$$P^2(t) = P^1(t) e^{-\int_0^t \eta^{1,2}(s) ds}, \quad t \geq 0$$

for a given, usually piecewise-constant, spread function $\eta^{1,2}(\cdot)$. With the discounting curve $P(\cdot)$ already constructed and $L^1(0, t_i^1, t_{i+1}^1)$ known for all i from the already-built index curve $P^1(\cdot)$, it is a simple exercise to obtain the spread function $\eta^{1,2}(\cdot)$ by solving (6.49) for different T 's.

Having built $P^2(\cdot)$, the remaining index curves $P^k(\cdot)$, $3 \leq k \leq K$ may be constructed in a similar fashion, always using floating-floating basis swap spreads for L^k versus L^1 basis swaps or, more generally, for whatever L^k versus L^i basis swaps are the most liquid with $i < k$. In particular, each index curve $P^k(\cdot)$ for $k > 1$ is built as a *spread*, or *basis*, curve to one of the previous curves.

In the presence of multiple curves, it is not entirely clear from the outset how to most sensibly define risk sensitivities. Fortunately, with this spread-based method of curve group construction, sensitivities to instruments used in the curve group have clear, and orthogonal, meaning:

- Perturbations to instruments used in building the base index curve, e.g. non-basis swaps and FRAs referencing L^1 , define risk sensitivities to the overall levels of interest rates. Clearly, with basis spreads for L^1 -versus- L^k floating-floating basis swaps, $k = 2, \dots, K$, kept constant, shifts to fixed-floating L^1 -versus-fixed swap rates will move all index curves together by the same amount in forward rate space. These sensitivities are the direct analog of the standard interest rate deltas in the traditional, single-curve, world, see Section 6.4.

- Perturbations to funding instruments define sensitivities to discounting.
- Perturbations to basis swap spreads for L^k -versus- L^1 floating-floating basis swaps define *basis risk*, i.e. the risk that index curves of different tenors do not move in lock step.

The parameterization allows us to naturally aggregate “similar” risks such as overall rate level risks, discounting risks, basis risks, while keeping different kinds separate for efficient risk management. Had we constructed all index curves in separation from each other (from multiple sets of vanilla fixed-floating swaps, say) such automatic aggregation would not be possible.

6.A Appendix: Spline Theory

6.A.1 Hermite Spline Theory

Consider a given set of data points (x_i, f_i) , $i = 1, \dots, N$, where $x_1 < x_2 < \dots < x_N$. We wish to apply an interpolation rule such that a continuous function $f(x)$, $x \in [x_1, x_N]$, is created. We require that f be piecewise cubic, be at least once differentiable (C^1), and be a true interpolating function, i.e. $f(x_i) = f_i$ for all $i = 1, \dots, N$.

In the Hermite spline description, tangents at points x_i , $i = 1, \dots, N$, are assumed exogenously specified. Let f'_i denote the tangent df/dx at $x = x_i$, $i = 1, \dots, N$. We write f as a piecewise cubic polynomial

$$f(x) = a_{3,i}(x - x_i)^3 + a_{2,i}(x - x_i)^2 + a_{1,i}(x - x_i) + a_{0,i}, \quad x \in [x_i, x_{i+1}],$$

with unknown coefficients $a_{j,i}$ specific to each interval $[x_i, x_{i+1}]$. Expressing that both f and f' should be continuous across each point x_i allows us, after a little rearrangement, to write the spline specification as simply

$$f(x) = \mathbf{D}_i(x)^\top \mathbf{M} \begin{pmatrix} f_i \\ f_{i+1} \\ f'_i h_i \\ f'_{i+1} h_i \end{pmatrix}, \quad x \in [x_i, x_{i+1}], \quad (6.50)$$

where $h_i \triangleq x_{i+1} - x_i$,

$$\mathbf{D}_i(x) = \begin{pmatrix} \delta_i^3 \\ \delta_i^2 \\ \delta_i \\ 1 \end{pmatrix}, \quad \delta_i \triangleq \frac{x - x_i}{h_i},$$

and \mathbf{M} is the *Hermite matrix*

$$\mathbf{M} = \begin{pmatrix} 2 & -2 & 1 & 1 \\ -3 & 3 & -2 & -1 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}.$$

One drawback of the Hermite specification is the need to directly specify tangents df/dx . A number of approaches exist to compute these directly from the given data points or by adding additional control points. For our purposes, we highlight the so-called *Catmull-Rom spline* (Catmull and Rom [1974]), where the derivatives are computed as²¹

$$f'_i = \frac{f_{i+1} - f_{i-1}}{x_{i+1} - x_{i-1}}, \quad i = 2, \dots, N-1. \quad (6.51)$$

At end points (x_1, f_1) and (x_N, f_N) forward and backward differences are used instead:

$$f'_1 = \frac{f_2 - f_1}{x_2 - x_1}; \quad f'_N = \frac{f_N - f_{N-1}}{x_N - x_{N-1}}. \quad (6.52)$$

Notice that with (6.51), the Hermite representation (6.50) can be rewritten in the derivative-free form

$$f(x) = \mathbf{D}_i(x)^\top \mathbf{A}_i \begin{pmatrix} f_{i-1} \\ f_i \\ f_{i+1} \\ f_{i+2} \end{pmatrix}, \quad x \in [x_i, x_{i+1}], \quad (6.53)$$

where

$$\mathbf{A}_i = \begin{pmatrix} -\alpha_i & 2 - \beta_i & -2 + \alpha_i & \beta_i \\ 2\alpha_i & \beta_i - 3 & 3 - 2\alpha_i & -\beta_i \\ -\alpha_i & 0 & \alpha_i & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}, \quad i = 2, \dots, N-2, \quad (6.54)$$

with

$$\alpha_i = \frac{h_i}{h_i + h_{i-1}}, \quad \beta_i = \frac{h_i}{h_{i+1} + h_i}.$$

As indicated, equation (6.53) only holds for $i = 2, \dots, N-2$, with external boundary conditions needed to establish the curve in the segments $[x_1, x_2]$ and $[x_{N-1}, x_N]$. Applying boundary conditions of the type (6.52), we get

$$\mathbf{A}_1 = \begin{pmatrix} 0 & 1 - \beta_1 & -1 & \beta_1 \\ 0 & -1 + \beta_1 & 1 & -\beta_1 \\ 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}, \quad (6.55)$$

and

²¹Variations exist which use more elaborate finite difference style derivatives, taking into account the fact that the grid may be non-equidistant; see Chapter 2. Given the semi-heuristic nature of the Catmull-Rom spline, it is doubtful that much is gained by such extensions.

$$\mathbf{A}_{N-1} = \begin{pmatrix} -\alpha_{N-1} & 1 & -1 + \alpha_{N-1} & 0 \\ 2\alpha_{N-1} & -2 & 2 - 2\alpha_{N-1} & 0 \\ -\alpha_{N-1} & 0 & \alpha_{N-1} & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}. \quad (6.56)$$

While Catmull-Rom splines shall suffice for the yield curve applications we have in mind here, let us note that it is possible to introduce further parameters to control local spline behavior. For completeness, let us quickly list one popular extension. First, we allow for the possibility that the curve is locally only C^0 by having incoming and outgoing tangents be different. That is, we define

$$f'_{i,I} = \lim_{\varepsilon \downarrow 0} \frac{f(x_i) - f(x_i - \varepsilon)}{\varepsilon}; \quad f'_{i,O} = \lim_{\varepsilon \downarrow 0} \frac{f(x_i + \varepsilon) - f(x_i)}{\varepsilon},$$

and rewrite the Hermite equation (6.50) as

$$f(x) = \mathbf{D}_i(x)^\top \mathbf{M} \begin{pmatrix} f_i \\ f_{i+1} \\ f'_{i,O} h_i \\ f'_{i+1,I} h_i \end{pmatrix}, \quad x \in [x_i, x_{i+1}]. \quad (6.57)$$

Only when $f'_{i,I} = f'_{i,O}$ for all i is the curve C^1 everywhere. The *Kochanek-Bartels spline* — also known as the *TCB spline* — is defined through the expressions

$$\begin{aligned} f'_{i,I} &= \frac{(1 - \sigma_i)(1 + c_i)(1 - b_i)}{2} \frac{f_{i+1} - f_i}{x_{i+1} - x_i} \\ &\quad + \frac{(1 - \sigma_i)(1 - c_i)(1 + b_i)}{2} \frac{f_i - f_{i-1}}{x_i - x_{i-1}}, \end{aligned} \quad (6.58)$$

$$\begin{aligned} f'_{i,O} &= \frac{(1 - \sigma_i)(1 - c_i)(1 - b_i)}{2} \frac{f_{i+1} - f_i}{x_{i+1} - x_i} \\ &\quad + \frac{(1 - \sigma_i)(1 + c_i)(1 + b_i)}{2} \frac{f_i - f_{i-1}}{x_i - x_{i-1}}, \end{aligned} \quad (6.59)$$

for parameters $\sigma_i, c_i, b_i \in [-1, 1]$, $i = 1, \dots, N$. The parameters σ_i , c_i , and b_i are used to control curve *tension*, *continuity*, and *bias*, respectively; clearly, when $\sigma_i = c_i = b_i = 0$, the Kochanek-Bartels spline reduces to the Catmull-Rom spline. A positive value of σ_i will tend to “tighten” the curve around the point (x_i, f_i) , and a negative value will generate “slack”. The parameters b_i are measures of over- and undershoot. To see this, set $\sigma_i = c_i = 0$ and note that when $b_i = 0$, left and right one-sided tangents are weighted equally producing a regular Catmull-Rom spline; when b_i is close to -1 (1), however, the outgoing (incoming) tangent dominates the path of the curve through the point (x_i, f_i) producing undershoot (overshoot). The parameters c_i control the degree of differentiability of the resulting spline: if a parameter $c_i \neq 0$,

the resulting spline will develop a corner (the direction of which depends on the sign of c_i) at point (x_i, f_i) , losing its differentiability. Kochanek-Bartels splines are used extensively in computer graphics applications.

6.A.2 C^2 Cubic Splines

The cubic splines in Section 6.A.1 are generally not twice differentiable, and their second derivatives will jump across each knot. We wish to remedy this, and now consider a twice differentiable C^2 cubic spline $f(x)$ interpolating a set of data points (x_i, f_i) , $i = 1, \dots, N$. By necessity, such a cubic spline interpolant is piecewise linear in its second derivative

$$f''(x) = \frac{x_{i+1} - x}{x_{i+1} - x_i} f_i'' + \frac{x - x_i}{x_{i+1} - x_i} f_{i+1}'', \quad x \in [x_i, x_{i+1}], \quad (6.60)$$

where we use primes to denote differentiation and where $f_i'' \triangleq f''(x_i)$. We emphasize that for a C^2 cubic spline, the second derivative is continuous across knot points: $\lim_{x \downarrow x_i} f''(x) = \lim_{x \uparrow x_i} f''(x) = f''(x_i)$. Integrating (6.60) twice and requiring the curve to pass through data points results in the classical spline equation

$$\begin{aligned} f(x) &= \frac{(x_{i+1} - x)^3}{6h_i} f_i'' + \frac{(x - x_i)^3}{6h_i} f_{i+1}'' + (x_{i+1} - x) \left(\frac{f_i}{h_i} - \frac{h_i}{6} f_i'' \right) \\ &\quad + (x - x_i) \left(\frac{f_{i+1}}{h_i} - \frac{h_i}{6} f_{i+1}'' \right), \quad x \in [x_i, x_{i+1}], \end{aligned} \quad (6.61)$$

where $h_i \triangleq x_{i+1} - x_i$. The second derivatives f_i'' , $i = 1, \dots, N$, needed to evaluate (6.61) can be obtained by requiring $f'(x)$ to be continuous across data points. The result is

$$\frac{h_{i-1}}{6} f_{i-1}'' + \frac{h_{i-1} + h_i}{3} f_i'' + \frac{h_i}{6} f_{i+1}'' = \frac{f_{i+1} - f_i}{h_i} - \frac{f_i - f_{i-1}}{h_{i-1}}, \quad i = 2, \dots, N-1. \quad (6.62)$$

The set of equations (6.62) is a tri-diagonal system for f_i'' that can be solved in $O(N - 2)$ operations once we have specified boundary conditions²² for f_1'' and f_N'' . A classical boundary condition is $f_1'' = f_N'' = 0$, leading to the so-called *natural cubic spline*.

While C^2 cubic splines have a number of useful features, they have, loosely speaking, a built-in aversion to make tight turns (since this will cause large values of f''). This, in turn, will often produce extraneous inflection points and non-local behavior, in the sense that perturbation of a single f_i will significantly affect the appearance of the curve for x -values far from x_i . Also, monotonicity and convexity properties of the original data-set will typically not be preserved.

²²Such boundary conditions may be indirect and can, among other choices, take the form of specification of a gradient f' at x_1 or x_N . By differentiation of (6.61), a gradient specification can always be turned into a condition on f'' at x_1 or x_N .

6.A.3 C^2 Exponential Tension Splines

An attractive remedy to the shortcomings of the cubic spline is to insert some *tension* in the cubic spline, that is, to apply a tensile force to the end-points of the spline. Formally, this can be accomplished (see Schweikert [1966]) by replacing the equation (6.60) with, $x \in [x_i, x_{i+1}]$,

$$f''(x) - \sigma^2 f(x) = \frac{x_{i+1} - x}{x_{i+1} - x_i} (f''_i - \sigma^2 f_i) + \frac{x - x_i}{x_{i+1} - x_i} (f''_{i+1} - \sigma^2 f_{i+1}), \quad (6.63)$$

where $\sigma > 0$ is a measure of the tension applied to the cubic spline²³. Notice that we have replaced the assumption of a piecewise linear second derivative with the assumption that the quantity $f''(x) - \sigma^2 f(x)$ is linear on each sub-interval $[x_i, x_{i+1}]$.

Integrating (6.63) twice and requiring that the curve pass through the given data points, one obtains (after some rearrangements)

$$\begin{aligned} f(x) = & \left(\frac{\sinh(\sigma(x_{i+1} - x))}{\sinh(\sigma h_i)} - \frac{x_{i+1} - x}{h_i} \right) \frac{f''_i}{\sigma^2} \\ & + \left(\frac{\sinh(\sigma(x - x_i))}{\sinh(\sigma h_i)} - \frac{x - x_i}{h_i} \right) \frac{f''_{i+1}}{\sigma^2} \\ & + f_i \frac{x_{i+1} - x}{h_i} + f_{i+1} \frac{x - x_i}{h_i}, \quad x \in [x_i, x_{i+1}], \end{aligned} \quad (6.64)$$

where $h_i = x_{i+1} - x_i$ as before. Requiring continuity of the first derivative we then get for the f''_i ,

$$\begin{aligned} & \left(\frac{1}{h_{i-1}} - \frac{\sigma}{\sinh(\sigma h_{i-1})} \right) \frac{f''_{i-1}}{\sigma^2} \\ & + \left(\frac{\sigma \cosh(\sigma h_{i-1})}{\sinh(\sigma h_{i-1})} - \frac{1}{h_{i-1}} + \frac{\sigma \cosh(\sigma h_i)}{\sinh(\sigma h_i)} - \frac{1}{h_i} \right) \frac{f''_i}{\sigma^2} \\ & + \left(\frac{1}{h_i} - \frac{\sigma}{\sinh(\sigma h_i)} \right) \frac{f''_{i+1}}{\sigma^2} = \frac{f_{i+1} - f_i}{h_i} - \frac{f_i - f_{i-1}}{h_{i-1}}. \end{aligned}$$

Again, this is a tri-diagonal system of equations that can be solved in $O(N - 2)$ operations once we specify f''_1 and f''_N .

From the representation (6.64), it is clear that on all intervals $[x_i, x_{i+1}]$ hyperbolic tension splines can be written as linear combinations of the basis functions 1, x , $e^{-\sigma x}$, $e^{\sigma x}$. The representation (6.64), however, has better behavior for large and small values of σ (see Renka [1987] and Rentrop [1980] for details about proper evaluation of the hyperbolic functions in (6.64) for large and small values of σ).

²³Extension to non-uniform tension parameter is straightforward and involves replacing σ with σ_i in (6.63), with σ_i then being a measure of the tension applied locally to the curve in the interval $[x_i, x_{i+1}]$.

We notice that when the tension parameter $\sigma = 0$, equations (6.63) and (6.60) are identical, i.e. the tension spline degenerates into a regular cubic spline. On the other hand, when $\sigma \gg 1$ (6.63) reduces to piecewise linear interpolation, as

$$\lim_{\sigma \rightarrow \infty} f(x) = \frac{x_{i+1} - x}{x_{i+1} - x_i} f_i + \frac{x - x_i}{x_{i+1} - x_i} f_{i+1}, \quad x \in [x_i, x_{i+1}]. \quad (6.65)$$

Evidently, the equation (6.63) defines a twice differentiable curve that is a hybrid between a cubic spline and a piecewise linear spline.

The convergence of the tension spline towards a piecewise linear curve as $\sigma \rightarrow \infty$ can be shown to be *uniform*, i.e. (6.65) holds uniformly in $[x_i, x_{i+1}]$ for $i = 1, \dots, N - 1$. Similarly

$$\lim_{\sigma \rightarrow \infty} f'(x) = \frac{f_{i+1} - f_i}{x_{i+1} - x_i} \quad \text{and} \quad \lim_{\sigma \rightarrow \infty} f''(x) = 0$$

uniformly in any closed subinterval of $[x_i, x_{i+1}]$. See Pruess [1976] for details and a proof. The uniform convergence is important as it guarantees that the monotonicity and convexity properties of the underlying discrete data set are preserved, simply by choosing a sufficiently high value of the tension factor. Due to this property, hyperbolic tension splines are said to be *shape-preserving*²⁴. As the tension factor increases, the resulting spline will also behave in increasingly local fashion towards input perturbations. In the limit $\sigma \rightarrow \infty$ each point $f(x)$ on the spline will only be associated with the two nearest-neighbor knots.

²⁴Generalizing, suppose we introduce constraints on function values, first derivatives, or second derivatives. As long as these constraints are satisfied by piecewise linear interpolation, there will exist some value of the tension parameter σ (possibly $\sigma = 0$) which will make the tension spline satisfy the constraints. This observation is key to algorithms for automatic selection of σ from externally specified function constraints. See, for instance, Lynch [1982] and Renka [1987] for details and efficient algorithms for automatic tension selection.

Vanilla Models with Local Volatility

We have shown in Section 5.10 that European swaptions (and caplets/floorlets which we equate with one-period swaptions) can be valued as European options on forward swap rates. As a consequence, a full term structure model that specifies the dynamics of the whole yield curve through time is essentially unnecessary for European swaption valuation. Instead, we only need a model for the evolution — in fact, just a terminal distribution — of a single swap rate in isolation. Models of this type shall be denoted *vanilla models*, to distinguish them from full term structure models. Vanilla models can be extended by copula methods to describe joint terminal distributions of more than one rate, as we discuss in Chapter 17. Ultimately, however, their primary purpose in this book is to serve as a foundation for development of more widely applicable full term structure models, that is, models which provide consistent dynamics for *all* points on the yield curve simultaneously. Term structure models are extensively covered later in the book.

In this chapter, we review one-factor diffusive models where our ability to alter the terminal distribution stems from a single source: a swap rate dependent diffusion function. Models of this type are often known as *deterministic volatility function* (DVF), or sometimes *local volatility function* (LVF), models. We first discuss the most common tractable specifications of such models — the CEV, displaced diffusion, and quadratic models — and then move on to efficient numerical or expansion-based methods for European option pricing within the general DVF model class. The listed techniques and results will be frequently referenced in later chapters, especially in the context of model calibration.

7.1 General Framework

7.1.1 Model Dynamics

Let $S(t)$ denote a forward Libor or swap rate, and let $W(t)$ be a one-dimensional Brownian motion under a measure P in which $S(\cdot)$ is a martingale. We assume that $S(t)$ follows the one-dimensional SDE

$$dS(t) = \lambda\varphi(S(t)) dW(t), \quad (7.1)$$

where λ is a positive constant¹ and $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ satisfies regularity conditions, such as those in Theorem 1.6.1. In most applications we would ideally want $S(t)$ to be non-negative, which is easily seen to impose the restriction

$$\varphi(0) = 0. \quad (7.2)$$

In some cases we may consciously decide to violate (7.2) for the sake of model tractability.

When dealing with vanilla models, we primarily work in the measure P , so we typically abbreviate E^P to E when there is no possibility of confusion.

7.1.2 Volatility Smile and Implied Density

The role of the function φ is to match the distribution of S to that observed through puts and calls traded in the market. Specifically, let $c(t, S; T, K)$ denote the (non-deflated) time t value of a T -maturity European call option struck at K with $S(t) = S$, i.e.

$$c(t, S(t); T, K) = E_t \left((S(T) - K)^+ \right). \quad (7.3)$$

The time t probability density of $S(T)$ can be derived from time t observed values of c , as (proceeding heuristically)

$$P_t(S(T) \in dK) / dK = E_t(\delta(S(T) - K)) \quad (7.4)$$

$$= E_t \left(\frac{\partial^2 c(T, S(T); T, K)}{\partial K^2} \right) = \frac{\partial^2 c(t, S(t); T, K)}{\partial K^2}, \quad (7.5)$$

where δ is the Dirac delta function. This classical result is due to Breeden and Litzenberger [1978] and allows us to construct the marginal density of $S(T)$ from prices of T -maturity call options for a continuum of strikes K .

In option markets, it is common to express the strike dependency of call (and put) options in terms of the so-called *implied volatilities*. Specifically,

¹We allow for time dependence later in the chapter, starting in Section 7.6.

for a given option price c at strike K and maturity T , we define the time t implied volatility function $\sigma_B(t, S; T, K)$ as the solution to

$$\begin{aligned} c(t, S; T, K) &= S\Phi(d_+) - K\Phi(d_-), \\ d_{\pm} &= \frac{\ln(S/K) \pm \frac{1}{2}\sigma_B(t, S; T, K)^2(T-t)}{\sigma_B(t, S; T, K)\sqrt{T-t}}. \end{aligned} \quad (7.6)$$

We recognize the right-hand side of (7.6) as the Black-Scholes-Merton formula for a martingale process, i.e. the Black model (see Remark 1.9.4), with constant volatility $\sigma_B(t, S; T, K)$. The mapping $K \mapsto \sigma_B(t, S; T, K)$ is known as the T -maturity *volatility smile*². In most established fixed income markets, the volatility smile is predominantly downward-sloping³ in K , although it is not uncommon for σ_B to eventually increase in K for sufficiently large values of K .

7.1.3 Choice of φ

If we allowed φ to depend on time, results by Dupire [1994] and Andersen and Brotherton-Ratcliffe [1998] demonstrate that any arbitrage-free marginal distribution of $S(T)$ can be realized by suitable choice of $\varphi = \varphi(t, S)$, $t \in [0, T]$. Indeed, non-parametric expressions exist to uniquely imply $\varphi(t, K)$ from observations of $\sigma_B(0, S(0); t, K)$ for the double continuum $(t, K) \in [0, T] \times [0, \infty)$. Unless the resulting φ happens to be monotonically increasing or decreasing in S , however, the resulting model will imply non-stationary volatility smile behavior, which is contrary to typical behavior of actual markets. To expand on this issue, consider setting

$$\varphi(S) = a + (S - S(0))^2, \quad (7.7)$$

where $a > 0$. The function $\varphi(S)$ is thus a *U*-shaped function with a minimum value of a at $S = S(0)$. Using formulas from Section 7.3 below, it can be verified (and is intuitively obvious) that the time 0 volatility smile σ_B produced by this parameterization is also *U*-shaped. Moving forward to time $t > 0$, consider the smile generated at t by (7.7) if $S(t) \gg S(0)$. At a large level of $S(t)$, $\varphi(S)$ will appear to be a strongly increasing function of S , causing (7.7) to produce a volatility smile no longer *U*-shaped, but instead monotonically increasing at all statistically relevant strikes. Conversely, if $S(t)$ diffuses below $S(0)$ such that $S(t) \ll S(0)$, a monotonically *decreasing* smile will arise at time t .

²In case the smile is monotonically downward or upward sloping, i.e. not *U*-shaped, it is often called a *volatility skew*. *Skew* then refers to the slope of the smile.

³This is not necessarily true for emerging markets where the volatility smile, when observed, can be significantly upward sloping or convex.

Strong level-dependence of the basic volatility smile shape is often at odds with observable market behavior, and non-monotonic specifications of $\varphi(S)$ — such as (7.7) — should consequently be approached with some care. As a consequence, the basic model (7.1) is most appropriate for markets where the volatility smile is (close to) a monotonic function of K . A classical monotonic choice for φ is the *constant elasticity of variance* (CEV) specification

$$\varphi(S) = S^p, \quad (7.8)$$

for some constant p . As we proceed to show, this specification is analytically tractable.

7.2 CEV Model

7.2.1 Basic Properties

In this section, we examine the CEV specification (7.8) in detail. We start out with the following proposition:

Proposition 7.2.1. *Consider the stochastic differential equation*

$$dS(t) = \lambda S(t)^p dW(t), \quad (7.9)$$

where $p > 0$ is constant and $W(t)$ is a one-dimensional Brownian motion. The following holds:

1. All solutions to (7.9) are non-explosive.
2. For $p \geq 1/2$, the SDE (7.9) has a unique solution.
3. For $0 < p < 1$, $S = 0$ is an attainable boundary for (7.9); for $p \geq 1$, $S = 0$ is an unattainable boundary for (7.9).
4. For $0 < p \leq 1$, $S(t)$ in (7.9) is a martingale; for $p > 1$, $S(t)$ is a strict supermartingale.

Proof. Property 1 follows from a remark on page 332 and equation (5.5.19) in Karatzas and Shreve [1997], and Property 2 follows from Example 5.2.14 in Karatzas and Shreve [1997]. Property 3 can be proven using the classical Feller boundary classification techniques based on speed/scale measure integral, see Section 5.5 of Karatzas and Shreve [1997]; Andersen and Andreasen [2000b] have the details. Property 4 is proven in Sin [1998]. \square More details on boundary characterization for CEV processes can be found in Davydov and Linetsky [2001].

Remark 7.2.2. For $p \geq 1/2$, the solution to (7.9) is unique. Hence, if the solution ever reaches the origin ($S = 0$), it stays there, i.e. is *absorbed*. For $0 < p < 1/2$, however, there are solutions that stay at origin if they reach it, and there are solutions that jump out if it. Hence, to define a unique

solution, a boundary condition at $S = 0$ must be specified for (7.9). In practice, we set $S = 0$ to be an *absorbing* barrier: if $S(t)$ hits 0 for the first time at $t = \tau$, $S(u) = 0$ for all $u \geq \tau$. This condition is not only imposed to be consistent with the case of $p \geq 1/2$, but is also the only boundary condition consistent with the absence of arbitrage.

Remark 7.2.3. While it is common to require the parameter p to be positive, the process is well-defined for negative p , $p < 0$, as well, with the same absorbing boundary condition at $S = 0$ as for the case $0 < p < 1/2$ above. This enlargement of the domain of applicability of the process is occasionally useful in the fixed income markets, although much less so than in equity or FX markets where the smiles can generally be much more downward sloping.

For $p < 1$ and $t > 0$, the time 0 probability that $S(t) = 0$ is non-zero. In fact, it can be shown (see, for example, Cox [1996]) that if τ , the first time $S(\cdot)$ hits 0, is greater than t , then

$$P_t(\tau < T | \tau > t) = G\left(|\vartheta|, \frac{X(t)}{2\lambda^2(T-t)}\right), \quad T > t,$$

where

$$\vartheta = \frac{1}{2(p-1)}, \quad (7.10)$$

$$X(t) = \frac{S(t)^{2(1-p)}}{(1-p)^2}, \quad (7.11)$$

and G is the *complementary Gamma function*

$$G(a, x) \triangleq \frac{\Gamma(a, x)}{\Gamma(a)},$$

with the incomplete Gamma function $\Gamma(a, x)$ given by

$$\Gamma(a, x) = \int_x^\infty u^{a-1} e^{-u} du, \quad \Gamma(a) = \Gamma(a, 0). \quad (7.12)$$

If the absorption probability is substantial, one may want to consider regularizing the process to prevent absorption; see Section 7.2.3 for this.

Due to the result in Proposition 7.2.1, Property 4, we normally prefer to avoid using $p > 1$. As $p > 1$ will produce volatility smiles increasing in K (and thereby different from those in fixed income markets), this restriction on p is often of little practical concern.

The transition density of $S(\cdot)$ in (7.9) is known in closed form and is listed below for reference, along with a short proof that highlights the relationship between CEV processes and squared Bessel processes.

Lemma 7.2.4. Consider the SDE (7.9) for any $p \neq 1$ (including $p < 0$ and $p > 1$), and let ϑ and $X(t)$ be as in (7.10)–(7.11). Let $q(X(T)|X(t))$ be the conditional P-density of $X(T)$ given $X(t) > 0$, $t < T$. If the level $S = 0$ is defined to be an absorbing boundary for (7.9) when $p \leq 1/2$, then

$$\begin{aligned} q(X(T)|X(t)) &= \frac{1}{2\lambda^2(T-t)} \exp\left(-\frac{X(T)+X(t)}{2\lambda^2(T-t)}\right) \\ &\quad \times \left(\frac{X(t)}{X(T)}\right)^{-\vartheta/2} I_{|\vartheta|}\left(\frac{\sqrt{X(T)X(t)}}{\lambda^2(T-t)}\right), \end{aligned}$$

where $I_a(x)$ is the modified Bessel function of the first kind of order a :

$$I_a(x) = \sum_{j=0}^{\infty} \frac{(x/2)^{a+2j}}{j!\Gamma(a+j+1)}.$$

Proof. According to Ito's lemma, the process $X(t)$ satisfies the SDE

$$dX(t) = \lambda^2 \frac{1-2p}{1-p} dt + 2\lambda \sqrt{X(t)} dW(t).$$

Define the process $Y(v)$ by $Y(v) = X(v/\lambda^2)$. Applying a time change, it follows that

$$dY(v) = \frac{1-2p}{1-p} dv + 2\sqrt{Y(v)} d\widetilde{W}(v),$$

where $\widetilde{W}(\cdot)$ is a Brownian motion, up to the absorption time $\inf\{v > 0 : Y(v) = 0\}$. The process for Y can be identified as a so-called *squared Bessel process of index ϑ* . Standard results for this process (see e.g. p. 117 of Borodin and Salminen [1996]) give the result in the lemma. \square

Remark 7.2.5. By the usual transformation rules for densities, the density for $S(T)$ conditional on $S(t)$ is

$$q(X(T)|X(t)) \cdot 2S(T)^{2(1-p)-1}/|1-p|.$$

7.2.2 Call Option Pricing

Consider now the valuation of European call options in the CEV model, requiring evaluation of the expectation

$$c_{\text{CEV}}(t, S(t); T, K) \triangleq \mathbb{E}_t \left((S(T) - K)^+ \right)$$

for $S(\cdot)$ that follows (7.9). Using the definition (7.11), we can rewrite this as

$$\begin{aligned} c_{\text{CEV}}(t, S(t); T, K) &= \mathbb{E}_t \left(\left([(1-p)^2 X(T)]^{-\vartheta} - K \right)^+ \right) \\ &= \int_0^\infty \left([(1-p)^2 X]^{-\vartheta} - K \right)^+ q(x|X(t)) dx, \end{aligned}$$

where we have assumed $p \neq 1$ and the density $q(x|X(t))$ is given in Lemma 7.2.4. A straightforward, but tedious, integration exercise (see e.g. Schroder [1989] or Andersen and Andreasen [2000b]) yields the following result:

Proposition 7.2.6. Consider the CEV model (7.9). Let $\chi_\nu^2(\gamma)$ be a non-central chi-square distributed variable with ν degrees of freedom and non-centrality parameter γ , and let $\Upsilon(x, \nu, \gamma) = P(\chi_\nu^2(\gamma) \leq x)$ be the cumulative distribution function for $\chi_\nu^2(\gamma)$. Also define

$$a = \frac{K^{2(1-p)}}{(1-p)^2 \lambda^2(T-t)}, \quad b = |p-1|^{-1}, \quad c = \frac{S^{2(1-p)}}{(1-p)^2 \lambda^2(T-t)}.$$

Then, for $0 < p < 1$ and an absorbing boundary at $S = 0$ we have, for $K > 0$,

$$c_{\text{CEV}}(t, S; T, K) = S(1 - \Upsilon(a, b+2, c)) - K\Upsilon(c, b, a). \quad (7.13)$$

Remark 7.2.7. The result above in fact holds for all $p < 1$, including negative p . A complimentary result holds for $p > 1$,

$$c_{\text{CEV}}(t, S; T, K) = S(1 - \Upsilon(c, b, a)) - K\Upsilon(a, b+2, c). \quad (7.14)$$

Remark 7.2.8. The special case $p = 1$ leads to the Black pricing formula with volatility λ , see (1.43) and Remark 1.9.4), so that

$$c_B(t, S; T, K; \lambda) = S\Phi(d_+) - K\Phi(d_-), \quad (7.15)$$

where

$$d_\pm = \frac{\ln(S/K) \pm \lambda^2(T-t)/2}{\lambda\sqrt{T-t}},$$

and $\Phi(\cdot)$ is the standard Gaussian CDF.

Remark 7.2.9. For the case $p = 0$, if we remove the assumption of an absorbing barrier at the origin, $S(t)$ is a Gaussian process. In this case, it is straightforward to compute that the option pricing formula, sometimes called the *Normal*, *Gaussian* or *Bachelier* pricing formula with (Normal) volatility⁴ λ , becomes

$$c_N(t, S; T, K; \lambda) = (S - K)\Phi(d) + \lambda\sqrt{T-t}\phi(d), \quad d = \frac{S - K}{\lambda\sqrt{T-t}}, \quad (7.16)$$

where $\Phi(\cdot)$ and $\phi(\cdot)$ are the standard Gaussian CDF and PDF, respectively.

Further details about the non-central chi-square distribution can be found in Chapter 3. A number of efficient numerical algorithms exist to compute $\Upsilon(x, \nu, \gamma)$; see Johnson et al. [1995] for a survey. A standard algorithm can be found in Ding [1992]. Figure 7.1 on page 298 gives some examples of volatility skews produced by the CEV model.

⁴Also known as Gaussian volatility; when applied to interest rates, Gaussian volatilities are often called *basis-point*, or *bp*, volatilities.

7.2.3 Regularization

As discussed earlier, the CEV process implies a positive probability of absorption at $S = 0$ (for $p < 1$). This phenomenon is not necessarily a problem for pricing of simple European call options, but is obviously not desirable from an empirical standpoint⁵, and might also create some difficulties in pricing of more exotic structures. To avoid absorption, we can specify a regularized version of the CEV model by letting,

$$\varphi(x) = x \min(\varepsilon^{p-1}, x^{p-1}), \quad \varepsilon > 0, p < 1. \quad (7.17)$$

Roughly speaking, when $S(t)$ crosses the level ε , the resulting process becomes (locally) a geometric Brownian motion with finite volatility ε^{p-1} . With $\varphi(x)$ now Lipschitz continuous, it is straightforward to verify that the process for $S(t)$ can no longer reach the origin. On the other hand, the specification (7.17) will not allow for closed-form call option pricing but will, in principle at least, require the usage of numerical methods such as the finite difference method (see Section 7.4). On the other hand, for small to moderate values of ε , we would expect the CEV pricing formulas from Proposition 7.2.6 to hold as a good approximation. Andersen and Andreasen [2000b] verify numerically that this holds quite robustly, for strikes not too far from the spot value of S . More formally, we have the following result:

Proposition 7.2.10. *For $p < 1$ and $\varepsilon > 0$, let*

$$\begin{aligned} dx(t) &= \lambda x(t)^p dW(t), \\ dy(t) &= \lambda y(t) \min(\varepsilon^{p-1}, y(t)^{p-1}) dW(t), \end{aligned}$$

where $x(0) = y(0) > 0$ and $W(t)$ is a one-dimensional Brownian motion in measure P . For $p < 1/2$, 0 is assumed to be an absorbing boundary for x . For some $T > t$ and some constant K , we then have

$$\lim_{\varepsilon \downarrow 0} |P(x(T) < h) - P(y(T) < h)| = 0,$$

$$\lim_{\varepsilon \downarrow 0} \left| E((x(T) - K)^+) - E((y(T) - K)^+) \right| = 0.$$

The result is intuitive, but the proof is somewhat technical, and we skip it. Details can be found in Andersen and Andreasen [2000b].

⁵As the measure P is equivalent to the real-life (statistical) measure, a non-zero probability of absorption under P implies a non-zero probability of absorption under the real-life measure.

7.2.4 Displaced Diffusion Models

An easy extension of the CEV model that is sometimes useful involves adding a displacement constant to the CEV specification. Specifically, we write

$$\varphi(x) = (\alpha + x)^p \quad (7.18)$$

for some constant α . In the process (7.1), (7.18), let us set $Z(t) = \alpha + S(t)$. By Ito's lemma, $Z(t)$ then satisfies the CEV SDE

$$dZ(t) = \lambda Z(t)^p dW(t).$$

With $Z(t)$ having an absorbing boundary at 0, $S(t)$ then must have an absorbing boundary at $-\alpha$. Call option pricing with (7.18) is straightforward:

Proposition 7.2.11. *Let*

$$c_{\text{DCEV}}(t, S(t); T, K, \alpha) = \mathbb{E}_t \left((S(T) - K)^+ \right)$$

be the call option price associated with the displaced CEV process (7.1), (7.18). Then

$$c_{\text{DCEV}}(t, S; T, K, \alpha) = c_{\text{CEV}}(t, S + \alpha; T, K + \alpha), \quad S, K > -\alpha, \quad (7.19)$$

where the right-hand side is given by Proposition 7.2.6.

Proof. The result follows directly from the observation that

$$\mathbb{E}_t \left((S(T) - K)^+ \right) = \mathbb{E}_t \left((Z(T) - (K + \alpha))^+ \right),$$

where $Z(t) = \alpha + S(t)$ follows a regular CEV process. \square

Introduction of the displacement constant α allows for a (somewhat) richer family of volatility smiles than those of the pure CEV specification. In practice, however, the main use of displacement constants is for the special case of the *displaced log-normal*, or *shifted log-normal*, process where $p = 1$. The call option price formula for this case is listed below, for later reference.

Proposition 7.2.12. *Consider the displaced log-normal process*

$$dS(t) = \lambda(\beta + \zeta S(t)) dW(t), \quad (7.20)$$

where $W(t)$ is a one-dimensional Brownian motion in measure P , and $\zeta, \lambda \neq 0$. Assuming $S(t), K > -\beta/\zeta$, we have

$$\begin{aligned} c_{\text{DLN}}(t, S(t); T, K) &\triangleq \mathbb{E}_t \left((S(T) - K)^+ \right) \\ &= \left(S(t) + \frac{\beta}{\zeta} \right) \Phi(d_+) - \left(K + \frac{\beta}{\zeta} \right) \Phi(d_-), \\ d_{\pm} &= \frac{\ln \left(\frac{S(t) + \beta/\zeta}{K + \beta/\zeta} \right) \pm \frac{1}{2} \zeta^2 \lambda^2 (T - t)}{\zeta \lambda \sqrt{T - t}}. \end{aligned}$$

Proof. The result follows directly from the Black-Scholes equation (see Section 1.9) and (7.19), after setting $\alpha = \beta/\zeta$ and writing $\lambda(\beta + \zeta S(t)) = \lambda\zeta(\alpha + S(t))$. \square

Remark 7.2.13. It is often convenient to rewrite the displaced log-normal process in a slightly different form

$$dS(t) = \sigma (bS(t) + (1 - b)L) dW(t). \quad (7.21)$$

The parameter L is often set to near, or exactly at, the initial value $S(0)$. In this parameterization, σ is expressed in the units of relative volatility, just like in the Black model, because $bS(0) + (1 - b)L \approx S(0)$. In particular, σ always has the same scale for all values of b . Moreover, the effects of σ and b are almost “orthogonal”, in the sense that the parameter σ changes the overall level of the implied volatility smile but not its slope, whereas b only changes the slope (skew) of the implied volatility smile but not its overall level (i.e. not the at-the-money implied volatility). We use the parameterization (7.21) extensively in later chapters.

Remark 7.2.14. Consider the general local volatility model (7.1). Expanding the local volatility function $\varphi(\cdot)$ around at-the-money to the first order, we obtain

$$dS(t) \approx \lambda (\varphi(S(0)) + \varphi'(S(0))(S(t) - S(0))) dW(t),$$

which we identify as being of the form (7.21) with

$$\sigma = \lambda \frac{\varphi(S(0))}{S(0)}, \quad b = \varphi'(S(0)) \frac{S(0)}{\varphi(S(0))}, \quad L = S(0).$$

Hence, a first-order approximation to any local volatility process is of displaced log-normal type. In view of this, displaced log-normal processes are extensively used in various types of approximations and asymptotic expansions.

The previous remark can be applied to the CEV process:

$$\sigma = \lambda S(0)^{p-1}, \quad b = p, \quad L = S(0).$$

The approximation of the CEV process with (7.21) turns out to be particularly close, and we later use it to increase the tractability of certain stochastic volatility models. We also use it as a justification to freely switch from one type of process to the other. It is worth noting, however, that (7.20) has certain drawbacks relative to a pure CEV process. First, the process for $S(\cdot)$ can become negative if β (as is usual) is positive. Second, in stochastic volatility applications the asymptotic linear growth of $\varphi(x)$ in x can sometimes lead to technical problems and unbounded second moments of $S(\cdot)$. We shall return to this issue shortly, in Chapter 8.

7.3 Quadratic Volatility Model

In practice, volatility smiles in fixed income markets are not always perfectly monotonic in strike; indeed, as mentioned earlier, for sufficiently high strikes it is not uncommon for the smile to reverse direction and start increasing in strike. This type of behavior is inconsistent with a pure CEV model, but can, to some extent, be captured by the displaced CEV specification $\varphi(x) = (\alpha + x)^p$. Often, however, this model is hard to fit to actual data. A more powerful approach involves overlaying the CEV process with stochastic volatility, something that we turn to in Chapter 8. If we here wish to stay within the realm of DVF processes, one way to generate arbitrarily convex smiles is to use a *quadratic volatility model*, where

$$\varphi(x) = \alpha + \beta x + \gamma x^2, \quad (7.22)$$

for constants α, β, γ . We develop some aspects of this model here, but remind the reader of the caveats discussed in Section 7.1.3; in particular, for the model to be realistic, γ should probably be small.

Before commencing with derivations, let us note that the behavior of a DVF model (7.1) equipped with volatility function (7.22) will depend strongly on the root configuration in the quadratic polynomial $\alpha + \beta x + \gamma x^2$. For instance, if φ has two real roots l, u , $l < u$, straddling the initial value $S(0)$, it is clear that $S(t)$ will itself be bound to this range, i.e. $S(t) \in [l, u]$. Specifically, whenever $S(t)$ gets close to either l or u , $\varphi(x)$ will approach zero and the diffusion for $S(t)$ will gradually slow down. As such range-bound dynamics are rather unrealistic for interest rate applications⁶, we do not consider it in the following.

7.3.1 Case 1: Two Real Roots to the Left of $S(0)$

We first consider the case where $\alpha + \beta x + \gamma x^2$ has two real roots l and u , $l < u$, both lying to the left of $S(0)$. Without loss of generality, we may then consider the normalized process

$$dS(t) = \frac{(S(t) - u)(S(t) - l)}{u - l} dW(t), \quad S(0) > u > l. \quad (7.23)$$

We start by listing a few lemmas.

Lemma 7.3.1. *The range for $S(t)$ in (7.23) is $S(t) \in (u, \infty)$. In particular, the process for $S(t)$ does not explode in measure P .*

⁶For an application of the range-bound quadratic model to FX markets (where currency controls may potentially create upper and lower bounds), see Ingwersen [1997].

Proof. That $S(t)$ cannot go below u is obvious; further, Feller's boundary criteria (e.g. Karlin and Taylor [1981], Chapter 15.6) establishes that u is not accessible when $S(0) > u$. As $S(t)$ is described by a time-homogeneous one-dimensional SDE, it cannot explode (Karatzas and Shreve [1997], p. 332). \square

While the process for $S(t)$ is non-explosive, the super-linear growth⁷ of $\varphi(x)$ causes some interesting technical problems. In particular, we have the following result, proved in Andersen [2010].

Lemma 7.3.2. *The process (7.23) is a strict supermartingale in measure P .*

As the process for S is not a martingale, the usual pricing results require some modifications. For the purpose of pricing puts and calls, we need use the following.

Lemma 7.3.3. *Suppose that $S(t)$ satisfies (7.23) in some measure P and assume that put-call parity holds. Then the prices at time 0 for the put (p) and call (c) are*

$$\begin{aligned} p(0, S(0); T, K) &= E \left((K - S(T))^+ \right), \\ c(0, S(0); T, K) &= p(0, S(0); T, K) + S(0) - K > E \left((S(T) - K)^+ \right). \end{aligned}$$

Proof. (Sketch only). In the absence of arbitrage, the put price is a local martingale in measure P . As a bounded local martingale is a martingale and the put payout is bounded between 0 and $K - u$, it follows then that the put price in fact must be a true P -martingale. The expression for $p(0, S(0); T, K)$ follows. Applying put-call parity (see Chapter 1) yields the result for $c(0, S(0); T, K)$, where the inequality follows from Lemma 7.3.2. \square

We emphasize the non-standard result $c(0, S(0); T, K) > E(c(T, S(T); T, K))$ which is a consequence of the supermartingale property of $S(t)$. The inequality holds for arbitrarily large strikes; indeed, rather counter-intuitively, $\lim_{K \rightarrow \infty} c(0, S(0); T, K) = S(0) - E(S(T)) > 0$. We should also note that our assumption of put-call parity being valid is critical here, as it allows us to produce unique prices of both puts and calls. As described in Heston et al. [2007] and Andersen [2010], it is, however, possible to work with other assumptions without violating no-arbitrage.

With 7.3.3 we are now ready to tackle the derivation of an option pricing formula. We will be using the shorthand

$$p(t) \triangleq p(t, S(t); T, K),$$

and so forth. First, notice the useful relationship

⁷A similar issue is present in CEV processes with $p > 1$, as noted earlier.

$$S - K = \frac{(S - u)(K - l) - (K - u)(S - l)}{u - l}, \quad (7.24)$$

which allows us to write

$$\begin{aligned} p(T) &= \frac{1}{u - l} ((K - u)(S(T) - l) - (S(T) - u)(K - l))^+ \\ &= \frac{(K - u)(S(T) - l)}{u - l} \mathbf{1}_{\{(K-u)(S(T)-l)-(S(T)-u)(K-l)>0\}} \\ &\quad - \frac{(S(T) - u)(K - l)}{u - l} \mathbf{1}_{\{(K-u)(S(T)-l)-(S(T)-u)(K-l)>0\}} \\ &\triangleq p_1(T) - p_2(T). \end{aligned} \quad (7.25)$$

The payouts p_1 and p_2 have identical structure, so it suffices to focus our attention on pricing one of them, e.g. p_1 .

From Lemma 7.3.3, we have $p_1(0) = \mathbb{E}(p_1(T))$, which we rewrite as

$$p_1(0) = \frac{K - u}{u - l} \mathbb{E}((S(T) - l) \mathbf{1}_{\{(S(T)-u)/(S(T)-l) < (K-u)/(K-l)\}}). \quad (7.26)$$

At this point our first instinct would be to perform a measure shift that eliminates that factor $S(T) - l$ in the expectation, i.e. we would like to introduce a new measure $\tilde{\mathbb{P}}$ such that

$$\tilde{\mathbb{P}}(B) = \frac{1}{S(0) - l} \mathbb{E}((S(T) - l)B),$$

for any \mathcal{F}_T -measurable event B . We recall, however, that $S(t)$ (and therefore $S(t) - l$) is not a martingale in \mathbb{P} , so such a measure shift cannot be performed outright. Let us nevertheless try. Proceeding mechanically as if $S(t)$ were a martingale, we would get, for the process $Y(t) \triangleq (S(t) - u)/(S(t) - l)$,

$$dY(t) \stackrel{?}{=} Y(t) d\tilde{W}(t), \quad Y(0) = \frac{S(0) - u}{S(0) - l} < 1, \quad (7.27)$$

where \tilde{W} is a Brownian motion in $\tilde{\mathbb{P}}$. Clearly, however, there are technical problems here: the range for $Y(t)$ in (7.27) is $[0, \infty)$, whereas we know that in measure \mathbb{P} we have $Y(t) \in (0, 1)$ (since $S(t) \in (u, \infty)$); the two measures therefore cannot be equivalent. For option pricing purposes, it turns out that the correct way to handle the technical conflict involves inserting an *absorbing boundary* at $Y = 1$ in (7.27).

Proposition 7.3.4. *Let*

$$dY(t) = Y(t) d\tilde{W}(t), \quad Y(0) = \frac{S(0) - u}{S(0) - l} < 1,$$

be geometric Brownian motion in $\tilde{\mathbb{P}}$. Define $\tau = \inf\{t > 0 : Y(t) = 1\}$, and let $K > u$. Then $p_1(0)$ in (7.26) is given by

$$p_1(0) = \frac{(K-u)(S(0)-l)}{u-l} \mathbb{E}^{\tilde{P}} \left(1_{\{Y(T) < (K-u)/(K-l)\}} 1_{\{\tau > T\}} \right). \quad (7.28)$$

Stated explicitly,

$$p_1(0) = K_1 \Phi \left(\frac{-\ln(X_1/K_1) + T/2}{\sqrt{T}} \right) - X_2 \Phi \left(\frac{\ln(X_2/K_2) + T/2}{\sqrt{T}} \right), \quad (7.29)$$

with Φ being the Gaussian cumulative distribution function, and

$$\begin{aligned} K_1 &= \frac{(K-u)(S(0)-l)}{u-l}, & X_1 &= \frac{(S(0)-u)(K-l)}{u-l}, \\ K_2 &= \frac{(K-l)(S(0)-l)}{u-l}, & X_2 &= \frac{(S(0)-u)(K-u)}{u-l}. \end{aligned}$$

Proof. The result (7.28) is proven in Andersen [2010]. The result (7.29) follows by direct calculations, similar to those leading to the Black-Scholes-Merton formula. \square

Following similar steps leads to an expression for $p_2(0)$, which in turn leads to the following result for $p(0) = p_1(0) - p_2(0)$.

Proposition 7.3.5. *Let K_i, X_i , $i = 1, 2$, be given as in Proposition 7.3.4. Assuming $K > u$, the put price $p(0)$ for the model (7.23) has the explicit representation*

$$\begin{aligned} p(0, S(0); T, K) &= K_1 \Phi \left(-d_-^{(1)} \right) - X_2 \Phi \left(d_+^{(2)} \right) - X_1 \Phi \left(-d_+^{(1)} \right) + K_2 \Phi \left(d_-^{(2)} \right), \\ d_{\pm}^{(i)} &= \frac{\ln(X_i/K_i) \pm T/2}{\sqrt{T}}, \quad i = 1, 2. \end{aligned}$$

An application of put-call parity then immediately gives the call price:

Corollary 7.3.6. *Assuming put-call parity holds, the call price for the model (7.23) is*

$$c(0, S(0); T, K) = S(0) - K + p(0, S(0); T, K),$$

with $p(0, S(0); T, K)$ given in Proposition 7.3.5.

We recall that Proposition 7.3.5 applies to (7.23), rather than our original process which, at the root configuration in question, is

$$dS(t) = \lambda \gamma (S(t) - u) (S(t) - l) dW(t) = q \frac{(S(t) - u)(S(t) - l)}{u-l} dW(t), \quad (7.30)$$

where $q = \lambda \gamma(u - l)$. The constant in front of the quadratic polynomial is easily handled by time-scaling: to price options in (7.30) we simply set the put price equal to $p(0, S(0); q^2 T, K)$, where $p(0, S(0); \cdot, K)$ is given by the formula in Proposition 7.3.5.

7.3.2 Case 2: One Real Root to the Left of $S(0)$

If we let the single root to $\alpha + \beta x + \gamma x^2$ be denoted u , $u < S(0)$, it suffices to consider the normalized process

$$dS(t) = (S(t) - u)^2 dW(t). \quad (7.31)$$

But this process is a special case, with power equal to 2, of the displaced CEV model in Section 7.2.4, and the option pricing formulas from that section then apply directly. As these formulas are rather complicated in their dependence on the non-central chi-square distribution, it is worthwhile noticing that simple expressions exist for the special case of power equal to 2. The result is listed below.

Proposition 7.3.7. *For the process (7.31), the put option price is*

$$\begin{aligned} p(0, S(0); T, K) &= (S(0) - u)(K - u)\sqrt{T} \\ &\quad \times \{d_+ \Phi(d_+) + \phi(d_+) - d_- \Phi(d_-) - \phi(d_-)\}, \end{aligned}$$

where $\phi(x)$ is the Gaussian density, and

$$d_{\pm} = \frac{\pm \frac{1}{S(0)-u} - \frac{1}{K-u}}{\sqrt{T}}.$$

Proof. We observe that for the process

$$dS(t) = (S(t) - u)(S(t) - l) dW(t), \quad l < u < S(0),$$

the put price can be computed from the result in Proposition 7.3.5, after a time-change, from T to $T(u-l)^2$; see the comments at the end of Section 7.3.1. Taking the limit of the put price as $l \uparrow u$ then establishes the result. \square

The call option price can, as before, be found by put-call parity. To establish put and call option prices for the original diffusion

$$dS(t) = \lambda(\alpha + \beta S(t) + \gamma S(t)^2) dW(t) = \lambda\gamma(S(t) - u)^2 dW(t),$$

we simply change T to $\lambda^2\gamma^2T$ in Proposition 7.3.7.

7.3.3 Extensions and Other Root Configurations

The results listed in Sections 7.3.1 and 7.3.2 have given a flavor of how to deal with quadratic volatility process, and shall suffice for the purposes of this book. Other root configurations are treated in detail in Andersen [2010], including the case where $\varphi(x)$ has no roots (in which case the put and call option price formulas are infinite sine-series). Andersen [2010] also discusses the case where an absorbing barrier has been inserted at the origin to prevent $S(t)$ from going negative.

7.4 Finite Difference Solutions for General φ

For general specifications of φ , closed-form solutions for European options will not exist. In such cases, we may instead rely on the finite difference methods discussed in Chapter 2. Consider again the evaluation of

$$c(t, S(t); T, K) = \mathbb{E}_t \left((S(T) - K)^+ \right),$$

with $S(t)$ following (7.1). With suitable regularity conditions on φ , the Feynman-Kac theorem of Section 1.8 shows that $c(t, S) = c(t, S; T, K)$ (with T, K fixed) satisfies the PDE

$$\frac{\partial c(t, S)}{\partial t} + \frac{1}{2} \lambda^2 \varphi(S)^2 \frac{\partial^2 c(t, S)}{\partial S^2} = 0, \quad (7.32)$$

subject to the terminal condition

$$c(T, S) = (S - K)^+. \quad (7.33)$$

This PDE can be solved numerically using, say, the Crank-Nicholson finite difference grid method in Chapter 2. A direct discretization of (7.32) is normally sufficient, but we note that it may occasionally be possible to take advantage of special forms of φ and introduce transformations of S to improve the properties of the finite difference scheme. For example, as we have already seen in Chapter 2, when $\varphi(S) = S$, it is customary (and appropriate) to introduce $y(S) = \ln S$ and discretize in y . More generally, for sufficiently regular φ , the transformation

$$y(S) = \int \frac{dS}{\varphi(S)} \quad (7.34)$$

(see (2.81)–(2.82)) might offer numerical advantages over a direct discretization provided, of course, that the inverse in (7.34) exists. The following semi-heuristic argument explains the rationale. With the transform (7.34), the SDE for $y(t) = y(S(t))$ is (ignoring the drift)

$$dy(t) = O(dt) + \lambda dW(t).$$

The diffusion coefficient in the process for y is independent of the state of S , suggesting that a differential operator expressed in terms of y may have better numerical properties than the one expressed in terms of S . Even if y is not used for discretization, the transformation (7.34) suggests the discretization grid in the S -domain. In particular, $\{S_n\}_{n=0}^{m+1}$ can be defined by the condition that $y_n = y(S_n)$, $n = 0, \dots, m + 1$, are equidistant over $[y(S_0), y(S_{m+1})]$. For $n = 0, \dots, m + 1$ this gives ($y^{-1}(\cdot)$ is the inverse transform of (7.34))

$$y_n = y(S_0) + \frac{n}{m+1} (y(S_{m+1}) - y(S_0)),$$

$$S_n = y^{-1}(y_n) = y^{-1} \left(y(S_0) + \frac{n}{m+1} (y(S_{m+1}) - y(S_0)) \right).$$

7.4.1 Multiple λ and T

In applications, we often need to compute the values of $c(t, S; T, K)$ for many different values of T and/or λ . This need arises, for instance, in a standard model calibration exercise where we use a root-search algorithm to determine the value of λ that will make the computed call prices at different maturities T equal to values observed in the market. In such cases, we note that one should *not* simply solve (7.32) over and over (at great computational expense), but instead rely on the following observation:

Proposition 7.4.1. *Let $g(\tau, x)$ solve the following PDE*

$$-\frac{\partial g(\tau, x)}{\partial \tau} + \frac{1}{2}\varphi(x)^2 \frac{\partial^2 g(\tau, x)}{\partial x^2} = 0, \quad (7.35)$$

with initial condition

$$g(0, x) = (x - K)^+. \quad (7.36)$$

Let $c(t, S)$ solve the backward PDE (7.32)–(7.33) for a given value of λ . Then

$$c(t, S; T, K) = g(\lambda^2(T - t), S). \quad (7.37)$$

Proof. Follows directly from a variable transformation $\tau(t) = \lambda^2(T - t)$ in (7.32)–(7.33), taking advantage of the time-homogeneity of φ . \square

Using finite difference techniques to solve the PDE (7.35), we can construct the function g on a (τ, S) -grid; once this grid is stored in memory, (7.37) is used to recover $c(t, S; T, K)$ for arbitrary choices of S , λ and T by simple lookup or interpolation. We emphasize that this approach involves the numerical solution of only a *single* PDE. Also note that PDE is solved forward in time from a known *initial* condition, rather than backwards from a *terminal* condition.

7.4.2 Forward Equation for Call Options

While the function g from (7.35) is conveniently independent of T and λ , it does depend on K through the initial condition (7.36). In some applications, we may wish to use different strikes for different values of T , in which case the approach in Section 7.4.1 requires us to numerically solve as many finite difference grids as there are different values of K . We can improve on this by replacing the backward equation (7.32) with the forward equation of Dupire [1994]. In this approach, calendar time t and the initial value of S are considered fixed, whereas maturity T and strike K are variable. In view of this, we define $c(T, K) = c(t, S; T, K)$ for fixed t, S . We need the following proposition:

Proposition 7.4.2. Define the function $c(T, K) \triangleq c(t, S; T, K)$ where t, S are fixed and $c(t, S; T, K)$ is defined by (7.3) for the model (7.1). Then $c(T, K)$ satisfies the forward PDE

$$-\frac{\partial c(T, K)}{\partial T} + \frac{1}{2}\lambda^2\varphi(K)^2\frac{\partial^2 c(T, K)}{\partial K^2} = 0, \quad (7.38)$$

for $T > t$, subject to the time t initial condition

$$c(t, K) = (S - K)^+.$$

Proof. In Dupire [1994], the result is proven by combining the Fokker-Planck equation (see Section 1.8) with the result (7.5), followed by a series of integrations. A more intuitive line of attack proceeds as follows. Consider the function $H(t) = (S(t) - K)^+$. While $H(t)$ clearly does not satisfy the smoothness requirements of Ito's lemma, the Tanaka extension nevertheless justifies the following result, obtained by formally applying Ito's lemma to H :

$$dH(t) = 1_{\{S(t)>K\}}\lambda\varphi(S(t)) dW(t) + \frac{1}{2}\delta(S(t) - K)\lambda^2\varphi(S(t))^2 dt. \quad (7.39)$$

That is,

$$\begin{aligned} H(T) &= H(t) + \int_t^T 1_{\{S(u)>K\}}\lambda\varphi(S(u)) dW(u) \\ &\quad + \frac{1}{2} \int_t^T \delta(S(u) - K)\lambda^2\varphi(S(u))^2 du \\ &= H(t) + M(T) + \frac{1}{2} \int_t^T \delta(S(u) - K)\lambda^2\varphi(K)^2 du, \end{aligned}$$

where δ is the Dirac delta function and $M(t)$ is a continuous martingale with $M(t) = 0$. From (7.3), we have that

$$\begin{aligned} c(t, S(t); T, K) &= \mathbb{E}_t(H(T)) \\ &= H(t) + \frac{1}{2} \int_t^T \mathbb{E}_t(\delta(S(u) - K))\lambda^2\varphi(K)^2 du \\ &= H(t) + \frac{1}{2}\lambda^2\varphi(K)^2 \int_t^T \frac{\partial^2 c(t, S(t); u, K)}{\partial K^2} du, \end{aligned}$$

where we have used the martingale property of M as well as the result (7.5). Differentiating this equation with respect to T gives the result in Proposition 7.4.2. \square

As mentioned in Chapter 1, the term $\int_t^T \delta(S(u) - K) du$ in the expression for $H(T)$ is known as the *local time* of $S(\cdot)$ at the level K . Local time and the Tanaka extension are deep subjects (see Karatzas and Shreve [1997] for a

formal discussion) and have many interesting applications in finance, see for instance Andersen et al. [2002], Andersen and Andreasen [2000a], Andersen and Buffum [2003], Henderson and Hobson [2000], Carr and Jarrow [1990], Carr and Wu [2003], among many others.

We emphasize that while the backward equation (7.32) holds for European derivative securities on S in general, the forward equation (7.38) is unique to calls and puts, as only put and call payouts allow for the basic result (7.5).

Equipped with Proposition 7.4.2, the following result immediately follows from the proof of Proposition 7.4.1. Notice the difference in the initial conditions (7.36) and (7.40).

Proposition 7.4.3. *Let $h(\tau, x)$ solve the following PDE*

$$-\frac{\partial h(\tau, x)}{\partial \tau} + \frac{1}{2}\varphi(x)^2 \frac{\partial^2 h(\tau, x)}{\partial x^2} = 0,$$

with initial condition

$$h(0, x) = (S - x)^+. \quad (7.40)$$

Then

$$c(t, S; T, K) = h(\lambda^2(T - t), K).$$

As long as the initial value of $S(t)$ is kept constant, the result in Proposition 7.4.3 allows us use a single finite difference grid to price call options with multiple maturities, strikes, and λ 's. We note, however, that in many applications $S(t)$ may in fact be T -dependent, as S will often represent, say, T -maturity Libor forward rates. In such cases, the question of whether Proposition 7.4.3 leads to a more efficient numerical scheme than Proposition 7.4.1 is settled by comparing the number of strikes and the number of spot levels involved.

7.5 Asymptotic Expansions for General φ

As we have shown, there are a number of “tricks” that can be employed to make the application of finite difference methods a computationally viable approach to pricing a large number of European call options. Nevertheless, there is significant convenience and computer code simplification associated with closed-form pricing formulas, so we now turn to the development of asymptotic approximations for the solution to the generic backward PDE (7.32). There are a number of approaches that can be taken, including the “most likely path” method in Gatheral [2001] (see also Gatheral [2006] and Section 22.1.7, and the singular perturbation techniques in Hagan and Woodward [1999b], Henry-Labordé [2008], Gatheral et al. [2009], to name a few. Our presentation here is based on a fairly straightforward, yet often highly accurate, asymptotic expansion in time to maturity.

7.5.1 Expansion around Displaced Log-Normal Process

As in Proposition 7.4.1, we start by writing $c(t, S; T, K) = g(\tau, S)$, where $\tau = \lambda^2(T - t)$ and g satisfies (7.35). Inspired by the known solution of (7.35) in Proposition 7.2.12 for the case $\varphi(x) = \beta + \zeta x$, $\zeta \neq 0$, let us guess at a solution of (7.35) of the form

$$g(\tau, S) = \left(S + \frac{\beta}{\zeta} \right) \Phi(z_+) - \left(K + \frac{\beta}{\zeta} \right) \Phi(z_-), \quad (7.41)$$

$$z_{\pm} = \frac{\ln \left(\frac{S+\beta/\zeta}{K+\beta/\zeta} \right) \pm \frac{1}{2}\Omega(\tau, S)^2}{\Omega(\tau, S)},$$

where the function $\Omega(\tau, S)$ is to be determined. In (7.41), note that we obviously must assume that $S, K > -\beta/\zeta$. Substituting (7.41) into (7.35) gives the following PDE for $\Omega(\tau, S)$:

$$\begin{aligned} & \left(S + \frac{\beta}{\zeta} \right)^2 \Omega \frac{\partial \Omega}{\partial \tau} \\ &= \frac{1}{2} \varphi(S)^2 \left[\left(S + \frac{\beta}{\zeta} \right)^2 \Omega \frac{\partial^2 \Omega}{\partial S^2} + (1 - h_{-3}) - h_1(1 - h_{-1}) \right], \end{aligned} \quad (7.42)$$

where

$$h_i \triangleq \left(S + \frac{\beta}{\zeta} \right) \frac{\partial \Omega}{\partial S} \left(\Omega^{-1} \ln \left(\frac{S + \beta/\zeta}{K + \beta/\zeta} \right) + \frac{1}{2} i \Omega \right), \quad i = -3, -1, 1.$$

The PDE (7.42) does not generally allow for an explicit solution, so we resort to an asymptotic expansion in τ .

Proposition 7.5.1. *An asymptotic expansion for the solution of (7.35) is given by (7.41), with*

$$\Omega(\tau, S) = \Omega_0(S)\tau^{1/2} + \Omega_1(S)\tau^{3/2} + O\left(\tau^{5/2}\right), \quad (7.43)$$

$$\Omega_0(S) = \ln \left(\frac{S + \beta/\zeta}{K + \beta/\zeta} \right) \left(\int_K^S \varphi(u)^{-1} du \right)^{-1}, \quad (7.44)$$

$$\Omega_1(S) = -\frac{\Omega_0(S)}{\left(\int_K^S \varphi(u)^{-1} du \right)^2} \ln \left(\Omega_0(S) \left(\frac{(S + \beta/\zeta)(K + \beta/\zeta)}{\varphi(S)\varphi(K)} \right)^{1/2} \right), \quad (7.45)$$

where the parameters β and ζ can be chosen arbitrarily, subject to the constraints $S, K > -\beta/\zeta$ and $\zeta \neq 0$.

Proof. In (7.41) we clearly require $\Omega(\tau, S) \sim \tau^{1/2}$ as $\tau \rightarrow 0$, so we seek a small-time solution of the form

$$\Omega(\tau, S) = \sum_{i \geq 0} \tau^{i+1/2} \Omega_i(S). \quad (7.46)$$

Notice that (7.46) omits all integer powers of τ — it turns out that their weights are all identically 0. Substituting (7.46) into (7.42) and matching terms of order $O(1)$ gives

$$(S + \beta/\zeta)^2 \Omega_0^2 = \varphi(S)^2 \left(1 - \frac{\Omega'_0}{\Omega_0} (S + \beta/\zeta) \ln \left(\frac{S + \beta/\zeta}{K + \beta/\zeta} \right) \right)^2, \quad (7.47)$$

where the prime denotes differentiation with respect to S . Taking the square root of the above equation and rearranging leads to two first-order ordinary differential equations of the Bernoulli type. Solving (7.47) subject to the boundary condition that the limit of Ω_0 must be finite for $S \rightarrow K$ (and discarding the negative solution) leads to (7.44).

Progressing now to the $O(\tau)$ term in (7.42), we get

$$\begin{aligned} 2(S + \beta/\zeta) \Omega_1 &= \frac{1}{2} \varphi(S)^2 \left((S + \beta/\zeta) \Omega_0'' + \Omega'_0 \right) \\ &\quad - \varphi(S) (S + \beta/\zeta) \ln \left(\frac{S + \beta/\zeta}{K + \beta/\zeta} \right) \left(\frac{\Omega'_1}{\Omega_0} - \frac{\Omega_1 \Omega'_0}{\Omega_0^2} \right). \end{aligned}$$

Inserting the result for Ω_0 and rearranging again leads to a Bernoulli-type ODE, the explicit solution of which is (7.45). As before, we have ensured that the limit $S \rightarrow K$ is finite. \square

Remark 7.5.2. We notice that $\Omega_0(K)$ and $\Omega_1(K)$ in Proposition 7.5.1 exist by construction. Taking the limit $S \rightarrow K$ explicitly, we get

$$\begin{aligned} \Omega_0(K) &= \frac{\varphi(K)}{K + \beta/\zeta}, \\ \Omega_1(K) &= \frac{1}{24} \Omega_0(K)^3 \left[1 + (K + \beta/\zeta)^2 \varphi(K)^{-2} (2\varphi(K)\varphi''(K) - \varphi'(K)^2) \right]. \end{aligned}$$

While Proposition 7.5.1 only includes two terms in the expansion for Ω , it is possible to compute further terms if necessary. Such terms become increasingly cumbersome however, and typically do not add much further accuracy.

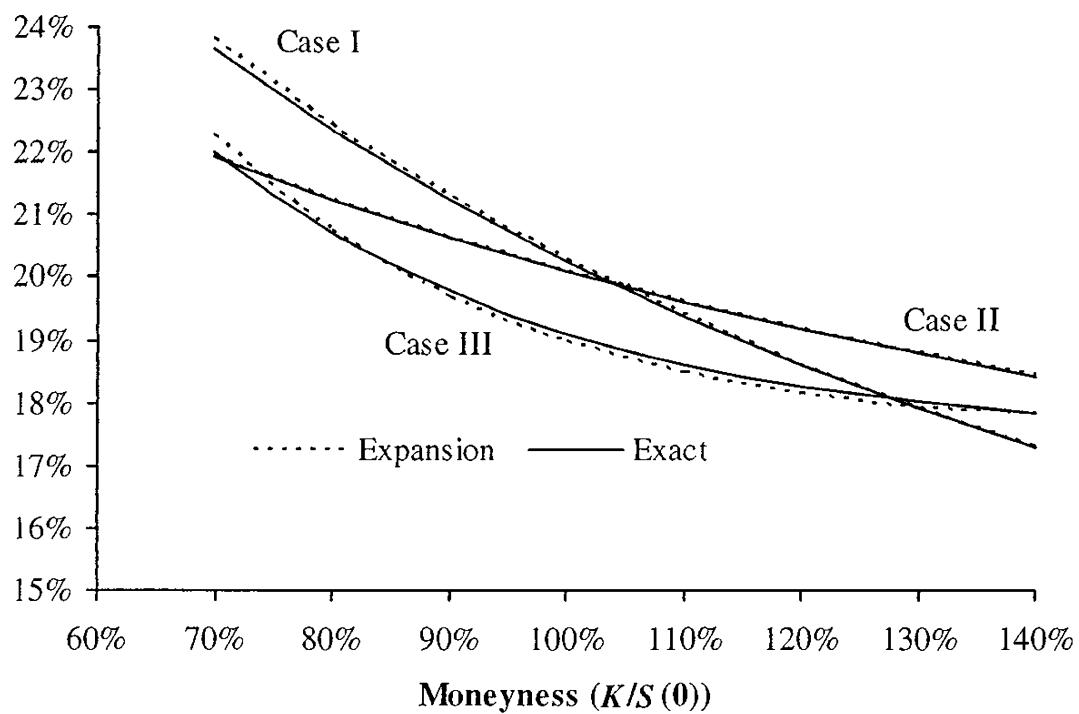
The best choice of the parameters β and ζ is not always obvious. One choice is to use Remark 7.2.14. Alternatively, we could think of a more global approach and, roughly speaking, set them in such a way that the straight line $\beta + \zeta x$ would provide as good a fit to $\varphi(x)$ as possible, over the

statistically relevant range of x . Sometimes, we can use a Taylor expansion around $x = (S + K)/2$, say, and set

$$\zeta = \varphi'((S + K)/2), \quad \beta = \varphi((S + K)/2) - \zeta(S + K)/2.$$

We note that when $\beta = 0$, $\Omega(\lambda^2(T-t), S(t))/\sqrt{T-t}$ in Proposition 7.5.1 conveniently becomes the time t implied Black volatility σ_B discussed earlier. For a few selected φ , Figure 7.1 below compares σ_B computed from the expansion in Proposition 7.5.1 (with $\beta = 0$) against exact results. Despite the long option maturity used in the figure, precision of the expansion is excellent, especially for the CEV case.

Fig. 7.1. Implied Volatility



Notes: The graph shows the implied volatility for a 10 year option, as a function of option moneyness $K/S(0)$. The initial value of the underlying is $S(0) = 6\%$. Three DVF models are considered in the figure. Case I: $\varphi(x) = x^{0.1}$, $\lambda = 1.59\%$. Case II: $\varphi(x) = x^{0.5}$, $\lambda = 4.90\%$. Case III: $\varphi(x) = x(1 + 30e^{-10x})$, $\lambda = 16.75\%$. The “Expansion” numbers in the graph were computed from the result in Proposition 7.5.1 with $\beta = 0$. For Case I and Case II, the “Exact” numbers were computed by the known CEV pricing formula in Proposition 7.2.6; for Case III the “Exact” numbers were computed in a Crank-Nicholson finite difference grid with 150 time steps and 250 spatial steps.

7.5.2 Expansion around Gaussian Process

For cases where φ is close to a constant, one might like to base the asymptotic expansion on $\varphi(x) = \beta$, for some constant β . In this case (which violates

one of the restrictions in Proposition 7.5.1), we use the Gaussian formula (7.16), and write

$$g(\tau, S) = (S - K) \Phi(w) - \Psi(\tau, S) \phi(w), \quad w = \frac{S - K}{\Psi(\tau, S)}. \quad (7.48)$$

For completeness, an asymptotic expansion of $\Psi(\tau, S)$ is given below.

Proposition 7.5.3. *An asymptotic expansion for the solution of (7.35) is given by (7.48), with*

$$\begin{aligned} \Psi(\tau, S) &= \Psi_0(S)\tau^{1/2} + \Psi_1(S)\tau^{3/2} + O\left(\tau^{5/2}\right), \\ \Psi_0(S) &= (S - K) \left(\int_K^S \varphi(u)^{-1} du \right)^{-1}, \\ \Psi_1(S) &= -\frac{\Psi_0(S)^3}{(S - K)^2} \ln \left(\Psi_0(S) (\varphi(S)\varphi(K))^{-1/2} \right). \end{aligned}$$

In Proposition 7.5.3, the limit $S \rightarrow K$ leads to the following expressions

$$\begin{aligned} \Psi_0(K) &= \varphi(K), \\ \Psi_1(K) &= \frac{1}{24} \varphi(K) (2\varphi(K)\varphi''(K) - \varphi'(K)^2). \end{aligned}$$

The proof of Proposition 7.5.3 is similar to that of Proposition 7.5.1 and is omitted. Note that $\Psi(\lambda^2(T-t), S(t))/\sqrt{T-t}$ can be interpreted as an *implied Normal volatility*.

7.6 Extensions to Time-Dependent φ

So far, we have limited our discussion to the case where the function φ is independent of calendar time t . While there is some danger in making φ a function of t — the model inevitably becomes less time-stationary — there are a number of applications where such an extension is necessary to improve the fit to market data. Unlike the non-parametric approaches in Dupire [1994], Derman and Kani [1994], and Andersen and Brotherton-Ratcliffe [1998] (and many others) where $\varphi(t, S)$ is calibrated to fit a double continuum of call option prices, the applications we have in mind are normally parametric, and are inspired by typical requirements of calibrating term structure models to swaptions and caplets.

By itself, swaption and caplet pricing does not require time-dependent parameters, as only the terminal distribution is relevant. From that point of view, vanilla models with time-dependent local volatility functions may appear to have limited use in fixed income modeling. However, they often arise as describing approximate dynamics of swap or Libor rates in term

structure models. Many examples of such approximations are given in later chapters (see Chapters 13 and 14, for instance), and handling time-dependent parameters in local volatility models is important for term structure model calibration.

7.6.1 Separable Case

Recall the basic SDE (7.1). Its simplest time-dependent extension specifies a time-dependent scaling volatility λ , $\lambda = \lambda(t)$:

$$dS(t) = \lambda(t)\varphi(S(t)) dW(t). \quad (7.49)$$

This is the so-called *separable case*, as the local volatility function is represented as a product of two functions: $\lambda(\cdot)$, a function of the time variable only, and $\varphi(\cdot)$, a function of the state variable only. The separable form allows for application of the following simple time change argument:

Proposition 7.6.1. *Define*

$$\tau(t) = \int_0^t \lambda(u)^2 du,$$

and define $s(\cdot)$ by $S(t) = s(\tau(t))$, with $S(t)$ following (7.49). Then

$$ds(\tau) = \varphi(s(\tau)) d\widetilde{W}(\tau), \quad s(0) = S(0), \quad (7.50)$$

where \widetilde{W} is a Brownian motion.

Proof. The result follows directly from standard results for time-changed Brownian motion, see e.g. Karatzas and Shreve [1997]. \square

Consider now the valuation of

$$c(t, S(t); T, K) = E_t \left((S(T) - K)^+ \right),$$

which in the notation of Proposition 7.6.1 can be written as

$$c(t, s(\tau(t)); T, K) = E \left((s(\tau(T)) - K)^+ \middle| \tilde{\mathcal{F}}_{\tau(t)} \right),$$

where $\tilde{\mathcal{F}}$ is the filtration generated by \widetilde{W} . As the process for $s(\cdot)$ in (7.50) is of the type (7.1) (with $\lambda = 1$), all results from previous sections hold unchanged after the simple substitutions $\lambda \mapsto 1$ and $(T-t) \mapsto (\tau(T)-\tau(t))$. Equivalently, whenever the European option price results for constant λ involve terms of the form $\lambda^2(T-t)$, they should be replaced with $\int_t^T \lambda(u)^2 du$ to accommodate a time-varying $\lambda(\cdot)$.

7.6.2 Skew Averaging

While the separable case can be handled quite easily, it is often too restrictive to be truly useful. Consider therefore the general case

$$dS(t) = \varphi(t, S(t)) dW(t), \quad (7.51)$$

for $\varphi(t, x)$ satisfying the standard regularity conditions. European options could be valued in this model by PDE methods without much difficulty. However, with calibration applications in mind, this may be too slow or insufficiently accurate.

In this section, we develop European option approximations based on the idea of *time averaging*. Given the SDE (7.51), we look for a model with a *time-independent* local volatility function that yields European option prices approximately matching prices from the time-dependent model. The time-independent local volatility function can then be interpreted as a time average of the time-dependent function. This reduces the problem to one we know how to solve.

We have already seen a flavor of the averaging results that we are looking for. As demonstrated in Section 7.6.1, the values of T -expiry European options in the model

$$dS(t) = \lambda(t)\varphi(S(t)) dW(t) \quad (7.52)$$

are the same as in the model

$$dS(t) = \bar{\lambda}\varphi(S(t)) dW(t),$$

where $\bar{\lambda}$ is given by

$$\bar{\lambda}^2 = \int_0^T \lambda(u)^2 du.$$

Thus, $\bar{\lambda}$ is an *effective volatility* for expiry T for the model (7.52).

Given the comments on *U-shaped* local volatility functions in Section 7.1.3, our initial focus shall be on functions that are monotonic in the state variable (see Section 7.6.3 for extensions). Such functions are typically well-described by two parameters, with the first parameter governing the overall level of volatility and the second the slope of the volatility smile (or skew). In the general case, both parameters are time-dependent. Let us concentrate on finding the averaging result for the time-dependent skew or, equivalently, on finding the *effective skew* formula.

We apply asymptotic expansion techniques with the *slope of the local volatility function* being the small parameter. Let us denote

$$X_0 = S(0), \quad \lambda(t) = \varphi(t, X_0), \quad g(t, x) = \frac{\varphi(t, x)}{\varphi(t, X_0)}.$$

Then (7.51) can be rewritten as

$$dS(t) = \lambda(t)g(t, S(t)) dW(t), \quad (7.53)$$

where

$$g(t, X_0) = 1. \quad (7.54)$$

Let us fix a time horizon $T > 0$ and attempt to derive conditions that a time-independent function $\bar{g}(x)$ needs to satisfy so that the SDE (7.53) can be replaced with

$$dS(t) = \lambda(t)\bar{g}(S(t)) dW(t) \quad (7.55)$$

for the purposes of valuing T -expiry European options of all strikes. Without loss of generality, the function $\bar{g}(x)$ is assumed to satisfy

$$\bar{g}(X_0) = 1.$$

Choose $\epsilon \geq 0$, the small slope parameter, and define

$$g^\epsilon(t, x) = g(t, X_0 + (x - X_0)\epsilon), \quad \bar{g}^\epsilon(x) = \bar{g}(X_0 + (x - X_0)\epsilon).$$

Next, define two sets of processes

$$\begin{aligned} dX^\epsilon(t) &= \lambda(t)g^\epsilon(t, X^\epsilon(t)) dW(t), \quad X^\epsilon(0) = X_0, \\ dY^\epsilon(t) &= \lambda(t)\bar{g}^\epsilon(Y^\epsilon(t)) dW(t), \quad Y^\epsilon(0) = X_0. \end{aligned}$$

The requirement that the prices of European options on $X^\epsilon(T)$ and $Y^\epsilon(T)$ across all strikes be close can be reformulated as the requirement that the distributions of $X^\epsilon(T)$ and $Y^\epsilon(T)$ be close. This can be formalized as finding $\bar{g}(\cdot)$ such that

$$q(\epsilon) \rightarrow \min,$$

where

$$q(\epsilon) \triangleq \mathbb{E} \left((X^\epsilon(T) - Y^\epsilon(T))^2 \right). \quad (7.56)$$

Considering the small slope limit $\epsilon \rightarrow 0$, we expand $q(\epsilon)$ in powers of ϵ to obtain

$$q(\epsilon) = q(0) + q'(0)\epsilon + \frac{1}{2}q''(0)\epsilon^2 + O(\epsilon^3).$$

As part of the proof below we will show that $q(0) = q'(0) = 0$. Hence, the minimization problem simplifies to

$$q''(0) \rightarrow \min.$$

The (necessary) minimum condition is given in the following result.

Proposition 7.6.2. *Any function \bar{g} that minimizes $q''(0)$ must satisfy the condition*

$$\frac{\partial \bar{g}(X_0)}{\partial x} = \int_0^T \frac{\partial g(t, X_0)}{\partial x} w_T(t) dt, \quad (7.57)$$

where

$$w_T(t) = \frac{v(t)^2 \lambda(t)^2}{\int_0^T v(t)^2 \lambda(t)^2 dt}, \quad v(t)^2 \triangleq \mathbb{E} \left((X^0(t) - X_0)^2 \right). \quad (7.58)$$

Proof. By Theorem 1.1.3, $q(\epsilon)$ as defined by (7.56) must equal

$$q(\epsilon) = \mathbb{E} \left(\int_0^T (g^\epsilon(t, X^\epsilon(t)) - \bar{g}^\epsilon(Y^\epsilon(t)))^2 \lambda(t)^2 dt \right).$$

Differentiating with respect to ϵ , we get (omitting arguments on g^ϵ and \bar{g}^ϵ for brevity)

$$q'(\epsilon) = 2\mathbb{E} \left(\int_0^T (g^\epsilon - \bar{g}^\epsilon) \times \left(\frac{\partial}{\partial \epsilon} g^\epsilon - \frac{\partial}{\partial \epsilon} \bar{g}^\epsilon \right) \lambda(t)^2 dt \right), \quad (7.59)$$

$$\begin{aligned} q''(\epsilon) &= 2\mathbb{E} \left(\int_0^T \left(\frac{\partial}{\partial \epsilon} g^\epsilon - \frac{\partial}{\partial \epsilon} \bar{g}^\epsilon \right)^2 \lambda(t)^2 dt \right) \\ &\quad + 2\mathbb{E} \left(\int_0^T (g^\epsilon - \bar{g}^\epsilon) \left(\frac{\partial^2}{\partial \epsilon^2} g^\epsilon - \frac{\partial^2}{\partial \epsilon^2} \bar{g}^\epsilon \right) \lambda(t)^2 dt \right). \end{aligned}$$

Since $g^0(t, x) = \bar{g}^0(x) = 1$, it follows that $q(0)$, $q'(0)$ and the second term in the expression for $q''(0)$ are zero. Hence,

$$q''(0) = 2\mathbb{E} \left(\int_0^T \left(\frac{\partial}{\partial \epsilon} g^\epsilon(t, X^\epsilon(t)) \Big|_{\epsilon=0} - \frac{\partial}{\partial \epsilon} \bar{g}^\epsilon(Y^\epsilon(t)) \Big|_{\epsilon=0} \right)^2 \lambda(t)^2 dt \right).$$

Note that

$$\begin{aligned} \frac{\partial}{\partial \epsilon} g^\epsilon(t, X^\epsilon(t)) &= \left[\epsilon \left(\frac{\partial X^\epsilon}{\partial \epsilon} \right)(t) + (X^\epsilon(t) - X_0) \right] \\ &\quad \times \frac{\partial g}{\partial x}(t, X_0 + \epsilon(X^\epsilon(t) - X_0)), \\ \frac{\partial}{\partial \epsilon} \bar{g}^\epsilon(Y^\epsilon(t)) &= \left[\epsilon \left(\frac{\partial Y^\epsilon}{\partial \epsilon} \right)(t) + (Y^\epsilon(t) - X_0) \right] \\ &\quad \times \frac{\partial \bar{g}}{\partial x}(X_0 + \epsilon(Y^\epsilon(t) - X_0)). \end{aligned}$$

In particular, as $Y^0(t) = X^0(t)$,

$$\begin{aligned} \frac{\partial}{\partial \epsilon} g^\epsilon(t, X^\epsilon(t)) \Big|_{\epsilon=0} &= (X^0(t) - X_0) \frac{\partial g}{\partial x}(t, X_0), \\ \frac{\partial}{\partial \epsilon} \bar{g}^\epsilon(Y^\epsilon(t)) \Big|_{\epsilon=0} &= (X^0(t) - X_0) \frac{\partial \bar{g}}{\partial x}(X_0). \end{aligned}$$

Thus,

$$\begin{aligned} q''(0) &= 2 \int_0^T \mathbb{E} \left((X^0(t) - X_0)^2 \lambda(t)^2 \right) \left(\frac{\partial g}{\partial x}(t, X_0) - \frac{\partial \bar{g}}{\partial x}(X_0) \right)^2 dt \\ &= 2 \int_0^T v(t)^2 \lambda(t)^2 \left(\frac{\partial g}{\partial x}(t, X_0) - \frac{\partial \bar{g}}{\partial x}(X_0) \right)^2 dt, \end{aligned}$$

with $v(t)^2$ defined in (7.58). Differentiating with respect to the slope $\partial \bar{g}(X_0)/\partial x$ and setting the resulting derivative to zero, we obtain a condition for the minimum of $q''(0)$. This gives (7.57). \square

It follows from the proposition that for the purposes of (approximately) pricing T -expiry European options, (7.53) can be replaced with (7.55), where $\bar{g}(\cdot)$ is a function whose slope (skew) at-the-money, $\partial \bar{g}(S(0))/\partial x$, is a weighted average of the time-dependent at-the-money slopes (skews) of the original function $\partial g(t, S(0))/\partial x$, $t \in [0, T]$. The weights $w_T(t)$ to be used in forming the slope-average are the weights $w_T(t)$ in (7.58). Once the SDE of the form (7.53) has been approximated with (7.55), various tools developed in the first part of the chapter become available, and European option prices can be computed efficiently.

7.6.2.1 Examples

The time-dependent local volatility function $g(t, x)$ is often defined to be a time-indexed collection of functions from the same family. Examples include the time-dependent displaced log-normal function

$$g(t, x) = b(t) \frac{x}{S(0)} + (1 - b(t)), \quad t \in [0, T], \quad (7.60)$$

or the time-dependent CEV function

$$g(t, x) = \left(\frac{x}{S(0)} \right)^{p(t)}, \quad t \in [0, T]. \quad (7.61)$$

Note that the functions in the formulas have been scaled to satisfy (7.54). The condition (7.57) does not define the function \bar{g} uniquely. To improve the accuracy of the approximation, it is often beneficial to choose \bar{g} from the same family as the functions they approximate. In particular, for g of the type (7.60), the function \bar{g} is best chosen to be of the same displaced log-normal type

$$\bar{g}(x) = \bar{b} \frac{x}{S(0)} + (1 - \bar{b}). \quad (7.62)$$

In the same vein, for the CEV case (7.61), a natural choice for \bar{g} is

$$\bar{g}(x) = \left(\frac{x}{S(0)} \right)^{\bar{p}}. \quad (7.63)$$

Both the displaced log-normal parameter b and the CEV parameter p are used as a measure of the skew in the implied volatility smile. The next

corollary expressed the averaging result directly in terms of these parameters, and also explicitly derives the averaging weights.

Corollary 7.6.3. *Over the time-horizon $[0, T]$, the effective skew \bar{b} in (7.62) for the model defined by the time-dependent local volatility function (7.60) is given by*

$$\bar{b} = \int_0^T b(t) w_T(t) dt, \quad (7.64)$$

where

$$w_T(t) = \frac{v(t)^2 \lambda(t)^2}{\int_0^T v(s)^2 \lambda(s)^2 ds}, \quad v(t)^2 = \int_0^t \lambda(s)^2 ds. \quad (7.65)$$

Proof. For $g(t, x)$ and $\bar{g}(x)$ given by (7.60) and (7.62), we have

$$\frac{\partial g}{\partial x}(t, S(0)) = \frac{b(t)}{S(0)}, \quad \frac{\partial \bar{g}}{\partial x}(S(0)) = \frac{\bar{b}}{S(0)}.$$

Thus, (7.64) follows from (7.57). The formula (7.65), and in particular the expression for $v(t)^2$, follows from the definition

$$v(t)^2 = E \left((X^0(t) - X_0)^2 \right)$$

and the fact that $X^0(t)$ satisfies

$$dX^0(t) = \lambda(t) g^0(t, X^0(t)) dW(t)$$

with

$$g^0(t, X^0(t)) \equiv 1.$$

□

Remark 7.6.4. An identical result holds for the effective CEV parameter \bar{p} ,

$$\bar{p} = \int_0^T p(t) w_T(t) dt,$$

where $p(\cdot)$ and \bar{p} are the parameters in (7.61) and (7.63), and $w_T(t)$ is as given in (7.65).

Example 7.6.5. Assuming constant volatility $\lambda(t) \equiv \lambda$, we obtain particularly simple formulas for the effective skew,

$$v(t)^2 = \lambda^2 t, \quad w_T(t) = \frac{t}{\int_0^T s ds} = \frac{t}{T^2/2},$$

so that

$$\bar{b} = \frac{1}{T^2/2} \int_0^T tb(t) dt.$$

This demonstrates that instantaneous skews $b(t)$ for larger t contribute more to \bar{b} than those for lower t . Intuitively, the process needs to build up its variance before the changes in the instantaneous slopes start having an effect on the effective slope of the local volatility.

7.6.2.2 A Caveat About the Process Domain

Even though the skew averaging result is obtained in the small slope limit, practical experience validates its broad applicability in option pricing problems. Some typical results can be found in Piterbarg [2005c] and Piterbarg [2006]. Still, the equivalence between the original time-dependent model and the time-averaged one should not be taken too far, as we now proceed to demonstrate. For this, we focus on the simple displaced diffusion model from the previous section, i.e. we consider the time-dependent SDE

$$dS(t) = \lambda (b(t)S(t) + (1 - b(t)) S(0)) dW(t), \quad (7.66)$$

and approximate it with

$$dS(t) = \lambda (\bar{b}S(t) + (1 - \bar{b}) S(0)) dW(t), \quad (7.67)$$

where \bar{b} is set as in Corollary 7.6.3. While the two SDEs (7.66) and (7.67) may have similar properties in the neighborhood of $S(0)$, they generally do not even have the same range for $S(t)$. For the constant parameter case (7.67) with $\bar{b} > 0$, the process $S(t)$ has a lower bound, the root of the local volatility function: $S(t) \in (S(0)(\bar{b} - 1)/\bar{b}, \infty)$. The same is not necessarily true for the time-dependent SDE (7.66), as should be reasonably clear from the following heuristic argument. If at a given time t , $S(t)$ is close to the root of the local volatility function but still above it, i.e.

$$S(t) \gtrsim S(0) (b(t) - 1) / b(t),$$

it may so happen that at $t + dt$, $S(t + dt)$ is actually *below* the root of the local volatility function,

$$S(t + dt) < S(0) (b(t + dt) - 1) / b(t + dt)$$

due to the change in the function $b(\cdot)$. The range

$$(-\infty, S(0) (b(t + dt) - 1) / b(t + dt))$$

will then be reachable by $S(\cdot)$. The following proposition provides formal justification.

Proposition 7.6.6. *Consider the SDE*

$$dX(t) = (a(t) + b(t)X(t)) dW(t) \quad (7.68)$$

with $X(0) \geq a(0)$. If $a'(u) \leq 0$ for all $u \in [0, t]$, then $X(t) > a(t)$ a.s. If there exists u , $0 < u < t$, such that $a'(u) > 0$, then $P(X(t) < l) > 0$ for any $l \in \mathbb{R}$.

Proof. Define

$$\zeta(t) = \int_0^t b(u) dW(u) - \frac{1}{2} \int_0^t b^2(u) du, \quad Z(t) = \exp(\zeta(t)).$$

Then the solution to the SDE (7.68) is given by

$$X(t) = Z(t) \left[X(0) + \int_0^t a(u) d(1/Z(u)) \right],$$

as can either be checked directly or obtained from Section 5.6.C of Karatzas and Shreve [1997]. Integrating by parts yields

$$X(t) = Z(t) (X(0) - a(0)) + a(t) - Z(t) \int_0^t \frac{a'(u)}{Z(u)} du.$$

With $X(0) \geq a(0)$,

$$Z(t) (X(0) - a(0)) + a(t)$$

is bounded from below by $a(t)$. If $a'(u) \leq 0$ for all $u \in [0, t]$ then the remaining term

$$-Z(t) \int_0^t \frac{a'(u)}{Z(u)} du$$

is positive and $X(t)$ is bounded from below by $a(t)$. If, however, there exists u such that $a'(u) > 0$, this term can be arbitrarily negative with positive probability. \square

In practice, the likelihood of actually breaching the lower boundary is typically small and we can often safely ignore this possibility. If needed, one can always “regularize” the time-dependent process to limit its range, along the same lines as done in Section 7.2.3.

7.6.3 Skew and Convexity Averaging by Small-Noise Expansion

The technique used in the previous section to derive Proposition 7.6.2 is not the only route to go. An alternative approach relies on *small-noise expansion*, a concept closely related to the Ito-Taylor expansion in Chapter 3. To illustrate the versatility of this method, we shall use it to derive not only the skew averaging result in Corollary 7.6.3, but also to demonstrate how to compute *average convexity* in a time-dependent quadratic model.

As our starting point, we define, for some constant X_0 , the quadratic form

$$\begin{aligned} \varphi(t, X(t)) &= \varphi(b(t), c(t), X(t)) \\ &= (1 - b(t)) X_0 + b(t) X(t) + \frac{1}{2} c(t) (X(t) - X_0)^2, \end{aligned}$$

and then introduce the following two processes:

$$dX(t) = \lambda(t)\varphi(b(t), c(t), X(t)) dW(t), \quad X(0) = X_0, \quad (7.69)$$

$$dY(t) = \lambda(t)\varphi(\bar{b}, \bar{c}, Y(t)) dW(t), \quad Y(0) = X_0, \quad (7.70)$$

where $W(t)$ is a Brownian motion in some probability measure. We can characterize the process for $X(t)$ as having quadratic local volatility with time-dependent slope $b(t)$ and time-dependent convexity $c(t)$; for a fixed value of T , we are interested in establishing how to set the constants \bar{b} in \bar{c} in the process for Y such that $Y(T)$ is a good approximation to $X(T)$.

We will answer the question posed above in the small-noise limit. For that, set

$$dX^\epsilon(t) = \epsilon\lambda(t)\varphi(b(t), c(t), X^\epsilon(t)) dW(t), \quad (7.71)$$

$$dY^\epsilon(t) = \epsilon\lambda(t)\varphi(\bar{b}, \bar{c}, Y^\epsilon(t)) dW(t), \quad (7.72)$$

with $Y^\epsilon(0) = X^\epsilon(0) = X_0$. Notice that $X^1(t) = X(t)$ and $Y^1(t) = Y(t)$, and that $X^0(t) = Y^0(t) = X_0$.

Lemma 7.6.7. *For the SDE (7.71), we have the formal expansion*

$$X^\epsilon(T) = X_0 + \epsilon A_X(T) + \frac{1}{2}\epsilon^2 B_X(T) + \frac{1}{6}\epsilon^3 C_X(T) + O(\epsilon^4),$$

where

$$\begin{aligned} A_X(T) &= X_0 \int_0^T \lambda(t) dW(t), \\ B_X(T) &= 2 \int_0^T \lambda(t)b(t)A_X(t) dW(t), \\ C_X(T) &= 3 \int_0^T \lambda(t)c(t)A_X(t)^2 dW(t) + 3 \int_0^T \lambda(t)b(t)B_X(t) dW(t). \end{aligned}$$

Proof. We rely on standard asymptotic expansion techniques (e.g. Yoshida [1992]) to construct a Taylor series of $X^\epsilon(T)$ around $\epsilon = 0$. Dropping the arguments of $\varphi(t) = \varphi(b(t), c(t), X^\epsilon(t))$ for brevity, we get

$$\begin{aligned} A_X(T) &= \left. \frac{\partial X^\epsilon(T)}{\partial \epsilon} \right|_{\epsilon=0} \\ &= \left. \left(\int_0^T \lambda(t)\varphi(t) dW(t) + \epsilon \int_0^T \lambda(t) \frac{\partial \varphi(t)}{\partial X^\epsilon(t)} \frac{\partial X^\epsilon(t)}{\partial \epsilon} dW(t) \right) \right|_{\epsilon=0} \\ &= \int_0^T \lambda(t)\varphi(b(t), c(t), X_0) dW(t) = X_0 \int_0^T \lambda(t) dW(t). \end{aligned}$$

Similarly,

$$\begin{aligned}
B_X(T) &= \frac{\partial^2 X^\epsilon(T)}{\partial \epsilon^2} \Big|_{\epsilon=0} \\
&= \left. \left(\int_0^T \lambda(t) \frac{\partial \varphi(t)}{\partial X^\epsilon(t)} \frac{\partial X^\epsilon(t)}{\partial \epsilon} dW(t) + \int_0^T \lambda(t) \frac{\partial \varphi(t)}{\partial X^\epsilon(t)} \frac{\partial X^\epsilon(t)}{\partial \epsilon} dW(t) \right) \right|_{\epsilon=0} \\
&\quad + \epsilon \int_0^T \lambda(t) \frac{\partial^2 \varphi(t)}{\partial X^\epsilon(t)^2} \left(\frac{\partial X^\epsilon(t)}{\partial \epsilon} \right)^2 dW(t) \Big|_{\epsilon=0} \\
&\quad + \epsilon \int_0^T \lambda(t) \frac{\partial \varphi(t)}{\partial X^\epsilon(t)} \frac{\partial^2 X^\epsilon(t)}{\partial \epsilon^2} dW(t) \Big|_{\epsilon=0} \\
&= 2 \int_0^T \lambda(t) b(t) A_X(t) dW(t),
\end{aligned}$$

where we have used the fact that $\partial \varphi(t)/\partial X^\epsilon(t) = b(t)$ when $X^\epsilon(t) = X_0$. The result for $C_X(T)$ follows in the same fashion. \square

For the variable Y^ϵ in (7.72), we get

$$Y^\epsilon(T) = X_0 + \epsilon A_X(T) + \frac{1}{2} \epsilon^2 B_Y(T) + \frac{1}{6} \epsilon^3 C_Y(T) + O(\epsilon^4),$$

where B_Y and C_Y are found by substituting \bar{b} for $b(t)$ and \bar{c} for $c(t)$ in the expressions for B_X and C_X in Lemma 7.6.7. We therefore immediately have the following result.

Lemma 7.6.8. *Consider the ϵ -indexed processes (7.71)–(7.72). Then, for $T > 0$,*

$$X^\epsilon(T) - Y^\epsilon(T) = \epsilon^2 I_1(\bar{b}; T) + \epsilon^3 I_2(\bar{b}, \bar{c}; T) + O(\epsilon^4),$$

where we have defined zero-mean random variables

$$\begin{aligned}
I_1(\bar{b}; T) &= \int_0^T \lambda(t) (b(t) - \bar{b}) A_X(t) dW(t), \\
I_2(\bar{b}, \bar{c}; T) &= \frac{1}{2} \int_0^T \lambda(t) (c(t) - \bar{c}) A_X(t)^2 dW(t) \\
&\quad + \frac{1}{2} \int_0^T \lambda(t) b(t) B_X(t) dW(t) - \frac{1}{2} \bar{b} \int_0^T \lambda(t) B_Y(t) dW(t).
\end{aligned}$$

There are numerous ways in which we can use the results of the previous section to determine the values of \bar{b} and \bar{c} that will make $Y^\epsilon(T)$ best approximate $X^\epsilon(T)$. Starting with \bar{b} , we here elect to set it such that the variance of the $O(\epsilon^3)$ term (the “skew term”) in Lemma 7.6.8 is minimized. That is, our optimal choice \bar{b}^* for \bar{b} is characterized by

$$\bar{b}^* = \operatorname{argmin}_{\bar{b}} E \left(I_1(\bar{b}; T)^2 \right). \quad (7.73)$$

Proposition 7.6.9. *The solution to (7.73) is*

$$\bar{b}^* = \int_0^T b(t) w_T(t) dt, \quad w_T(t) = \frac{\lambda(t)^2 v(t)^2}{\int_0^T \lambda(t)^2 v(t)^2 dt},$$

where $v(\cdot)^2$ is defined in (7.65).

Proof. First, we need to establish the expectation of the random variable $I_1(\bar{b}; T)^2$. From elementary properties of the Ito integral (see Theorem 1.1.3), we know that

$$\begin{aligned} \mathbb{E}(I_1(\bar{b}; T)^2) &= \mathbb{E}\left(\left(\int_0^T \lambda(t)(b(t) - \bar{b}) A_X(t) dW(t)\right)^2\right) \\ &= \int_0^T \lambda(t)^2 (b(t) - \bar{b})^2 \mathbb{E}(A_X(t)^2) dt. \end{aligned}$$

Since $A_X(t)$ is a Gaussian random variable with mean 0 and variance $X_0^2 v(t)^2$, it follows that

$$\mathbb{E}(I_1(\bar{b}; T)^2) = X_0^2 \int_0^T \lambda(t)^2 (b(t) - \bar{b})^2 v(t)^2 dt.$$

The (necessary) condition for a minimum is

$$\frac{1}{X_0^2} \frac{\partial \mathbb{E}(I_1(\bar{b}; T)^2)}{\partial \bar{b}} = 2\bar{b} \int_0^T \lambda(t)^2 v(t)^2 dt - 2 \int_0^T \lambda(t)^2 b(t) v(t)^2 dt = 0,$$

from which the result in Proposition 7.6.9 follows. \square

As advertised, the result of Proposition 7.6.9 is identical to that of Corollary 7.6.3.

It remains to find \bar{c} . We fundamentally wish to fix it such that the variance of the $O(\epsilon^4)$ term (the ‘convexity term’) in Lemma 7.6.8 is minimized, given $\bar{b} = \bar{b}^*$. When $\bar{b} = \bar{b}^*$, however, we can observe that

$$I_2(\bar{b}^*, \bar{c}; T)^2 \approx \frac{1}{2} \int_0^T \lambda(t)(c(t) - \bar{c}) A_X(t)^2 dW(t),$$

which suggests the simplified condition⁸

$$\bar{c}^* = \operatorname{argmin}_{\bar{c}} \mathbb{E}\left(\left(\frac{1}{2} \int_0^T \lambda(t)(c(t) - \bar{c}) A_X(t)^2 dW(t)\right)^2\right). \quad (7.74)$$

⁸More rigorous results can be found in Andersen and Hutchings [2010], but the accuracy of (7.74) is typically sufficient for applications.

Proposition 7.6.10. *The value \bar{c}^* that satisfies (7.74) is*

$$\bar{c}^* = \int_0^T c(t) q_T(t) dt, \quad q_T(t) \triangleq \frac{\lambda(t)^2 v(t)^4}{\int_0^T \lambda(t)^2 v(t)^4 dt},$$

where $v(\cdot)^2$ is defined in (7.65).

Proof. We note that

$$\begin{aligned} \mathbb{E} \left(\left(\frac{1}{2} \int_0^T \lambda(t) (c(t) - \bar{c}) A_X(t)^2 dW(t) \right)^2 \right) \\ = \frac{1}{4} \int_0^T \lambda(t)^2 (c(t) - \bar{c})^2 \mathbb{E}(A_X(t)^4) dt. \end{aligned}$$

From a standard property of Gaussian random variables, we have

$$\mathbb{E}(A_X(t)^4) = 3X_0^4 v(t)^4.$$

Applying this result, the (necessary) condition for a minimum is

$$\frac{1}{2}\bar{c} \int_0^T \lambda(t)^2 v(t)^4 dt - \frac{1}{2} \int_0^T \lambda(t)^2 c(t) v(t)^4 dt = 0.$$

The Proposition 7.6.10 follows. \square

Remark 7.6.11. For the special case where λ is constant, we have $v(t)^2 = \lambda^2 t$ and therefore

$$\bar{b}^* = \frac{2 \int_0^T b(t) t dt}{T^2}, \quad \bar{c}^* = \frac{3 \int_0^T c(t) t^2 dt}{T^3}.$$

Note that the contribution of the instantaneous convexity $c(t)$ to the effective local volatility convexity grows with t at a faster rate ($O(t^2)$) than the contribution of $b(t)$ to the effective local volatility skew ($O(t)$).

7.6.4 Numerical Example

A brief numerical example is now in order. To provide a simple setup in which we can test our averaging results, we consider a two-period case where

$$\lambda(t) = \begin{cases} \lambda_0, & t \in [0, T'], \\ \lambda', & t \in (T', T], \end{cases} \quad b(t) = \begin{cases} b_0, & t \in [0, T'], \\ b', & t \in (T', T], \end{cases} \quad c(t) = \begin{cases} 0, & t \in [0, T'], \\ c', & t \in (T', T]. \end{cases} \quad (7.75)$$

The advantage of this setup is that it allows for high precision call option pricing without the need for finite difference grids or Monte Carlo methods. In particular, by having $c(t) = 0$ for $t \in [0, T']$, it follows that

$$dX(t) = (b_0 X(t) + (1 - b_0) X_0) \lambda_0 dW(t), \quad t \in [0, T']$$

such that, from the fact that these dynamics are those of a simple displaced log-normal process,

$$X(T') = \begin{cases} b_0^{-1} (X_0 \exp(-\frac{1}{2} b_0^2 \lambda_0^2 T' + b_0 \lambda_0 W(T')) - (1 - b_0) X_0), & b_0 \neq 0 \\ X_0 + X_0 \lambda_0 W(T'), & b_0 = 0. \end{cases} \quad (7.76)$$

Let⁹ $C(t, x; K, T)$ be the time t price of a K -strike, T -maturity call option when $X(t) = x$. Clearly (assuming zero interest rates)

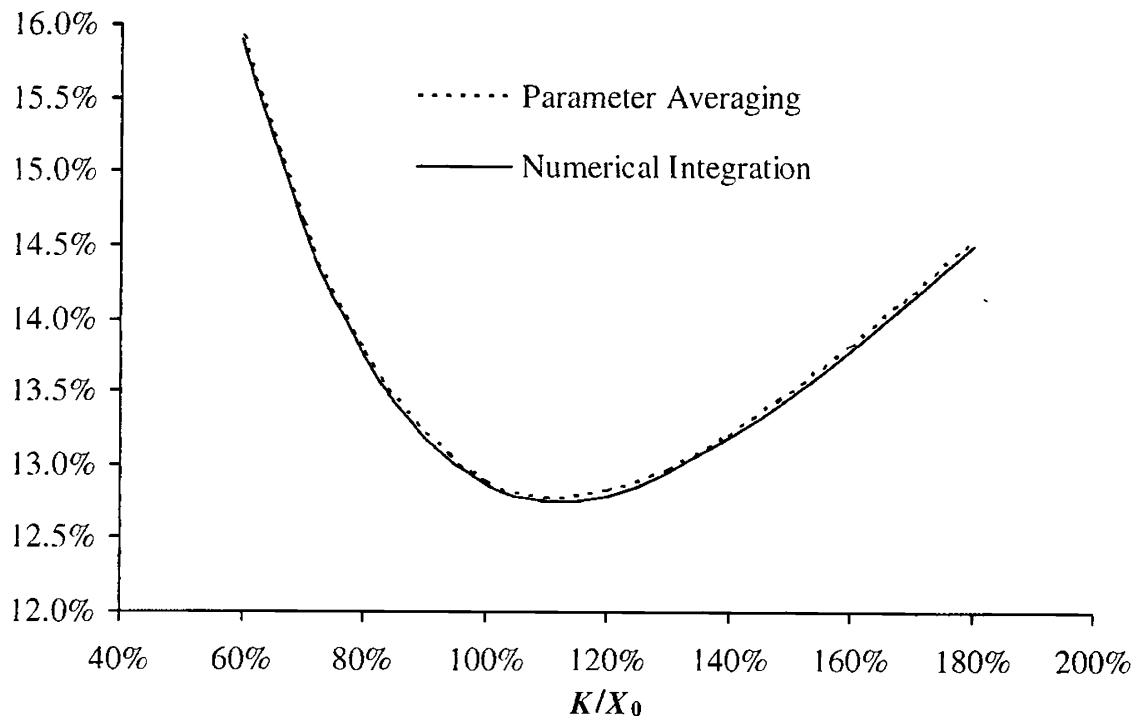
$$C(0, X_0; K, T) = \mathbb{E}(C(T', X(T'); K, T)).$$

At time T' , process parameters switch to constant values λ' , c' , b' so for any value of $X(T')$ computation of $C(T', X(T'); K, T)$ can be done using the formulas for call options in the quadratic model (see Section 7.3 and Andersen [2010]). From (7.76), computation of $C(0, X_0; K, T)$ can then easily be performed by numerical integration. Figure 7.2 below shows a sample fit for a high-convexity case ($X_0 = 1$, $c' = 4$).

The constant-parameter approximation here does an excellent job of matching the volatility smile of the true model. For even higher precision — especially for the (rare) case where convexity is very large and rapidly changing in time — additional correction terms may be required; see Andersen and Hutchings [2010] for the details and more numerical tests.

⁹We temporarily use notation C (rather than the usual c) for a call option, to distinguish it from the convexity function $c(t)$.

Fig. 7.2. Implied Volatility Smile



Notes: Parameters are as in (7.75), with $T' = 1$, $T = 2$, $\lambda_0 = 10\%$, $\lambda' = 15\%$, $b_0 = 0$, $b' = 0.75$, $c' = 4$. The x -axis denotes relative strike K/X_0 , with $X_0 = 1$. The “Numerical Integration” graph is the 2 year implied volatility smile for the time-dependent model, computed as outlined in the text (100 integration nodes). The “Parameter Averaging” graph computes the 2 year volatility smile from a constant-parameter quadratic model with parameters set as in Propositions 7.6.9 and 7.6.10.

Vanilla Models with Stochastic Volatility I

In Chapter 7 we introduced and studied diffusive single-factor vanilla models where the volatility is a deterministic function of the underlying rate. While such level-dependence of volatility is observable in interest rate markets — implied Black volatility is normally higher when rates are low — there is strong empirical evidence for additional sources of randomness in interest rate volatilities. To make our model setup more realistic, and to improve our ability to fit models to market-implied volatility smiles, we continue our investigation of vanilla models by enlarging the DVF models from the previous chapter to allow the volatility to be driven by a separate Brownian motion. The resulting models are said to have *stochastic volatility*.

Beyond raising the dimension of our models dynamics from one to two, the introduction of stochastic volatility brings with it a number of technical complications and, for many important models, the need to work with Fourier transforms when pricing options. We discuss these issues in detail in this chapter, paying particular attention to the *displaced log-normal Heston model* which has good analytical tractability and often provides an excellent fit to market observations.

Stochastic volatility constitutes a large and important topic in contemporary fixed income modeling, and we shall need two chapters of this book to lay the proper foundation for later work. In this chapter, our focus is on basic material and on the development of Fourier integration methods in a time-homogeneous setting. More advanced topics — including many numerical methods and the extension to time-dependent parameters — are postponed to Chapter 9.

8.1 Model Definition

As in Chapter 7, let $S(t)$, the “underlying” as we shall often call it, denote a forward Libor or swap rate. Also, let $Z(t)$, $W(t)$ be two different one-dimensional Brownian motions under a measure P in which $S(t)$ is a

martingale; we assume that $Z(t)$ and $W(t)$ are correlated with constant correlation ρ . As before, we use E instead of E^P for the expected value operator under measure P if there is no possibility of confusion. A fairly general family of stochastic volatility models¹ is obtained by specifying

$$dS(t) = \lambda\varphi(S(t)) \sqrt{z(t)} dW(t), \quad (8.1)$$

$$dz(t) = \theta(m(t) - z(t)) dt + \eta\psi(z(t)) dZ(t), \quad z(0) = z_0, \quad (8.2)$$

where λ, θ, η are positive constants, $m(\cdot)$ a positive deterministic function of time, and $\varphi(\cdot)$ and $\psi(\cdot)$ two smooth deterministic functions. In these SDEs, $z(\cdot)$ is a *stochastic variance* process, the square root of which scales a DVF diffusion term similar to that discussed in Chapter 7.

We notice that the drift term of $z(\cdot)$ is such that $z(t)$ gets pulled towards the level $m(t)$ at an exponential rate of θ , known as the *mean reversion speed* (or sometimes just the *mean reversion*). The parameter η is the *volatility of variance*, and $\psi(z)$ is a skew function for the stochastic variance. We shall later discuss in more detail the roles and effects of the individual parameters in the dynamics for $z(t)$, but before doing so let us try to indicate what constitutes a reasonable model specification. First, since the effect of $z(\cdot)$ on the volatility of $S(\cdot)$ is multiplicative, the initial value z_0 and the value $m(t)$ to which it mean-reverts can be scaled to arbitrary level; for convenience² we typically set $m(t) \equiv z_0 = 1$. As for the functions $\psi(\cdot)$ and $\varphi(\cdot)$, there are many empirically reasonable choices, but convenience and efficiency of available valuation algorithms for European options need to be considered. Typically, the function $\psi(\cdot)$ is chosen to be the square root function, making the process for $z(t)$ *affine* and improving analytical tractability. That said, other power functions, nevertheless, can be used and sometimes may be preferred, for reasons explained later (see, e.g., the end of Section 8.3). Analytical tractability also suggests using a linear function for $\varphi(\cdot)$, such that the underlying DVF model is a standard displaced log-normal model, see Section 7.2.4.

It only remains to comment on the correlation parameter ρ . In interest rate applications, the correlation ρ between the Brownian motions driving the stochastic variance and the underlying is often set to 0 due to undesirable effects of common measure changes on the stochastic variance process when correlation is non-zero, see Proposition 8.3.9. This is rarely a limitation,

¹For non-linear functions $\varphi(x)$ or $\varphi(t, x)$ such models are sometimes called *local stochastic volatility*, or LSV, models. Occasionally the name is also used for models with linear φ .

²Note that setting $m(\cdot)$ to a constant different from z_0 defines a model with constant coefficients that has a somewhat richer term volatility structure than with $m(\cdot) \equiv z_0$. The utility of this is limited as we are ultimately interested in time-dependent model extensions anyway.

as the effect of correlation on option prices and their implied volatilities³ can typically be captured in parameters of the function $\varphi(\cdot)$. From the perspective of matching the implied volatility smile, non-zero correlation is thus largely superfluous. Provided that we define our hedge sensitivities in a certain, natural way, this observation also holds for hedging, a point we shall return to in Section 8.9.2. To keep our discussion general, we nevertheless keep correlation non-zero for much of the discussion that follows.

With the parameter specializations described above, the simplified model we shall concentrate most of our efforts on is defined as

$$dS(t) = \lambda (bS(t) + (1 - b)L) \sqrt{z(t)} dW(t), \quad (8.3)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \eta \sqrt{z(t)} dZ(t), \quad z(0) = z_0 = 1, \quad (8.4)$$

with $\langle dZ(t), dW(t) \rangle = \rho dt$. Going forward, this model will be referred to as simply the *SV model*. For the case where $b = 1$, the model becomes identical to the so-called *Heston model*; see Heston [1993]. To avoid degenerate situations, we make the following assumption:

Assumption 8.1.1. All parameters b , θ , η , λ are strictly positive, and $|\rho| < 1$.

8.2 Model Parameters

We proceed to a more detailed discussion of the parameters in the model (8.3)–(8.4). First, recall that in the local volatility model of the displaced log-normal type (7.21), the parameter λ is responsible for the overall level of the implied volatility smile, while the parameter b is responsible for its slope. This interpretation of the parameters carries over to the stochastic volatility case (8.3)–(8.4), and we often refer to λ and b as the *SV volatility* and the *skew*, respectively.

The volatility of variance parameter η controls the curvature of the volatility smile, see Section 8.7. The effect of η on the volatility smile is similar to that of the second-order, or convexity, term in a quadratic DVF model of Section 7.3, although the dynamics of the volatility smile are quite different in the two models, a point we shall return to later, in Section 8.8.

The mean reversion of variance, θ , controls the speed at which deviations of $z(\cdot)$ away from z_0 are pulled back towards this level. Increasing θ decreases the long-term variance of $z(\cdot)$ and limits the effect of the stochastic variance process on the volatility smile for medium- and long-dated maturities. In essence, θ controls the speed of decay of the volatility smile convexity.

³If the correlation is negative — i.e. if $z(t)$ tends to be high when $S(t)$ is low — the model will imply a downward-sloping volatility smile, as should be intuitively clear.

The local volatility function $\varphi(x) = bx + (1 - b)L$ involves a quantity L , the *level* parameter. As discussed in the previous chapter, we normally set this to a number equal or close to⁴ $S(0)$, to ensure that λ will have the dimension of implied Black volatility, irrespective of the setting of b . This decoupling of parameters is particularly convenient in a calibration context.

As in the (local volatility) displaced log-normal model, λ is expressed in the units of relative volatility, while the skew b is typically confined to a range between 0% to 100%, although the “super-Normal” ($b < 0$) and “super-log-normal” ($b > 1$) settings may occasionally be useful. For $b < 0$ or $b > 1$, our earlier discussion in Section 7.6.2.2 shows that if $L > 0$, the state space for $S(\cdot)$ is bounded (above or below depending on the sign of b) by the value $-L(1 - b)/b$. The existence of such a bound is somewhat unrealistic; however, the advantages of being able to use values of b outside of $[0, 1]$ usually outweigh this concern.

The parameter η is expressed in the units of annualized relative volatility of *variance*. Sometimes it is more natural to think in terms of the volatility of *volatility*, i.e. the volatility of the process for $\sqrt{z(t)}$. By Ito’s lemma,

$$d\sqrt{z(t)} = O(dt) + \frac{\eta}{2} dZ(t).$$

When $z(t)$ has unit magnitude, $\eta/2$ can loosely be thought of as the volatility of volatility. For example, a value of 100% for η associates the implied Black volatility of the model with an instantaneous relative annualized volatility of about 50%. The related parameter θ , the speed of mean reversion, is expressed in percentage points per year. The inverse quantity θ^{-1} is measured in years and is related to the time over which a volatility shock dissipates. Specifically, the half-life of a volatility shock is

$$t_{1/2} = \frac{\ln 2}{\theta}.$$

All major interest rate markets exhibit high volatility of variance/low mean reversion of variance parameters, with $\eta = 150\%$ and $\theta = 10\%$ being typical parameter settings. While a half-life of a volatility shock of $10 \ln 2 \approx 7$ years may appear quite unrealistic, one should not forget that the pricing measure P will rarely represent real-world probabilities whereby the drift in the process for $z(\cdot)$ will likely contain strong market price of risk adjustment. The impact of measure changes on the speed of mean reversion for the variance is highlighted in Proposition 8.3.9.

⁴The rationale for not letting $L = S(0)$ always is that computation of delta sensitivities $\partial/\partial S(0)$ would then perturb the constant in the linear form $\varphi(x)$ which may or may not be desirable. See Sections 16.1.1 and 16.1.2 for more details.

8.3 Basic Properties

In this section we collect several important facts about the distribution and other relevant characteristics of $z(\cdot)$ and $S(\cdot)$ in the model (8.3)–(8.4). First, we look at the regularity properties of the process for the stochastic variance $z(\cdot)$; the results below should be compared to Proposition 7.2.1.

Proposition 8.3.1. *The SDE (8.4) has a unique solution. If $2z_0\theta \geq \eta^2$, i.e. the so-called Feller condition holds, $z = 0$ is unattainable. If the Feller condition is violated, $2z_0\theta < \eta^2$, then $z = 0$ is an attainable boundary but is strongly reflecting.*

Proof. See Revuz and Yor [1999] or Andersen and Piterbarg [2007]. \square

The transition distribution for the variance process $z(\cdot)$ given by (8.4) was derived in Cox et al. [1985] is listed below.

Proposition 8.3.2. *Let $\Upsilon(z; \nu, \gamma)$ be the cumulative distribution function for the non-central chi-square distribution with ν degrees of freedom and non-centrality parameter γ :*

$$\Upsilon(z; \nu, \gamma) = e^{-\gamma/2} \sum_{j=0}^{\infty} \frac{(\gamma/2)^j}{j! 2^{\nu/2+j} \Gamma(\nu/2 + j)} \int_0^z y^{\nu/2+j-1} e^{-y/2} dy. \quad (8.5)$$

For the process (8.4) define

$$d = 4\theta z_0 / \eta^2, \quad n(t, T) = \frac{4\theta e^{-\theta(T-t)}}{\eta^2 (1 - e^{-\theta(T-t)})}, \quad T > t. \quad (8.6)$$

Let $T > t$. Conditional on $z(t)$, $z(T)$ is distributed as $e^{-\theta(T-t)} / n(t, T)$ times a non-central chi-square distributed random variable with d degrees of freedom and non-centrality parameter $z(t)n(t, T)$,

$$P(z(T) < x | z(t)) = \Upsilon \left(\frac{x \cdot n(t, T)}{e^{-\theta(T-t)}}; d, z(t)n(t, T) \right).$$

Of particular importance, especially in Monte Carlo methods discussed later in Section 9.5, are the conditional moments of $z(\cdot)$. From the known properties of the non-central chi-square distribution, the following corollary easily follows⁵:

Corollary 8.3.3. *For $T > t$, $z(T)$ has the following first two conditional moments:*

$$E(z(T)|z(t)) = z_0 + (z(t) - z_0) e^{-\theta(T-t)}, \quad (8.7)$$

$$\text{Var}(z(T)|z(t)) = \frac{z(t)\eta^2 e^{-\theta(T-t)}}{\theta} \left(1 - e^{-\theta(T-t)} \right) + \frac{z_0\eta^2}{2\theta} \left(1 - e^{-\theta(T-t)} \right)^2. \quad (8.8)$$

⁵In Appendix A.A, p.1150, we also derive an expression for, and a numerical approximation to, $E(\sqrt{z(t)})$.

The transition distribution is useful for setting numerical bounds for PDE and Monte Carlo methods. Because it is somewhat complicated, we often find it convenient to use the stationary distribution of $z(t)$ (that is, the distribution of $z(\infty)$) instead, as an approximation.

Proposition 8.3.4. *The stationary distribution of $z(\cdot)$ in (8.4) is a Gamma distribution, see (3.9), and the stationary density $\pi(z)$ is given by*

$$\pi(z) = \frac{z^{\alpha-1} e^{-\beta z}}{\Gamma(\alpha) \beta^{-\alpha}},$$

where

$$\alpha = 2\theta z_0 / \eta^2, \quad \beta = 2\theta / \eta^2.$$

In particular, the mean of the stationary distribution is given by

$$\int_0^\infty z \pi(z) dz = \frac{\alpha}{\beta} = z_0,$$

and the variance by

$$\int_0^\infty (z - z_0)^2 \pi(z) dz = \frac{\alpha}{\beta^2} = \frac{z_0 \eta^2}{2\theta}.$$

Proof. Follows directly from Proposition 8.3.2 and Corollary 8.3.3, by taking the limit $T - t \rightarrow \infty$.

Now let us look at the properties of the process $S(\cdot)$ for the underlying. The martingale property for $S(\cdot)$ should not be taken for granted in stochastic volatility models, but fortunately holds in our case:

Proposition 8.3.5. *The process $S(\cdot)$ given by (8.3)–(8.4) is a proper martingale.*

Proof. See Andersen and Piterbarg [2007]. \square

The SDE (8.3) for $S(\cdot)$ can be integrated explicitly:

Proposition 8.3.6. *In the model (8.3)–(8.4), we have*

$$S(t) = \frac{1}{b} [(bS(0) + (1-b)L) X(t) - (1-b)L],$$

where

$$dX(t)/X(t) = \lambda b \sqrt{z(t)} dW(t), \quad X(0) = 1,$$

i.e.,

$$\ln X(t) = \lambda b \int_0^t \sqrt{z(s)} dW(s) - \frac{1}{2} \lambda^2 b^2 \int_0^t z(s) ds. \quad (8.9)$$

Proof. Follows by applying Ito's lemma to $\ln(bS(t) + (1 - b)L)$. \square

The moment-generating function of $\ln X(t)$ in (8.9) is of fundamental importance for European option pricing in the model (8.3)–(8.4), and is linked to the moment-generating function of the integrated variance process, as the following proposition demonstrates.

Proposition 8.3.7. *Define*

$$\Psi_X(u; t) \triangleq E\left(e^{u \ln X(t)}\right) = E(X(t)^u). \quad (8.10)$$

In the model (8.3)–(8.4), for any $u \in \mathbb{C}$ for which the right-hand side exists, we have

$$\Psi_X(u; t) = \Psi_{\bar{z}}\left(\frac{1}{2}(\lambda b)^2 u(u-1), u; t\right),$$

where we have denoted

$$\Psi_{\bar{z}}(v, u; t) \triangleq E^{\tilde{P}}\left(e^{v \bar{z}(t)}\right), \quad \bar{z}(t) \triangleq \int_0^t z(s) ds. \quad (8.11)$$

Under the new probability measure \tilde{P} the process for $z(\cdot)$ is

$$dz(t) = (\theta(z_0 - z(t)) + \rho\eta\lambda buz(t)) dt + \eta\sqrt{z(t)} d\tilde{Z}(t), \quad z(0) = z_0, \quad (8.12)$$

with \tilde{Z} a \tilde{P} -Brownian motion. If $\rho = 0$, then $\tilde{P} = P$ and $z(\cdot)$ in (8.11) follows (8.4) rather than (8.12).

Proof. From (8.9) we get

$$E\left(e^{u \ln X(t)}\right) = E\left(\varsigma(t) \exp\left(\frac{1}{2}u(u-1)\lambda^2 b^2 \int_0^t z(s) ds\right)\right),$$

where $\varsigma(t)$ is the exponential martingale

$$\begin{aligned} \varsigma(t) &= \mathcal{E}\left(u\lambda b \int_0^t \sqrt{z(s)} dW(s)\right) \\ &= \exp\left(u\lambda b \int_0^t \sqrt{z(s)} dW(s) - \frac{1}{2}u^2\lambda^2 b^2 \int_0^t z(s) ds\right). \end{aligned}$$

Letting $\varsigma(t)$ be the density process for a measure change, Proposition 8.3.7 follows from Girsanov's theorem, see Theorem 1.5.1. \square

A version of the proposition above also holds for a more general process (8.2) for $z(\cdot)$, see Andersen and Piterbarg [2007]. What makes the specification (8.4) particularly useful is the availability of a closed-form expression for $\Psi_{\bar{z}}(v, u; t)$.

Proposition 8.3.8. For $\Psi_{\bar{z}}(v, u; t)$ defined by (8.11) we have that

$$\begin{aligned}\ln \Psi_{\bar{z}}(v, u; t) &= A(v, u) + B(v, u) z_0, \\ A(v, u) &= \frac{\theta z_0}{\eta^2} \left[2 \ln \left(\frac{2\gamma}{\theta' + \gamma - e^{-\gamma T} (\theta' - \gamma)} \right) + (\theta' - \gamma) T \right], \\ B(v, u) &= \frac{2v(1 - e^{-\gamma T})}{(\theta' + \gamma)(1 - e^{-\gamma T}) + 2\gamma e^{-\gamma T}}, \\ \gamma &= \gamma(v, u) = \sqrt{(\theta')^2 - 2\eta^2 v}, \\ \theta' &= \theta'(u) = \theta - \rho\eta\lambda bu.\end{aligned}$$

Proof. The process (8.12) is of the form

$$dz(t) = (\theta z_0 - \theta' z(t)) dt + \eta \sqrt{z(t)} d\tilde{Z}(t), \quad \theta' = \theta - \rho\eta\lambda bu,$$

which is of the same form as the short rate process in Cox et al. [1985], see Section 10.2. As demonstrated by, e.g., Dufresne [2001], the discount bond pricing result from Cox et al. [1985] (derived via PDE methods) immediately establishes the moment-generating function of the time integral of $z(\cdot)$. \square

Beyond being useful in the proof of Proposition 8.3.7, measure changes are of primary importance in interest rate modeling, where a stochastic volatility model would typically be “embedded” in a full term structure model. To get a feel for issues that arise in this context, let us consider the impact of measure changes on the stochastic variance process. For this, let $V(t, X(t))$ be the numeraire-deflated price process for some asset in the model (8.3)–(8.4), where $V(t, x)$ is a deterministic function. Implicit in the notation is the assumption that the price does not depend on the stochastic variance process $z(\cdot)$, an assumption that holds true in the cases of interest to us. Assuming the price process is positive, it can be used as a numeraire, defining a new measure \tilde{P} , see Section 1.3. Since we have

$$dV(t, X(t))/V(t, X(t)) = \lambda bX(t) \frac{\partial \ln(V(t, X(t)))}{\partial x} \sqrt{z(t)} dW(t),$$

the process

$$\begin{aligned}\left(d\tilde{W}(t), d\tilde{Z}(t) \right)^T &= (dW(t), dZ(t))^T \\ &- \left(\lambda bX(t) \frac{\partial \ln(V(t, X(t)))}{\partial x} \sqrt{z(t)}, \rho \lambda bX(t) \frac{\partial \ln(V(t, X(t)))}{\partial x} \sqrt{z(t)} \right)^T dt\end{aligned}$$

is a two-dimensional Brownian motion under the measure \tilde{P} , see Theorem 1.5.1, and we obtain the following result.

Proposition 8.3.9. In the model (8.3)–(8.4), the dynamics of the stochastic variance process $z(\cdot)$ under a measure \tilde{P} defined by a numeraire $V(t, X(t))$ are given by

$$dz(t) = \tilde{\theta}(t, X(t)) \left(\tilde{f}(t, X(t)) - z(t) \right) dt + \eta \sqrt{z(t)} d\tilde{Z}(t),$$

where

$$\tilde{\theta}(t, x) = \theta - \eta \rho \lambda b x \frac{\partial \ln(V(t, x))}{\partial x}, \quad \tilde{f}(t, x) = \frac{\theta z_0}{\tilde{\theta}(t, x)},$$

and $\tilde{Z}(\cdot)$ is a \tilde{P} -Brownian motion.

We note that if $\rho \neq 0$, not only do the speed of mean reversion and the mean reversion level get altered by the measure change, they become dependent on the process for the underlying $S(\cdot)$ itself. As mentioned before, this makes it difficult, if not impossible, to relate statistically-observed stochastic variance parameters to the risk-neutral ones. Additionally, non-zero value of the correlation ρ introduces technical complications in interest rate modeling due to the heavy use of measure change machinery, complications that we normally avoid by setting ρ to 0.

Returning to the examination of the properties of the S -process, we note that while $S(\cdot)$ in (8.3)–(8.4) is always a martingale (see Proposition 8.3.5), some of its higher-order moments may become infinite with time. This has important implications in interest rate modeling where values of some common types of contracts require finite second-order moments, see Chapter 16 on convexity derivatives. The following proposition gives sharp conditions on moment existence.

Proposition 8.3.10. Consider the model (8.3)–(8.4). For a given $u > 1$, set $v = (\lambda b)^2 u(u - 1)/2 \geq 0$ and define

$$\beta = 2v/\eta^2 > 0, \quad \alpha = 2(\rho\eta\lambda bu - \theta)/\eta^2, \quad D = \alpha^2 - 4\beta.$$

The moment $E(S(T)^u)$ will be finite for $T < T^*$ and infinite for $T \geq T^*$, where T^* is given by

1. $D \geq 0, \alpha < 0$:

$$T^* = \infty;$$

2. $D \geq 0, \alpha > 0$:

$$T^* = \gamma_+^{-1} \eta^{-2} \ln \left(\frac{\alpha/2 + \gamma_+}{\alpha/2 - \gamma_+} \right), \quad \gamma_+ \triangleq \frac{1}{2} \sqrt{D};$$

3. $D < 0$:

$$T^* = 2\gamma_-^{-1} \eta^{-2} \times (\pi 1_{\{\alpha < 0\}} + \arctan(2\gamma_-/\alpha)), \quad \gamma_- \triangleq \frac{1}{2} \sqrt{-D}.$$

Proof. See Andersen and Piterbarg [2007]. \square

The problem of moment explosions in the SV model (8.3)–(8.4) can be resolved by replacing (8.4) with a slightly more general specification (8.2) with $\psi(z) = z^p$ for $p < 1/2$, at a cost of losing some analytical tractability.

There are a number of subtle but important issues related to stochastic volatility processes with $\psi(z) = z^p$; the reader is referred to Andersen and Piterbarg [2007] for a comprehensive discussion. While somewhat outside the main focus of our exposition, we list some relevant results in Appendix 8.A.

8.4 Fourier Integration

Having covered the basics, we now turn to the problem of establishing of accurate pricing methods for the SV model. The method we present here is based on the application of Fourier integration methods, and is largely taken from Lewis [2000], with some modifications. Carr and Madan [1999], Lipton [2002], and Lee [2004], among many others, can be consulted for additional details.

8.4.1 General Theory

The following general result shows how to calculate call option prices when a moment-generating function is available for the logarithm of the underlying.

Theorem 8.4.1. *Let ξ be a random variable, and define its moment-generating function by $\chi(u)$,*

$$\chi(u) = \mathbb{E}(e^{u\xi}).$$

Then for $k \in \mathbb{R}$,

$$\mathbb{E}((e^\xi - e^k)^+) = \chi(1) - \frac{e^k}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-k(\alpha+i\omega)} \chi(\alpha+i\omega)}{(\alpha+i\omega)(1-\alpha-i\omega)} d\omega \quad (8.13)$$

for any $0 < \alpha < 1$ for which the right-hand side exists.

Proof. Let

$$c(k) = \mathbb{E}((e^\xi - e^k)^+).$$

To improve regularity of our eventual numerical scheme, we split out a bounded component $\min(e^{\xi-k}, 1)$ from the unbounded function $(e^\xi - e^k)^+$, writing

$$\begin{aligned} c(k) &= \mathbb{E}(\max(e^\xi - e^k, 0)) \\ &= \mathbb{E}(e^\xi - e^k \min(e^{\xi-k}, 1)) \\ &= \chi(1) - e^k \mathbb{E}(\min(e^{\xi-k}, 1)). \end{aligned}$$

Our intention is now to apply Fourier transforms in the computation of $\mathbb{E}(\min(e^{\xi-k}, 1))$. While the function $\min(e^{\xi-k}, 1)$ is bounded by design, it

is not integrable — it equals 1 for all $x \geq k$. To work around this, we can follow Carr and Madan [1999] and write, with $p(x)$ being the density of ξ ,

$$\mathbb{E}(\min(e^{\xi-k}, 1)) = e^{-\alpha k} \int_{-\infty}^{\infty} [\min(e^{-(k-x)}, 1) e^{\alpha(k-x)}] [e^{\alpha x} p(x)] dx,$$

where $\alpha > 0$ is a classical *dampening constant*. Note that this integral is a convolution

$$(f_1 * f_2)(k) \triangleq \int_{-\infty}^{\infty} f_1(k-x) f_2(x) dx$$

of two functions,

$$f_1(x) = \min(e^{-x}, 1) e^{\alpha x}$$

and

$$f_2(x) = e^{\alpha x} p(x),$$

evaluated at k . Let \mathcal{F} be Fourier transform and \mathcal{F}^{-1} its inverse, i.e.,

$$(\mathcal{F}f)(\omega) \triangleq \int_{-\infty}^{\infty} e^{i\omega x} f(x) dx, \quad (8.14)$$

$$(\mathcal{F}^{-1}g)(x) \triangleq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-i\omega x} g(\omega) d\omega. \quad (8.15)$$

As is well known, the Fourier transform of a convolution is a product of Fourier transforms, so

$$\begin{aligned} & \int_{-\infty}^{\infty} [\min(e^{-(k-x)}, 1) e^{\alpha(k-x)}] [e^{\alpha x} p(x)] dx \\ &= (f_1 * f_2)(k) = (\mathcal{F}^{-1}(\mathcal{F}(f_1 * f_2)))(k) = (\mathcal{F}^{-1}(g_1(\omega) g_2(\omega)))(k), \end{aligned}$$

where

$$\begin{aligned} g_1(\omega) &= \int_{-\infty}^{\infty} e^{i\omega x} \min(e^{-x}, 1) e^{\alpha x} dx, \\ g_2(\omega) &= \int_{-\infty}^{\infty} e^{i\omega x} e^{\alpha x} p(x) dx. \end{aligned}$$

Simple calculations lead to

$$\begin{aligned} g_1(\omega) &= \int_{-\infty}^0 e^{x(\alpha+i\omega)} dx + \int_0^{\infty} e^{x(-1+\alpha+i\omega)} dx \\ &= \frac{1}{\alpha+i\omega} - \frac{1}{\alpha-1+i\omega} \\ &= \frac{1}{(\alpha+i\omega)(1-\alpha-i\omega)}, \\ g_2(\omega) &= \chi(\alpha+i\omega), \end{aligned}$$

where the convergence of integrals follows from the fact that $0 < \alpha < 1$. Therefore,

$$\mathbb{E}(\min(e^{\xi-k}, 1)) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-k(\alpha+i\omega)} \chi(\alpha+i\omega)}{(\alpha+i\omega)(1-\alpha-i\omega)} d\omega$$

and the theorem follows. \square

Remark 8.4.2. The formula (8.13) from Theorem 8.4.1 can be re-written as

$$\mathbb{E}((e^\xi - e^k)^+) = \chi(1) - \frac{e^k}{\pi} \int_0^{\infty} \operatorname{Re} \left(\frac{e^{-k(\alpha+i\omega)} \chi(\alpha+i\omega)}{(\alpha+i\omega)(1-\alpha-i\omega)} \right) d\omega,$$

a form that is used in, say, Attari [2004] and may yield computational benefits.

Proof. Let \bar{x} be the complex conjugate of x , $x \in \mathbb{C}$. If $H(\omega)$ is such that

$$H(-\omega) = \overline{H(\omega)}, \quad (8.16)$$

then

$$\begin{aligned} \int_{-\infty}^{\infty} H(\omega) d\omega &= \int_{-\infty}^0 H(\omega) d\omega + \int_0^{\infty} H(\omega) d\omega \\ &= \int_0^{\infty} H(-\omega) d\omega + \int_0^{\infty} H(\omega) d\omega \\ &= \overline{\int_0^{\infty} H(\omega) d\omega} + \int_0^{\infty} H(\omega) d\omega \\ &= 2\operatorname{Re} \left(\int_0^{\infty} H(\omega) d\omega \right). \end{aligned}$$

Since

$$\overline{\chi(\alpha+i\omega)} = \mathbb{E}(\overline{e^{(\alpha+i\omega)\xi}}) = \mathbb{E}(e^{(\alpha-i\omega)\xi}) = \chi(\alpha-i\omega),$$

the integrand in (8.13) satisfies (8.16) and the result follows. \square

A result complimentary to Theorem 8.4.1 holds for a call option on ξ rather than e^ξ .

Theorem 8.4.3. *In the notations of Theorem 8.4.1,*

$$\mathbb{E}((\xi - k)^+) = \left. \frac{d\chi(k)}{dk} \right|_{k=0} - k + \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-k(-\alpha+i\omega)} \chi(-\alpha+i\omega)}{(-\alpha+i\omega)^2} d\omega$$

for any $\alpha > 0$ for which the right-hand side exists.

Proof. As in the proof of Theorem 8.4.1, denote

$$c(k) = \mathbb{E}((\xi - k)^+).$$

While not strictly necessary, to keep the presentation consistent with the proof of Theorem 8.4.1, we manipulate this expression to obtain a bounded payoff inside the expected value,

$$\begin{aligned} c(k) &= \mathbb{E}(\max(\xi - k, 0)) \\ &= \mathbb{E}(\xi - \min(\xi, k)) \\ &= \chi'(0) - k - \mathbb{E}(\min(\xi - k, 0)), \end{aligned}$$

where $\chi'(\cdot)$ is the first-order derivative of the moment-generating function. Choosing $\alpha > 0$ and dampening the integrand with an exponential function, we obtain

$$\begin{aligned} \mathbb{E}((\xi - k)^+) &= \chi'(0) - k \\ &\quad - e^{\alpha k} \int_{-\infty}^{\infty} [\min(-(k-x), 0) e^{-\alpha(k-x)}] [e^{-\alpha x} p(x)] dx, \end{aligned}$$

where $p(x)$ is the density of ξ . By the same arguments as in the proof of Theorem 8.4.1,

$$\mathbb{E}((\xi - k)^+) = \chi'(0) - k - e^{\alpha k} (\mathcal{F}^{-1}(g_1(\omega) g_2(\omega))) (k),$$

where

$$\begin{aligned} g_1(\omega) &= \int_{-\infty}^{\infty} e^{i\omega x} \min(-x, 0) e^{-\alpha x} dx, \\ g_2(\omega) &= \int_{-\infty}^{\infty} e^{i\omega x} e^{-\alpha x} p(x) dx. \end{aligned}$$

Simple calculations lead us to

$$\begin{aligned} g_1(\omega) &= - \int_0^{\infty} x e^{x(-\alpha+i\omega)} dx = -\frac{1}{(-\alpha + i\omega)^2}, \\ g_2(\omega) &= \chi(-\alpha + i\omega), \end{aligned}$$

and the theorem follows. \square

8.4.2 Applications to SV Model

Combining Theorem 8.4.1 with the closed-form expression for the moment-generating function in the SV model (Propositions 8.3.6, 8.3.7, and 8.3.8), we obtain an efficient formula for pricing European call and put options in

the model (8.3)–(8.4). As suggested in Andersen and Andreasen [2002], its numerical properties can be enhanced by a type of control variate method where we add the Black formula and subtract its Fourier representation, reducing the discretization errors in the process. We present the call price result in this form.

Theorem 8.4.4. *The price of a call option $c_{\text{SV}}(0, S; T, K)$ in the SV model (8.3)–(8.4) is given by*

$$\begin{aligned} c_{\text{SV}}(0, S; T, K) &= \frac{1}{b} c_B(0, S'; T, K', \lambda b) \\ &\quad - \frac{K'}{2\pi b} \int_{-\infty}^{\infty} \frac{e^{(1/2+i\omega)\ln(S'/K')}}{\omega^2 + 1/4} q(1/2 + i\omega) d\omega, \end{aligned} \quad (8.17)$$

where $c_B(0, S'; T, K', \sigma)$ is the Black formula for spot S' , strike K' , expiry T and volatility σ , with

$$S' = bS + (1 - b)L, \quad K' = bK + (1 - b)L.$$

Also, we have defined

$$q(u) = \Psi_{\bar{z}} \left(\frac{1}{2} (\lambda b)^2 u(u-1), u; T \right) - e^{\frac{1}{2} \lambda^2 b^2 z_0 T u(u-1)}, \quad (8.18)$$

with $\Psi_{\bar{z}}$ given in Proposition 8.3.8.

Remark 8.4.5. In (8.17) we use volatility λb in the Black model. As a further refinement, one can use the ATM volatility implied by the SV model instead. The ATM volatility can, for instance, be approximated by an expansion approach, as explained in Sections 8.7 and 9.2.

Proof. From Proposition 8.3.6,

$$\begin{aligned} c_{\text{SV}}(0, S; T, K) &= E(S(T) - K)^+ \\ &= \frac{1}{b} E \left(S' e^{\ln X(T)} - K' \right)^+ \\ &= \frac{S'}{b} E \left(e^{\ln X(T)} - e^{\ln(K'/S')} \right)^+. \end{aligned}$$

By Theorem 8.4.1 and the definition (8.10) of $\Psi_X(u; t)$,

$$c_{\text{SV}}(0, S; T, K) = \frac{1}{b} \left(S' - \frac{K'}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-(\alpha+i\omega)\ln(K'/S')}\Psi_X(\alpha+i\omega; T)}{(\alpha+i\omega)(1-\alpha-i\omega)} d\omega \right), \quad (8.19)$$

where we have used the fact that $\Psi_X(1; T) = 1$. Applying this result to the SV model with $\eta = 0$, we find that the value of the option in the displaced log-normal model $c_{\text{DLN}}(0, S; T, K)$ is given by

$$c_{\text{DLN}}(0, S; T, K) = \frac{1}{b} \left(S' - \frac{K'}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-(\alpha+i\omega)\ln(K'/S')} \Psi_X^0(\alpha+i\omega; T)}{(\alpha+i\omega)(1-\alpha-i\omega)} d\omega \right), \quad (8.20)$$

where

$$\Psi_X^0(u; T) \triangleq \mathbb{E} \left(e^{u(\lambda b \sqrt{z_0} W(T) - \frac{1}{2} \lambda^2 b^2 z_0 T)} \right) = e^{\frac{1}{2} \lambda^2 b^2 z_0 T (u^2 - u)}. \quad (8.21)$$

On the other hand,

$$c_{\text{DLN}}(0, S; T, K) = \frac{1}{b} c_B(0, S'; T, K', \lambda b),$$

so that

$$\begin{aligned} & \frac{1}{b} c_B(0, S'; T, K', \lambda b) \\ & - \frac{1}{b} \left(S' - \frac{K'}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-(\alpha+i\omega)\ln(K'/S')} \Psi_X^0(\alpha+i\omega; T)}{(\alpha+i\omega)(1-\alpha-i\omega)} d\omega \right) = 0. \end{aligned}$$

Adding the left-hand side of this identity, which is zero, to the right-hand side of (8.19), we obtain

$$\begin{aligned} c_{\text{SV}}(0, S; T, K) &= \frac{1}{b} c_B(0, S'; T, K', \lambda b) \\ &- \frac{K'}{2\pi b} \int_{-\infty}^{\infty} \frac{e^{-(\alpha+i\omega)\ln(K'/S')} q(\alpha+i\omega)}{(\alpha+i\omega)(1-\alpha-i\omega)} d\omega, \end{aligned}$$

where

$$q(u) = \Psi_X(u; T) - \Psi_X^0(u; T).$$

Using Propositions 8.3.7 and 8.3.8 for $\Psi_X(u; T)$ and (8.21) for $\Psi_X^0(u; T)$, and setting $\alpha = 1/2$, the result follows. \square

Remark 8.4.6. The choice of $\alpha = 1/2$ in Theorem 8.4.4 is common in practice (see Lipton [2002]) and appears to give robust and stable results in most situations. As first pointed out by Lewis [2001], the value of α can be seen to define an integration contour in the complex plane, and values of α other than $1/2$ can be used as long as $\alpha + i\omega$ for all $\omega \in \mathbb{R}$ lie in the so-called *strip of convergence*⁶. One can attempt to optimize α to improve the numerical properties of the integral, see, e.g., Lee [2004] or Lord and Kahl [2007] for details. Moreover, integration contours are not restricted to straight lines. Lucic [2007] shows that all singularities of the function $q(u)$ are real (for our definition of q), paving the way for finding better — curvilinear — contours.

⁶The region of $u \in \mathbb{C}$ for which the moment-generating function $\chi(u)$ exists. Heston [1993] and Lewis [2000] establish the strip of convergence for the Heston model. The strip is directly related to moment existence, for the latter see Proposition 8.3.10.

Remark 8.4.7. Integrating complex values functions, such as $q(\alpha + i\omega)$, in a complex domain typically requires some care. Particularly troublesome are multi-valued functions such as the complex logarithm, as present in the expression for $\Psi_{\bar{z}}$ in Proposition 8.3.8. Should an integration contour cross a branching cut of such a function, the value will jump to a different branch, typically leading to wrong results. Fortunately the moment-generating function as presented in Proposition 8.3.8 is free of such problems. This is not the case for other, mathematically equivalent, expressions, such as, say, the one given in the original paper Heston [1993] — the reader is referred to Albrecher et al. [2007] for proofs and a detailed discussion of related issues.

Remark 8.4.8. By Assumption 8.1.1, Theorem 8.4.4 does not cover the case $b = 0$. If needed, this case can be handled by utilizing Theorem 8.4.3 instead of Theorem 8.4.1.

8.4.3 Numerical Implementation

The Fourier integral in (8.17) can be evaluated directly by any numerical integration scheme, in what is sometimes called the *direct integration approach*, see Kilin [2007]. With suitable restrictions on the integration technique and the integration grid spacing, one can formulate the pricing formula as a *discrete Fourier transform* (DFT), allowing for the usage of the *Fast Fourier Transform* (FFT) method, see Press et al. [1992]. The FFT method is developed in Section 8.4.5 below for applications requiring calculations of option prices for multiple strikes — such as volatility smile calibration or evaluation of European payoffs beyond simple puts and calls. The FFT method is certainly not competitive for calculating a *single* call option price, so here we focus on the direct integration method.

A direct numerical integration of (8.17) involves a scheme to discretize the integral and to handle the infinite integration domain. Many algorithms of varying degrees of sophistication have been proposed, some of which involve adaptive error control, optimal choice of dampening parameter α , and the mapping of the infinite integration domain on to a finite one. Lee [2004], Kilin [2007], Kahl and Jäckel [2005], Lord and Kahl [2007] contain sample algorithms, none of which employ the Black control variate inherent in our formulation (Theorem 8.4.4). As the control variate produces powerful error cancellations, we find that its inclusion allows for excellent results even when much simpler integration schemes are employed. We outline one such approach here.

Turning first to the integration bounds, we focus on the behavior of the integrand in (8.17) for large $|\omega|$; in fact, by Remark 8.4.2, only the limit $\omega \rightarrow +\infty$ needs to be explored. It turns out that the function $q(1/2 + i\omega)$ decays exponentially for large ω . In particular, as we can write

$$|q(1/2 + i\omega)| = e^{\operatorname{Re}(\ln(q(1/2 + i\omega)))},$$

we have the following result for $\ln(q(1/2 + i\omega))$.

Proposition 8.4.9. *Under our standard assumption that $|\rho| < 1$, for $q(\cdot)$ defined as in Theorem 8.4.4 we have*

$$\lim_{\omega \rightarrow +\infty} \frac{1}{\omega} \ln(q(1/2 + i\omega)) = -q_\infty,$$

where we have defined

$$q_\infty \triangleq \frac{\lambda b z_0}{\eta} \left(\sqrt{1 - \rho^2} + i\rho \right) (1 + \theta T). \quad (8.22)$$

Proof. The proof is obtained by applying simple calculus to formulas from Proposition 8.3.8; here we merely sketch it following the ideas of Kahl and Jäckel [2005]. We consider the limit of large positive ω . Let us denote

$$u(\omega) = 1/2 + i\omega, \quad v(\omega) = \frac{1}{2} (\lambda b)^2 u(\omega) (u(\omega) - 1) = -\frac{1}{2} (\lambda b)^2 (\omega^2 + 1/4).$$

Using the notations of Proposition 8.3.8, we have (we use “ \sim ” to denote equivalence in the limit $\omega \rightarrow +\infty$),

$$\theta'(u(\omega)) \sim -i\rho\eta\lambda b\omega, \quad \gamma(v(\omega), u(\omega)) \sim \rho^c\eta\lambda b\omega,$$

where

$$\rho^c \triangleq (1 - \rho^2)^{1/2}. \quad (8.23)$$

From the asymptotic behavior of $\gamma(\cdot, \cdot)$ it follows that the term $e^{-\gamma T}$ in the expressions for $A(\cdot, \cdot)$, $B(\cdot, \cdot)$ in Proposition 8.3.8 tends to zero as $\omega \rightarrow +\infty$. Therefore,

$$B(v(\omega), u(\omega)) \sim -\frac{\lambda b}{\eta} (\rho^c + i\rho) \omega,$$

and the logarithm in the definition of $A(\cdot, \cdot)$ tends to a constant,

$$\lim_{\omega \rightarrow +\infty} \ln \left(\frac{2\gamma}{\theta' + \gamma - e^{-\gamma T}(\theta' - \gamma)} \right) = \ln \left(\frac{2\rho^c}{\rho^c - i\rho} \right).$$

Therefore, only the term $(\theta' - \gamma)T$ in the expression for $A(\cdot, \cdot)$ grows with ω , and thus

$$A(v(\omega), u(\omega)) \sim -\frac{\lambda b z_0}{\eta} \theta(i\rho + \rho^c) T \omega.$$

Hence,

$$\begin{aligned} -\frac{1}{\omega} \ln(\Psi_X(1/2 + i\omega; T)) &= -\frac{1}{\omega} \ln(\Psi_{\bar{z}}(v(\omega), u(\omega); T)) \\ &= \frac{1}{\omega} (A(v(\omega), u(\omega)) + z_0 B(v(\omega), u(\omega))) \\ &\rightarrow \frac{\lambda b z_0}{\eta} (\rho^c + i\rho) (1 + \theta T) \end{aligned}$$

as $\omega \rightarrow +\infty$. Clearly, $\Psi_X^0(1/2 + i\omega; T)$ decays faster than that, as $e^{-\text{const} \times \omega^2}$, so $q(\cdot)$ inherits its tail behavior from $\Psi_X(\cdot; T)$, and the result follows. \square

The indefinite integral in Theorem 8.4.4 needs to be truncated before it can be evaluated numerically. Let $\omega_{\max} > 0$ be the upper truncation limit. We have the following simple tail estimate,

$$\begin{aligned} & \left| \int_{\omega_{\max}}^{\infty} \frac{e^{(1/2+i\omega) \ln(S'/K')}}{\omega^2 + 1/4} q(1/2 + i\omega) d\omega \right| \\ & \leq \int_{\omega_{\max}}^{\infty} \left| e^{(1/2+i\omega) \ln(S'/K')} \right| \frac{|q(1/2 + i\omega)|}{\omega^2} d\omega \\ & \leq \sqrt{\frac{S'}{K'}} e^{-\text{Re}(q_\infty)\omega_{\max}} \int_{\omega_{\max}}^{\infty} \frac{d\omega}{\omega^2} \\ & = \sqrt{\frac{S'}{K'}} \frac{e^{-\text{Re}(q_\infty)\omega_{\max}}}{\omega_{\max}}. \end{aligned}$$

If $\varepsilon_\omega > 0$ is the absolute tolerance for computing the option price via (8.17) (a value of $\varepsilon_\omega = 10^{-3}$ to 10^{-6} is a reasonable choice), then we set the upper truncation limit ω_{\max} by the condition

$$\frac{e^{-\text{Re}(q_\infty)\omega_{\max}}}{b\omega_{\max}} = \varepsilon_\omega, \quad (8.24)$$

where q_∞ is as given in Proposition 8.4.9. With Remark 8.4.2 in mind and a computational budget of N_ω points (N_ω is usually of the order of 100), we proceed to discretize uniformly over $[0, \omega_{\max}]$ and apply the rectangular rule

$$\begin{aligned} & \text{Re} \left(\int_0^{\infty} \frac{e^{(1/2+i\omega) \ln(S'/K')}}{\omega^2 + 1/4} q(1/2 + i\omega) d\omega \right) \\ & \approx \frac{\omega_{\max}}{N_\omega} \sum_{n=0}^{N_\omega-1} \frac{e^{\ln(S'/K')/2}}{\omega_n^2 + 1/4} \text{Re} \left(e^{i\omega_n \ln(S'/K')} q(1/2 + i\omega_n) \right), \quad (8.25) \end{aligned}$$

where

$$\omega_n = \omega_{\max} n / N_\omega, \quad n = 0, \dots, N_\omega - 1.$$

Other quadrature rules (e.g. the trapezoidal rule) can, of course, be used instead of the rectangular one.

8.4.4 Refinements of Numerical Implementation

While the method of Section 8.4.3 is simple and robust, numerical experiments show that the integration interval $[0, \omega_{\max}]$, with ω_{\max} obtained in (8.24), is often too wide, in the sense that a large proportion of the N_ω integration points are located in the region of integration where the integrand is so small that contributions to the integral are immaterial. To rectify

this, we can contemplate using an adaptive integration scheme, which by design would focus the computational work in regions where the integrand is material. Alternatively, we can refine our analysis of the integrand to provide guidance for where an ordinary integration scheme should spend its time. The latter is more involved but also more illuminating, so we pursue this approach here. Much of the material is based on Bang [2009], which can be consulted for additional details. As noted earlier, the ultimate benefit of sophisticated integration schemes (including the one proposed here) tends to be rather limited in practice, as long as the Black-Scholes control variate is properly employed.

We start by stating the following refinement of Proposition 8.4.9.

Proposition 8.4.10. *Let $q(\cdot)$ be defined as in Theorem 8.4.4 and assume, as always, that $|\rho| < 1$. Then for any $\epsilon > 0$ there exists $\Omega_\epsilon > 0$ such that, for any ω that satisfies*

$$\omega \geq \max \left(\Omega_\epsilon, \frac{5}{\eta \lambda b \rho^c T} \right),$$

we have

$$\frac{1 - \epsilon}{\omega^2} \leq \left| \ln(q(1/2 + i\omega)) - \left(-q_\infty \omega + q_0 - \frac{q_{-1}}{\omega} \right) \right| \leq \frac{1 + \epsilon}{\omega^2}, \quad (8.26)$$

where (compare to (8.22))

$$\begin{aligned} q_\infty &= \frac{\lambda b z_0}{\eta} (\rho^c + i\rho) (1 + \theta T), \\ q_0 &= \frac{z_0}{\rho^c \eta^2} (\rho^c + i\rho) \widehat{\theta}(1 + \theta T) + \frac{2\theta z_0}{\eta^2} \left(\ln(2\rho^c) + i \arctan\left(\frac{\rho}{\rho^c}\right) \right), \\ q_{-1} &= \frac{\theta z_0}{\eta^2} \left(T \mu \eta \lambda b + 2\widehat{\theta} \frac{\rho^c + i\rho}{(\rho^c)^2 \eta \lambda b} \right) + \mu \frac{\lambda b z_0}{\eta}, \\ \mu &= \frac{\widehat{\theta}^2}{2\eta^2 (\lambda b)^2 (\rho^c)^3} + \frac{1}{8\rho^c}. \end{aligned}$$

Here $\rho^c = (1 - \rho^2)^{1/2}$ is given by (8.23) and $\widehat{\theta} = \theta'(1/2)$, where $\theta'(u) = \theta - \rho \eta \lambda b u$ is defined in Proposition 8.3.8.

Proof. The proof is by expanding $\ln(q(1/2 + i\omega))$ into a series in $1/\omega$ for small values of $1/\omega$, along the lines of the proof of Proposition 8.4.9. Full details are available in Bang [2009]. \square

Let us denote

$$r(\omega) = \ln(q(1/2 + i\omega))$$

and by $r_\infty(\omega)$ its expansion to the zeroth order for large ω (see (8.26)),

$$r_\infty(\omega) = -q_\infty \omega + q_0.$$

Consider the integral on the left-hand side of (8.25), and let us split out a part that covers the region of (approximate) validity for the asymptotic approximation $\ln(q(1/2 + i\omega)) \approx r_\infty(\omega)$. To define this region, let us choose $\varepsilon'_\omega > 0$ reasonably small (of the order 10^{-2}) and pick $\omega'_{\max} > 0$ such that the following two conditions are simultaneously met:

$$\omega'_{\max} > \max \left(\frac{5}{\eta \lambda b \rho^c T} \right) \quad (8.27)$$

and, for any $\omega > \omega'_{\max}$,

$$\frac{|q_{-1}|}{\omega} \leq |r_\infty(\omega)| \varepsilon'_\omega. \quad (8.28)$$

Then, from Proposition 8.4.10,

$$\frac{|\ln(q(1/2 + i\omega)) - r_\infty(\omega)|}{|r_\infty(\omega)|} \approx \frac{|q_{-1}|}{\omega |r_\infty(\omega)|} \leq \varepsilon'_\omega$$

and, thus, for $\omega > \omega'_{\max}$, the function $\ln(q(1/2 + i\omega))$ is indeed well-approximated by $r_\infty(\omega)$. Accordingly, we write

$$\int_0^\infty \frac{e^{(1/2+i\omega) \ln(S'/K')} q(1/2 + i\omega)}{\omega^2 + 1/4} d\omega = I_1 + I_2 + I_3, \quad (8.29)$$

where

$$I_1 = \int_0^{\omega'_{\max}} \frac{e^{(1/2+i\omega) \ln(S'/K')} q(1/2 + i\omega)}{\omega^2 + 1/4} d\omega, \quad (8.30)$$

$$I_2 = \int_{\omega'_{\max}}^\infty \frac{e^{(1/2+i\omega) \ln(S'/K')}}{\omega^2 + 1/4} \left(q(1/2 + i\omega) - e^{r_\infty(\omega)} \right) d\omega, \quad (8.31)$$

$$I_3 = \int_{\omega'_{\max}}^\infty \frac{e^{(1/2+i\omega) \ln(S'/K')}}{\omega^2 + 1/4} e^{r_\infty(\omega)} d\omega. \quad (8.32)$$

As it turns out, the integral I_3 in (8.32) can be expressed through special functions. Let $E_1(z)$ be the so-called *exponential integral* (see Abramowitz and Stegun [1965]), i.e. an analytic continuation of the integral

$$E_1(z) = \int_1^{+\infty} \frac{e^{-zk}}{z} dk$$

to the complex plane. We then have the following result.

Lemma 8.4.11. *Let a and c be two non-negative real numbers and let z be a complex number such that $\operatorname{Re}(z) > 0$. Then*

$$\begin{aligned} R(z, a, c) &\triangleq \int_c^\infty \frac{e^{-zk}}{k^2 + a^2} dk \\ &= \frac{1}{2ia} (e^{-iaz} E_1(z(c - ia)) - e^{iaz} E_1(z(c + ia))). \end{aligned}$$

Proof. Follows by standard contour integration methods of complex analysis. Details are in Bang [2009]. \square

Remark 8.4.12. The function $E_1(\cdot)$ can be evaluated numerically using an algorithm from Press et al. [1992]. Bang [2009] also recommends an efficient algorithm available from <http://jin.ece.uiuc.edu>.

With the help of Lemma 8.4.11, we can rewrite I_3 in (8.32) as

$$\begin{aligned} \int_{\omega'_{\max}}^\infty \frac{e^{(1/2+i\omega)\ln(S'/K')}}{\omega^2 + 1/4} e^{r_\infty(\omega)} d\omega &= e^{q_0 + \ln(S'/K')/2} \\ &\quad \times R(q_\infty - i \ln(S'/K'), 1/2, \omega'_{\max}), \end{aligned}$$

and calculate it efficiently using Remark 8.4.12.

Turning next to the integral I_2 in (8.29), we wish to employ a quadrature rule designed to handle the oscillations of the integrand in (8.31). To that end, and following Bang [2009], we introduce a step size

$$\delta_\omega = \frac{(\lambda b z_0 \sqrt{T})^{-1}}{2N_{\text{stdev}}},$$

where N_{stdev} is a user-specified range in standard deviations⁷ (typically 5–6), set the number of points to be N''_ω (to be specified shortly), define

$$\omega''_n = \omega'_{\max} + \delta_\omega n, \quad n = 0, \dots, N''_\omega,$$

and write

$$\begin{aligned} I_2 &\approx e^{q_0} \sqrt{\frac{S'}{K'}} \int_{\omega'_{\max}}^{\omega'_{\max} + \delta_\omega N''_\omega} \frac{e^{\omega(-q_\infty + i \ln(S'/K'))}}{\omega^2 + 1/4} (e^{r(\omega) - r_\infty(\omega)} - 1) d\omega \\ &= e^{q_0} \sqrt{\frac{S'}{K'}} \sum_{n=0}^{N''_\omega - 1} \int_{\omega''_n}^{\omega''_{n+1}} \frac{e^{\omega(-q_\infty + i \ln(S'/K'))}}{\omega^2 + 1/4} (e^{r(\omega) - r_\infty(\omega)} - 1) d\omega, \end{aligned}$$

so that

⁷This step size in Fourier space is inspired by a Fourier transform of a Gaussian distribution. If the “width” of a Gaussian PDF is given by its standard deviation σ , then the “width” of its characteristic function is given by $1/\sigma$.

$$\begin{aligned}
I_2 &\approx e^{q_0} \sqrt{\frac{S'}{K'}} \sum_{n=0}^{N''_\omega - 1} \frac{e^{r(\omega''_n) - r_\infty(\omega''_n)} - 1}{(\omega''_n)^2 + 1/4} \int_{\omega''_n}^{\omega''_{n+1}} e^{\omega(-q_\infty + i \ln(S'/K'))} d\omega \\
&= e^{q_0} \sqrt{\frac{S'}{K'}} \sum_{n=0}^{N''_\omega - 1} \frac{e^{r(\omega''_n) - r_\infty(\omega''_n)} - 1}{(\omega''_n)^2 + 1/4} \\
&\quad \times \frac{e^{\omega''_{n+1}(-q_\infty + i \ln(S'/K'))} - e^{\omega''_n(-q_\infty + i \ln(S'/K'))}}{-q_\infty + i \ln(S'/K')}.
\end{aligned} \tag{8.33}$$

Note how we integrated analytically the oscillatory part of the integrand on the last step. With this scheme in place, we calculate I_2 using the quadrature rule (8.33) with N''_ω terms of the sum where we choose N''_ω adaptively by stopping when incremental changes from new terms in the sum are small enough.

Finally, let us discuss the term I_1 in (8.29), defined by (8.30). Here nothing special⁸ is needed and we can just use a quadratic or trapezoidal rule with a given budget of N'_ω points (say, around 50 or so) along the same lines as we did in (8.25).

In conclusion, let us summarize the complete algorithm for calculating the integral in (8.29). First we choose a small $\varepsilon'_\omega > 0$ (of the order 10^{-2}) and find the cutoff point ω'_{\max} that satisfies (8.27)–(8.28). Then we decompose the integral in (8.29) into three parts. The first integral I_1 is calculated by the standard quadratic or trapezoidal rule, similarly to (8.25). The second integral I_2 is calculated by the quadrature rule (8.33) with the number of points determined by the convergence criteria (relative or absolute). Finally the term I_3 is calculated per Remark 8.4.12. We note that while this scheme is more complex than what we described in Section 8.4.3, it does result in a faster and more accurate algorithm with a better utilization of the computational budget.

8.4.5 Fourier Integration for Arbitrary European Payoffs

Consider the problem of computing prices of European-style options with arbitrary payoffs. In particular, let $f(x)$ be a payoff function, and consider the problem of computing the following expected value,

$$\mathbb{E}(f(S(T))).$$

Clearly,

⁸Of the two terms in the definition of $q(1/2 + i\omega)$ in (8.18), the (second) one related to the Gaussian distribution decays much faster than the (first) one related to the SV model, as we already noted. Hence, we can stop sampling the second term for smaller values of ω , to save a bit on calculation time. This is described in Bang [2009].

$$\mathbb{E}(f(S(T))) = \int f(K) \mathbb{P}(S(T) \in dK)$$

and, by (7.5),

$$\mathbb{E}(f(S(T))) = \int f(K) \frac{\partial^2 c(0, S(0); T, K)}{\partial K^2} dK, \quad (8.34)$$

where $c(0, S; T, K)$ is the European call option value for the process $S(\cdot)$. Integrating by parts, we obtain a useful representation of a general European payoff in terms of European calls and puts.

Proposition 8.4.13. *For any twice-continuously differentiable⁹ $f(x)$, the value of a European option with payoff $f(\cdot)$ and expiry T is equal to the weighted integral of call and put options with weights equal to the second derivative of $f(\cdot)$,*

$$\begin{aligned} \mathbb{E}(f(S(T))) &= f(K^*) + f'(K^*)(S(0) - K^*) \\ &+ \int_{-\infty}^{K^*} p(0, S(0); T, K) f''(K) dK + \int_{K^*}^{\infty} c(0, S(0); T, K) f''(K) dK, \end{aligned} \quad (8.35)$$

for any K^* .

Proof. Follows by integration by parts of (8.34). \square

A combination of the suitably-discretized integral representation from Proposition 8.4.13 and Theorem 8.4.4 gives us an algorithm for computing values of European-style options with arbitrary payoffs. With the need to simultaneously compute call option prices of different strikes, the FFT method may deserve a closer look. In order to apply FFT, the discretization scheme of the integrals in (8.35) should be chosen carefully. From Theorem 8.4.4, the integrals to evaluate are

$$I(K') = \int_{-\infty}^{\infty} \frac{e^{(1/2+i\omega)\ln(S'/K')}}{\omega^2 + 1/4} q(1/2 + i\omega) d\omega \quad (8.36)$$

for various K' . We set $K^* = S(0)$ in (8.35) and discretize K in such a way that $\ln(S'/K')$ in (8.36) are equidistant. In particular, we choose $\delta > 0$, the discretization step, and define

$$x_n = \delta n, \quad K'_n = S' e^{x_n}, \quad n = -N, \dots, N.$$

This leads to

$$bK_n + (1-b)L = (bS + (1-b)L)e^{x_n},$$

or

⁹But see Section 16.6.1 for extensions.

$$K_n = \left(S + \frac{1-b}{b} L \right) e^{x_n} - \frac{1-b}{b} L.$$

Then

$$\begin{aligned} I_n &\triangleq I(K'_n) = \int_{-\infty}^{\infty} \frac{e^{-(1/2+i\omega)\delta n}}{\omega^2 + 1/4} q(1/2 + i\omega) d\omega = e^{-0.5\delta n} J_n, \\ J_n &\triangleq \int_{-\infty}^{\infty} e^{-i\omega\delta n} \frac{q(1/2 + i\omega)}{\omega^2 + 1/4} d\omega. \end{aligned}$$

At a computational effort of $O(N \ln N)$, all J_n 's can now be evaluated by applying (inverse) FFT to the function

$$\frac{q(1/2 + i\omega)}{\omega^2 + 1/4}.$$

Once the J_n are computed, all

$$p_{\text{SV}}(0, S; T, K_n), \quad c_{\text{SV}}(0, S; T, K_n), \quad n = -N, \dots, N,$$

can be calculated easily. The value of the option with any payoff $f(\cdot)$ is then obtained by discretizing the integrals in (8.35). We state the result as a proposition.

Proposition 8.4.14. Fix $\delta > 0$. Let K_n, K'_n , $n = -N, \dots, N$, be defined by

$$K_n = \left(S + \frac{1-b}{b} L \right) e^{\delta n} - \frac{1-b}{b} L, \quad K'_n = S' e^{\delta n}.$$

Then the value of a call option with payoff $f(\cdot)$ at time T in the SV model (8.3)–(8.4), is approximately given by

$$\begin{aligned} \mathbb{E}(f(S(T))) &\approx f(S(0)) + \sum_{n=-N}^{-1} p_{\text{SV}}(0, S; T, K_n) f''(K_n) (K_{n+1} - K_n) \\ &\quad + \sum_{n=0}^{N-1} c_{\text{SV}}(0, S; T, K_n) f''(K_n) (K_{n+1} - K_n), \end{aligned}$$

where

$$\begin{aligned} c_{\text{SV}}(0, S; T, K_n) &= \frac{1}{b} c_B(0, S'; T, K'_n, \lambda b) - \frac{K'_n}{2\pi b} e^{-0.5\delta n} J_n, \\ p_{\text{SV}}(0, S; T, K_n) &= \frac{1}{b} p_B(0, S'; T, K', \lambda b) - \frac{K'_n}{2\pi b} e^{-0.5\delta n} J_n, \end{aligned}$$

with $\{J_n\}_{n=-N}^N$ evaluated by an inverse FFT transform of the function

$$\frac{q(1/2 + i\omega)}{\omega^2 + 1/4},$$

and $q(u)$ given in Theorem 8.4.4.

Using FFT to compute the $2N + 1$ J_n -integrals improves numerical effort of a direct integration scheme, from $O(N^2)$ to $O(N \ln N)$. On the other hand, FFT has several potential drawbacks, including the fact that it imposes quite onerous requirements on the discretization of the strike domain, requiring that N be a power of 2 and that the grid be equidistant in $\ln(S'/K')$. Also, by the nature of FFT, an equidistant grid of the same size is then used to discretize the frequency domain. Both choices are often suboptimal — for example, we may want to choose a strike grid to take into account particular features of the payoff $f(\cdot)$, and we may want to discretize the frequency domain with a different number of grid points and/or non-equidistant spacing. In fact, Kilin [2007] observes that the integration effort is dominated by the calculation of the values of $q(1/2 + i\omega)$ for different ω and that they, critically, do not depend on strike. Kilin [2007] convincingly demonstrates that a careful implementation of the direct integration method of (8.17), even for multiple strikes, is often more efficient than FFT, provided that i) the values of $q(\cdot)$ are cached and reused when valuing different options, and ii) better discretization schemes are employed in the strike/frequency domains than those required by the FFT method.

8.5 Integration in Variance Domain

Under the assumption $\rho = 0$, a well-known “mixing” result (see e.g. Hull and White [1987]) represents the value of a European call option in the SV model (8.3)–(8.4) as an integral of the values of call options under the displaced log-normal model against the distribution of integrated variance. Specifically, the following lemma holds.

Lemma 8.5.1. *In the SV model (8.3)–(8.4) with $\rho = 0$, the value of a call option is given by*

$$c_{\text{SV}}(0, S; T, K) = \frac{1}{b} E \left(c_B \left(0, S; T, K, \lambda b \sqrt{\bar{z}(T)/T} \right) \right), \quad (8.37)$$

where (see (8.11))

$$\bar{z}(T) = \int_0^T z(t) dt$$

and $c_B(\cdot, \cdot; \cdot, \cdot, \sigma)$ is the value of a call option in the Black model with volatility σ .

Proof. Follows by conditioning on the trajectory of $z(\cdot)$ and using the independence of the Brownian motion $W(\cdot)$ of $z(\cdot)$. \square

Remark 8.5.2. An extension of this result to non-zero correlation ρ exists, see Proposition A.3.7 and in particular equation (A.39). Unfortunately it cannot be used for our purposes here, as the more general formula involves not only the time integral of $z(\cdot)$ but also other random variables.

It is natural to treat the function under the expected value operator in (8.37) as a function of $\bar{z}(T)$,

$$c_{\text{SV}}(0, S; T, K) = \mathbb{E}(C(\bar{z}(T))), \quad C(U) = \frac{1}{b} c_B \left(0, S; T, K, \lambda b \sqrt{U/T}\right). \quad (8.38)$$

As the moment-generating function $\Psi_{\bar{z}}(u, 0; T)$ of $\bar{z}(T)$ is known from Proposition 8.3.8, the expected value in (8.38) can be computed by Fourier integration. In particular, denoting by $p_{\bar{z}}(U)$ the probability density function of $\bar{z}(T)$, consider using (8.38) to argue that

$$\begin{aligned} c_{\text{SV}}(0, S; T, K) &= \int_0^\infty C(U)p_{\bar{z}}(U) dU \\ &= \frac{1}{2\pi} \int_0^\infty C(U) \int_{-\infty}^\infty e^{-i\omega U} \Psi_{\bar{z}}(i\omega, 0; T) d\omega dU \\ &= \frac{1}{2\pi} \int_{-\infty}^\infty \Psi_{\bar{z}}(i\omega, 0; T) \left(\int_0^\infty C(U)e^{-i\omega U} dU \right) d\omega \\ &= \frac{1}{2\pi} \int_{-\infty}^\infty \Psi_{\bar{z}}(i\omega, 0; T) (\mathcal{F}C)(-\omega) d\omega, \end{aligned}$$

where

$$(\mathcal{F}C)(\omega) \triangleq \int_0^\infty e^{i\omega U} C(U) dU \quad (8.39)$$

is the Fourier transform of $C(U)$ and we have used in the second equality the fact that $\Psi_{\bar{z}}$ is the Fourier transform of $p_{\bar{z}}$.

This argument demonstrates the main idea behind Fourier integration in the variance domain, but suffers from the fundamental problem that the function $C(\cdot)$ is not integrable, whereby the Fourier transform (8.39) is not well-defined. Fortunately we can solve the problem by the standard remedy of introducing a dampening function $e^{-\alpha U}$ in the integrand, as the following proposition demonstrates.

Proposition 8.5.3. *For $\alpha > 0$ such that $\Psi_{\bar{z}}(\alpha, 0; T)$ exists, the following holds,*

$$c_{\text{SV}}(0, S; T, K) = \frac{1}{2\pi} \int_{-\infty}^\infty \Psi_{\bar{z}}(\alpha + i\omega, 0; T) (\mathcal{F}\hat{C})(-\omega) d\omega,$$

where

$$\hat{C}(U) = C(U)e^{-\alpha U}, \quad (8.40)$$

and $\Psi_{\bar{z}}(u, 0; T)$ is given in Proposition 8.3.8.

Proof. We have

$$\begin{aligned}
c_{\text{SV}}(0, S; T, K) &= \int_0^\infty C(U) e^{-\alpha U} (e^{\alpha U} p_{\bar{z}}(U)) dU \\
&= \frac{1}{2\pi} \int_0^\infty C(U) e^{-\alpha U} \left(\int_{-\infty}^\infty e^{-i\omega U} \Psi_{\bar{z}}(\alpha + i\omega, 0; T) d\omega \right) dU \\
&= \frac{1}{2\pi} \int_{-\infty}^\infty \Psi_{\bar{z}}(\alpha + i\omega, 0; T) \left(\int_0^\infty C(U) e^{-\alpha U} e^{-i\omega U} dU \right) d\omega \\
&= \frac{1}{2\pi} \int_{-\infty}^\infty \Psi_{\bar{z}}(\alpha + i\omega, 0; T) (\mathcal{F}\hat{C})(-\omega) d\omega.
\end{aligned}$$

□

It is probably the case that the numerical method based on the result of Proposition 8.5.3 is not as speedy as the direct integration method in Section 8.4, but it allows for interesting generalizations to arbitrary payoff functions and arbitrary skew functions, a setup where it compares favorably to Monte Carlo or PDE methods. With this generalization in mind, consider the general model specification (8.1)–(8.2), where we have the following counterpart to Lemma 8.5.1.

Lemma 8.5.4. *For a positive constant v , let $g(t, S; v)$ satisfy the PDE*

$$\frac{\partial g(t, S; v)}{\partial S} + \frac{1}{2} v \varphi(S)^2 \frac{\partial^2 g(t, S; v)}{\partial S^2} = 0, \quad (8.41)$$

subject to the terminal boundary condition $g(T, S; v) = f(S)$. For the general stochastic volatility model dynamics (8.1)–(8.2) with $\rho = 0$ we have

$$\mathbb{E}(f(S(T))) = \mathbb{E}(g(0, S(0); T^{-1}\lambda^2 \bar{z}(T))). \quad (8.42)$$

Consistent with (8.39) and (8.40), we proceed to introduce a Fourier transform of a damped function g ,

$$(\mathcal{F}\hat{g})(\omega) = \int_{-\infty}^\infty e^{i\omega U} e^{-\alpha U} g(0, S(0); T^{-1}\lambda^2 U) dU, \quad (8.43)$$

where $\alpha > 0$ is as in Proposition 8.5.3. Then we have the following generalization of Proposition 8.5.3.

Proposition 8.5.5. *Consider the system (8.1)–(8.2), with $\psi(z) = \sqrt{z}$. Let $g(t, S; v)$ be as in (8.41) and $(\mathcal{F}\hat{g})$ as in (8.43) for $\alpha > 0$ such that $\Psi_{\bar{z}}(\alpha + i\omega, 0; T)$ is finite for all ω . Then*

$$\mathbb{E}(f(S(T))) = \frac{1}{2\pi} \int_{-\infty}^\infty (\mathcal{F}\hat{g})(-\omega) \Psi_{\bar{z}}(\alpha + i\omega, 0; T) d\omega,$$

where $\Psi_{\bar{z}}(u, 0; T)$ is given in Proposition 8.3.8.

The proposition gives us a way to compute values of arbitrary European options in a model with an essentially arbitrary volatility function $\varphi(\cdot)$. In calculating the integral in (8.43), we need a way to efficiently compute the function $g(0, S(0); v)$ from (8.41) for many different values of v . Fortunately, in Chapter 7 we developed many such methods, ranging from analytical expressions, to expansions and finite difference methods¹⁰. We note that if the function $\varphi(\cdot)$ is complicated enough to require finite difference methods, it is crucial that we use the “trick” of Section 7.4.1 to ensure that only a single finite difference grid is solved.

Remark 8.5.6. It can be verified that the moment-generating function $\Psi_{\bar{z}}(u, 0; T)$ is finite in a neighborhood around $u = 0$. Moments of arbitrary order of $\bar{z}(T)$ consequently exist and can be computed by differentiation

$$\mathbb{E}(\bar{z}(T)^n) = \left. \frac{d^n}{du^n} \Psi_{\bar{z}}(u, 0; T) \right|_{u=0}, \quad n = 1, 2, \dots.$$

Among other things, these moments can be used to dimension the U -grid used for the integration algorithm. For instance, for a given confidence multiplier γ (e.g. 5 or 10) we can, somewhat crudely, set

$$U_{\max} = \mathbb{E}(\bar{z}(T)) + \gamma \sqrt{\text{Var}(\bar{z}(T))}, \quad U_{\min} = \left(\mathbb{E}(\bar{z}(T)) - \gamma \sqrt{\text{Var}(\bar{z}(T))} \right)^+.$$

More elaborate schemes are also possible.

We note that Proposition 8.5.5 can also be applied to the case $\psi(z) = \sqrt{z - v}$, where $v > 0$ is a constant and where we enforce the additional constraint that $v < z_0$. To see this, consider the SDE

$$dz(t) = \theta(z_0 - z(t)) dt + \eta \sqrt{z(t) - v} dZ(t),$$

and set $z^*(t) = z(t) - v$. Then

$$dz^*(t) = \theta(z_0^* - z^*(t)) dt + \eta \sqrt{z^*(t)} dZ(t), \quad z_0^* = z_0 - v > 0, \quad (8.44)$$

and

$$\begin{aligned} \mathbb{E}\left(e^{u \int_0^T z(t) dt}\right) &= \mathbb{E}\left(e^{u \int_0^T (z^*(t) + v) dt}\right) \\ &= e^{uvT} \mathbb{E}\left(e^{u \int_0^T z^*(t) dt}\right) = e^{uvT} \Psi_{\bar{z}}(u, 0; T), \end{aligned}$$

where $\Psi_{\bar{z}}(u, 0; T)$ is computed as in Proposition 8.3.8 with the substitution $z_0 \rightarrow z_0 - v$. The form $\psi(x) = \sqrt{x - v}$ is useful if we wish to keep the process

¹⁰Many of the methods in Chapter 7 were specific to calls, for which the boundary condition on the PDE is $f(S(T)) = (S(T) - K)^+$. Not only is this case by far the most important in practice, but also helps with pricing of other payouts via the replication approach (Proposition 8.4.13).

$z(t)$ away from $z = 0$: it easily follows from (8.44) and $z(t) = z^*(t) + v$ that $z(t)$ will never go below v . According to Proposition 8.A.1, another way to keep $z(\cdot)$ away from the origin is to use $\psi(x) = x^p$, $1/2 < p < 1$. This case, however, has no analytical tractability.

For general $\psi(\cdot)$, let us consider ways to characterize the function $\Psi_{\bar{z}}(u, 0; T)$ that we now define by (8.11) for a general $z(\cdot)$ in (8.2). A useful starting point is the following result, easily proven from the Feynman-Kac formula in Section 1.8.

Lemma 8.5.7. *Let*

$$dz(t) = \theta(z_0 - z(t)) dt + \eta\psi(z(t)) dZ(t).$$

Then $\Psi_{\bar{z}}(u, 0; T) = L(0, z_0; u)$, where $L(t, z; u)$ satisfies the PDE

$$\frac{\partial L}{\partial t} + \theta(z_0 - z) \frac{\partial L}{\partial z} + \frac{1}{2}\eta^2\psi(z)^2 \frac{\partial^2 L}{\partial z^2} + uzL = 0,$$

subject to the boundary condition $L(T, z; u) = 1$.

Solution of the PDE in Lemma 8.5.7 can, of course, be done by finite difference methods, but at considerable numerical expense. An asymptotic expansion approach with decent precision is possible, however, and shall be demonstrated in Section 9.2 for the more general case of time-dependent λ . As it turns out, for many choices of $\psi(\cdot)$ — most notably for $\psi(z) = z^p$ — naively writing

$$\psi(z(t)) \approx \sqrt{z(t)}\psi(z(0)) / \sqrt{z(0)}$$

and then using the expression for $\Psi_{\bar{z}}(u, 0; T)$ from Proposition 8.3.8 often gives good results. Indeed, as shown in Andersen and Brotherton-Ratcliffe [2005], for call options, the dependence of option values on p in the specification $\psi(z) = z^p$ is quite mild across a reasonably wide range of strikes.

For complicated functions $\varphi(\cdot)$ and $\psi(\cdot)$ — and for the case where $\rho \neq 0$ — we always have the option of abandoning Fourier methods altogether and instead opting for more generally applicable numerical techniques, such as Monte Carlo and two-dimensional finite difference methods. We cover the application of these schemes to stochastic volatility models later on, in Sections 9.5 and 9.4, respectively.

8.6 CEV-Type Stochastic Volatility Models and SABR

As discussed earlier, certain choices of $\varphi(\cdot)$ and $\psi(\cdot)$ introduce technical problems, such as exploding higher-order moments of $S(\cdot)$, non-zero probability of generating negative $S(\cdot)$, or non-zero probability of the variance process $z(\cdot)$ being absorbed at zero. In practice, moment explosion is often the thorniest of these issues, as it has the potential to produce severe errors

for certain common securities (see Section 16.9). As it turns out, a simple switch from a linear function for $\varphi(\cdot)$ to a CEV-type specification prevents moment explosions that exist (Proposition 8.3.10) in the SV model. This is a useful result, so let us state it formally below. The proof is in Andersen and Piterbarg [2007].

Proposition 8.6.1. *Consider the model (8.1)–(8.2) with $\varphi(x) = x^c$ and $\psi(z) = z^p$, with $0 < c < 1$ and $p > 0$. Then for all $T \geq 0$ and $u \geq 0$,*

$$\mathbb{E}(S(T)^u) < \infty.$$

A particular CEV-type stochastic volatility model that has gained popularity with many practitioners is the so-called *SABR model*, see Hagan et al. [2002]. In Hagan et al. [2002], the SABR model is defined as

$$dS(t) = S(t)^c u(t) dW(t), \quad (8.45)$$

$$du(t) = \nu u(t) dZ(t), \quad (8.46)$$

with $\langle dW(t), dZ(t) \rangle = \rho dt$ and $0 < c < 1$. Note that the stochastic volatility $u(\cdot)$ is here modeled as simple geometric Brownian motion with zero drift. To translate the SDE (8.45)–(8.46) into more familiar terms, set $u(t) = \lambda \sqrt{z(t)}$, where $\lambda = u(0)/\sqrt{z_0}$. Then, with $\eta = 2\nu$,

$$\begin{aligned} dS(t) &= \lambda S(t)^c \sqrt{z(t)} dW(t), \\ dz(t) &= \frac{1}{4} \eta^2 z(t) dt + \eta z(t) dZ(t). \end{aligned}$$

We recognize this as a special case of our set-up (8.1)–(8.2) with $\psi(z) = z$, $m(t) = 0$, and *negative* mean reversion speed $\theta = -\eta^2/4$. The drift term in the process for $z(\cdot)$ is rather unattractive but allows for some tractability, as we shall see below. While higher-order moments can be very large in the SABR model, it follows from Proposition 8.6.1 that all positive moments of $S(t)$ exist (the fact that the mean reversion is negative can be shown to not influence the result in the proposition). Notice also that in the SABR model $S(\cdot)$ cannot go negative (although absorption at zero is a possibility) and that the variance process is strictly positive.

The main justification for the form of the equations (8.1)–(8.2) is that it allows for fairly accurate asymptotic expansions for European option prices. Hagan et al. [2002] obtained the first such expansion result by combining classical perturbation methods with, in the words of Obloj [2008], “impressive intuition”. Still, the result in Hagan et al. [2002] suffers from an internal inconsistency as $c \rightarrow 1$ and has later been revised by authors relying on more formal approaches. The result we list below is proven in Obloj [2008], based on earlier theoretical results in Berestycki et al. [2004] and Henry-Labordére [2005]. A similar result has been proven by Osajima [2007], using the small-noise expansion technique that we employed in Section 7.6.3.

Proposition 8.6.2. *For the model (8.45)–(8.46), the implied volatility smile is*

$$\sigma_B(t, S(t); K, T) = I^0 (1 + (T - t)I^1) + O((T - t)^2),$$

where

$$I^0 = \frac{-\nu \ln(K/S(t))}{\ln\left(\frac{\sqrt{1-2\rho q+q^2}+q-\rho}{1-\rho}\right)}, \quad q = \frac{\nu}{u(t)} \frac{S(t)^{1-c} - K^{1-c}}{1-c},$$

$$I^1 = \frac{(c-1)^2}{24} \frac{u(t)^2}{(S(t)K)^{1-c}} + \frac{1}{4} \frac{\rho\nu u(t)c}{(S(t)K)^{(1-c)/2}} + \frac{2-3\rho^2}{24} \nu^2.$$

Due to its lack of a mean reversion parameter, the SABR model often has difficulty matching smiles at different maturities when only a single set of calibration parameters $(\nu, c, \rho, u(0))$ is used. In practice, many financial institutions therefore maintain T -indexed vectors of these parameters, using the model primarily as a tool to interpolate and extrapolate the volatility smile. Some care must be exercised here, since the expansion listed above is not necessarily arbitrage-free; indeed, it is known that the expansion above may imply negative state price densities for low strikes and large maturities¹¹. These issues could potentially be rectified by ad-hoc methods for modifying the density, see Section 16.9 for an example.

8.7 Numerical Examples: Volatility Smile Statics

Having established a valuation formula for European options in the SV model, let us proceed to put it to work on some concrete model parameterizations. In doing so, we pay special attention to the way the various parameters of the SV model effect the implied volatility smile $\sigma_B(0, S(0); K, T)$. The results here provide additional color to the qualitative parameter discussion in Section 8.2. To aid our discussion, we start by listing a small- T expansion for the implied volatility of the SV model. The expansion is not particularly precise for medium and long-dated securities, but it suffices for the largely qualitative analysis in this section. As the expansion relies on techniques that we discuss in detail later (in Section 9.2) we skip the proof and also omit, for now, a precise characterization of the approximation convergence as $T \rightarrow 0$.

Lemma 8.7.1. *Define log-moneyness $\chi \triangleq \ln(K/S(0))$ and consider writing the implied Black volatility as*

¹¹Relative to the original SABR expansion in Hagan et al. [2002], the expansion in Proposition 8.6.2 is more robust in the low-strike tail; see Obloj [2008] for some numerical comparisons.

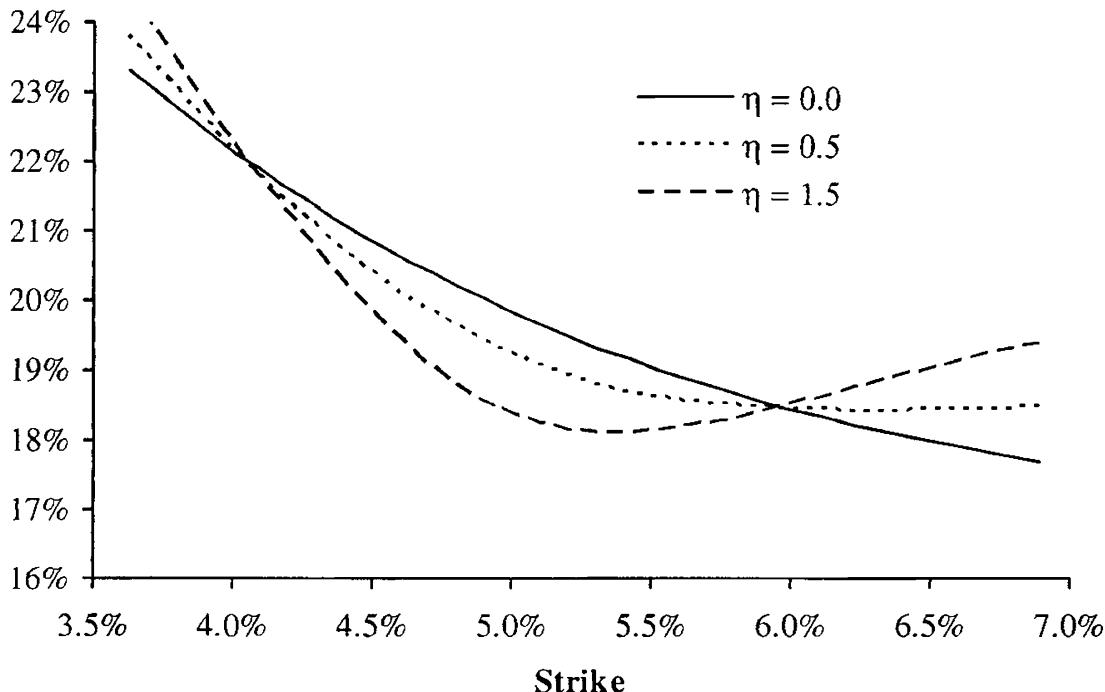
$$\sigma_B(0, S(0); T, K) = \sigma_{\text{ATM}} + R \cdot \chi + \frac{1}{2} B \cdot \chi^2 + \dots$$

for certain constants R and B . For small T and small χ , in the SV model (8.3)–(8.4) with $L = S(0)$ we have

$$\begin{aligned}\sigma_{\text{ATM}} &\approx \lambda, \quad R \approx \frac{\lambda}{2} \left(-(1 - b) + \frac{\eta\rho}{2\lambda} \right), \\ B &\approx \lambda \left(\frac{1 - b^2}{6} + \frac{\eta^2 (2 - 5\rho^2)}{24\lambda^2} \right).\end{aligned}$$

Armed with Lemma 8.7.1, we start out with an example of how the volatility of variance parameter η affects the convexity of the volatility smile. As discussed previously, η serves to generate convexity in the volatility smile, an effect that is obvious from the approximation for B in Lemma 8.7.1 and also clearly visible in Figure 8.1.

Fig. 8.1. 1 Year Volatility Smile

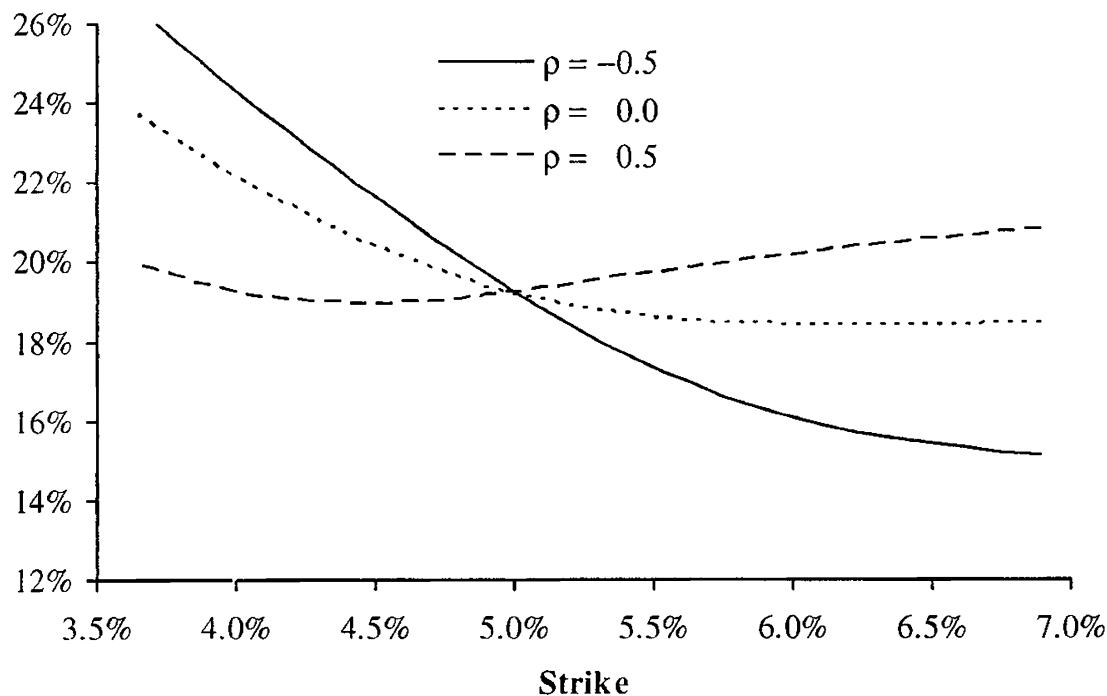


Notes: Implied volatility smile for SV model with $T = 1$, $S(0) = L = 5\%$, $z_0 = 1$, $b = 0.1$, $\lambda = 20\%$, $\theta = 0.1$, and $\rho = 0$. The volatility of variance parameter η varies as shown in the graph.

In Figure 8.1, the variance process is uncorrelated to the rate process, whereby Lemma 8.7.1 tells us that the slope (or skew) of the volatility smile at the at-the-money strike (5%) is generated solely by the slope parameter $b = 0.1$ in the local volatility function of the SV model. The stochastic volatility process can, of course, contribute to the skew if we use non-zero

correlation; see Figure 8.2 for a numerical example. As expected, lowering correlation rotates the smile clockwise, qualitatively similar to the impact of b . Another effect is also evident in Figure 8.2: when ρ moves away from zero, the convexity of the smile around the ATM strike is reduced. This effect is consistent with the expression for B in Lemma 8.7.1 which shows that the convexity (approximately) scales with¹² $2 - 5\rho^2$.

Fig. 8.2. 1 Year Volatility Smile



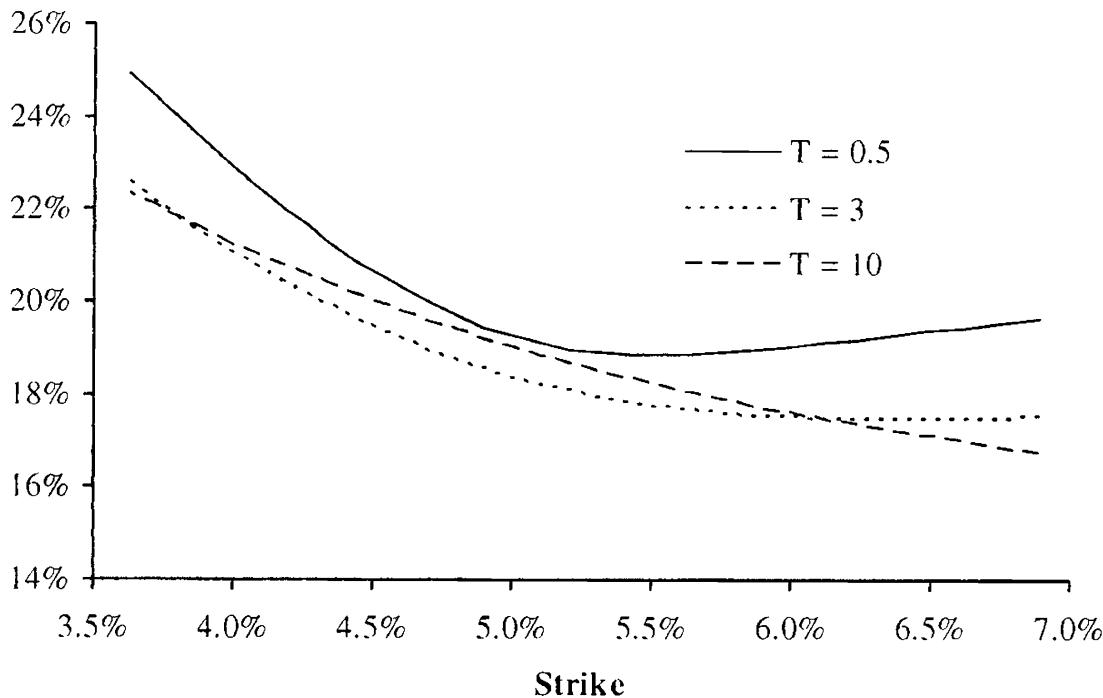
Notes: Implied volatility smile for SV model with $T = 1$, $S(0) = L = 5\%$, $z_0 = 1$, $b = 0.1$, $\lambda = 20\%$, $\theta = 0.1$, and $\eta = 1$. The correlation parameter ρ varies as shown in the graph.

The examples shown in Figures 8.1 and 8.2 both list the 1 year volatility smile only. To examine how the volatility smile $\sigma_B(0, S(0); K, T)$ in the SV model depends on T , consider first the case where $\rho = 0$; representative data are shown in Figure 8.3. The convexity of the smile, which originates with the stochastic volatility process, here clearly decays away as maturity is increased. As hinted at by Lemma 8.5.1, the convexity of the smile at time T is roughly proportional to the variance of the normalized realized variance $T^{-1} \int_0^T z(t) dt$. The convexity decay can therefore be interpreted as a mean reversion effect, since the variance of the normalized realized variance itself decays to a long-term (stationary) level, as can be seen from Corollary 8.3.3.

¹²Indeed, according to Lemma 8.7.1 the (short-maturity) smile convexity originating from stochastic volatility can become *negative* is $|\rho| > \sqrt{2/5} \approx 0.632$. This is easily verified numerically.

The speed of the decay is controlled by manipulating mean reversion speed θ ; the higher θ is, the quicker the smile convexity decays in the T -direction.

Fig. 8.3. Term Structure of Volatility Smiles



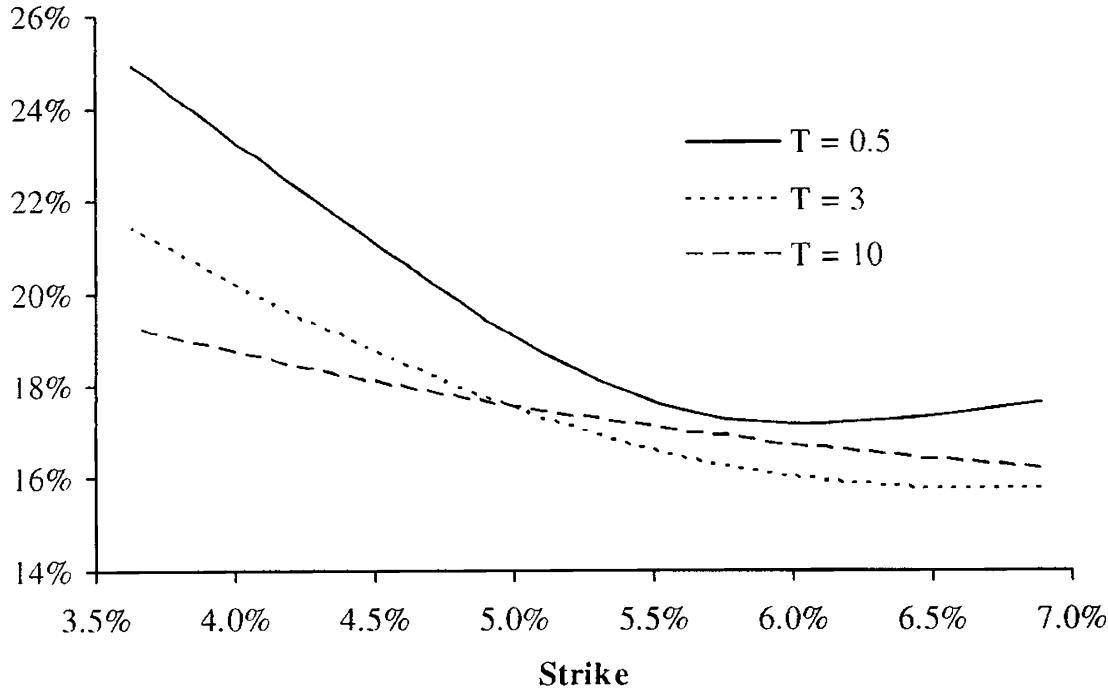
Notes: Implied volatility smile for SV model with $S(0) = L = 5\%$, $z_0 = 1$, $\rho = 0.0$, $\lambda = 20\%$, $\theta = 0.5$, $b = 0.1$, and $\eta = 1.5$. The smile maturity T varies as shown in the graph.

We note in passing that the ATM volatility of a constant parameter SV model is not a monotonic function of option maturity, as a quick glance at Figure 8.3 will confirm. For an analysis of the ATM volatility level and its dependence on maturity, see Lewis [2000].

In Figure 8.3 the slope of the smile around the ATM point is generated only from the parameter b in the local volatility function and consequently shows little decay in T . If, on the other hand, we had used a negative variance-spot correlation to generate the skew, we would expect the volatility smile to flatten out in T , for the same reason that the smile convexity decays. Figure 8.4 confirms this intuition.

8.8 Numerical Examples: Volatility Smile Dynamics

As we mentioned earlier, one rationale for introducing stochastic volatility into an LV model is the desire to generate realistic *smile dynamics*. In Section 7.1.3, we listed some qualitative reasons for the failure of LV models to generate reasonable model dynamics in certain cases; we are now in a

Fig. 8.4. Term Structure of Volatility Smiles

Notes: Implied volatility smile for SV model with $S(0) = 5\%$, $z_0 = 1$, $\rho = -0.5$, $\lambda = 20\%$, $\theta = 0.5$, $b = 1$, and $\eta = 1.5$. The smile maturity T varies as shown in the graph.

position to expand on this discussion and to show some concrete results. Specifically, we here wish to compare how the volatility smile moves with the underlying rate process, for two models: i) an ordinary (log-normal) Heston model obtained by setting $b = 1$ in the SV model (8.3)–(8.4); and ii) a pure LV model with quadratic volatility,

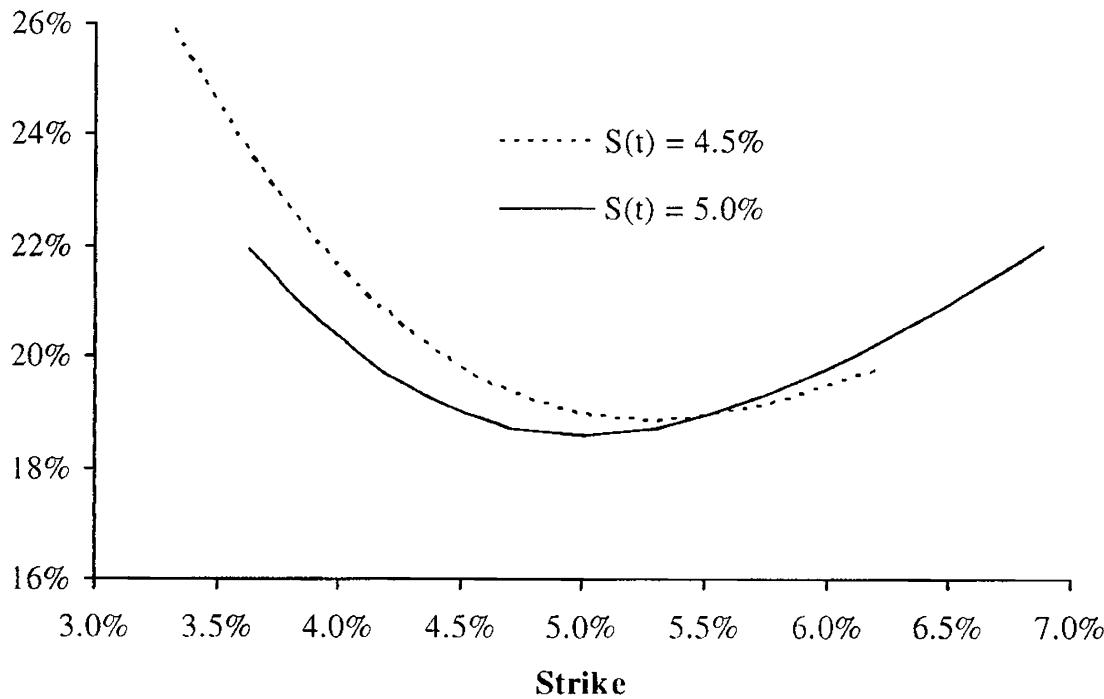
$$dS(t) = \lambda \left((1-b)L + bS(t) + \frac{1}{2}c(S(t) - L)^2 \right) dW(t). \quad (8.47)$$

For our numerical experiments, we move calendar time forward to some arbitrary value t and examine how the smile looks for several levels of $S(t)$. In performing this analysis for the Heston model, we shall initially assume that $z(t)$ stays equal to its initial value z_0 , but we relax this assumption later.

First, we consider the case of a (near) symmetric smile which in the local volatility model (8.47) can be obtained by setting $b = 1$. The effect of a 50 bps downward move in $S(0)$ (i.e. $S(t) = S(0) - 0.5\%$) on a specific LV model is shown in Figure 8.5. Starting from a symmetric smile when the forward rate $S(t) = S(0) = 5\%$, a shift down to 4.5% causes an overall increase in volatility levels, as well as a clock-wise tilt of the previously symmetric smile. This is readily understood, as the quadratic local volatility function

will itself increase and loose its symmetry when $S(t)$ is reduced from 5% to 4.5%.

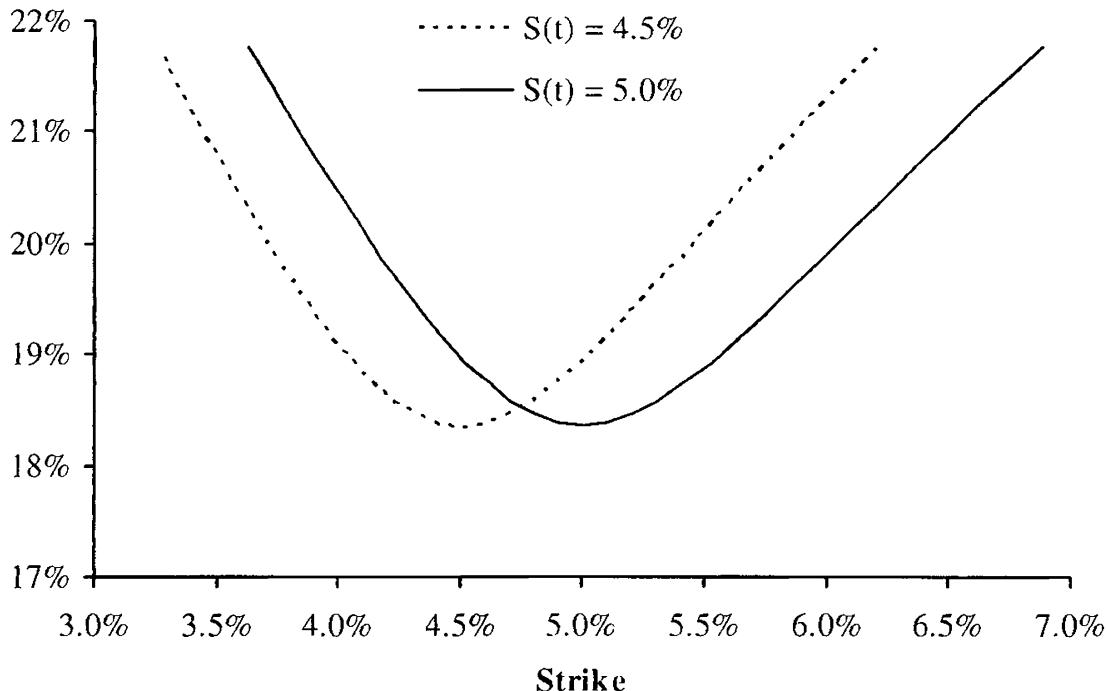
Fig. 8.5. Volatility Smile Dynamics in Quadratic LV Model



Notes: Time t implied volatility smile for quadratic LV model with $T = t + 1$, $S(0) = L = 5\%$, $b = 1$, $\lambda = 18\%$, and $c = 0.6$. Two different values for the forward rate $S(t)$ are used, as indicated in the graph.

Turning now to the Heston model, we first make the observation from Theorem 8.4.4 that European put and call option values normalized by spot S in both the Heston and Black models — and thereby the implied volatility smile of the Heston model — depend on strike K and forward rate $S(t)$ only through the ratio $K/S(t)$. Specifically, we have $\sigma_B(t, S(t); K, T) = g(K/S(t), T - t)$, for some function $g(\cdot, \cdot)$. In trader lingo, this is known as a “sticky delta” volatility smile¹³, and implies that the $T = t + \Delta$ volatility smile expressed in moneyness $K/S(t)$, or log-moneyness $\ln(K/S(t))$, is independent of t and $S(t)$, as long as $z(t)$ remains unchanged at its initial value z_0 . This fact makes it easy to construct the Heston model dynamics of the volatility smile in strike space; Figure 8.6 shows an example for a case where the correlation ρ has been set to zero to make the smile is symmetric in log-moneyness. Notice that as $S(t)$ drops from 5% to 4.5%, the volatility smile floats to the left, in tandem with the move in $S(t)$ such that the bottom of the smile remains centered at the forward rate.

¹³A reflection of the fact that the delta in the Black model, i.e. $\partial c_B / \partial S$, only depends on K/S .

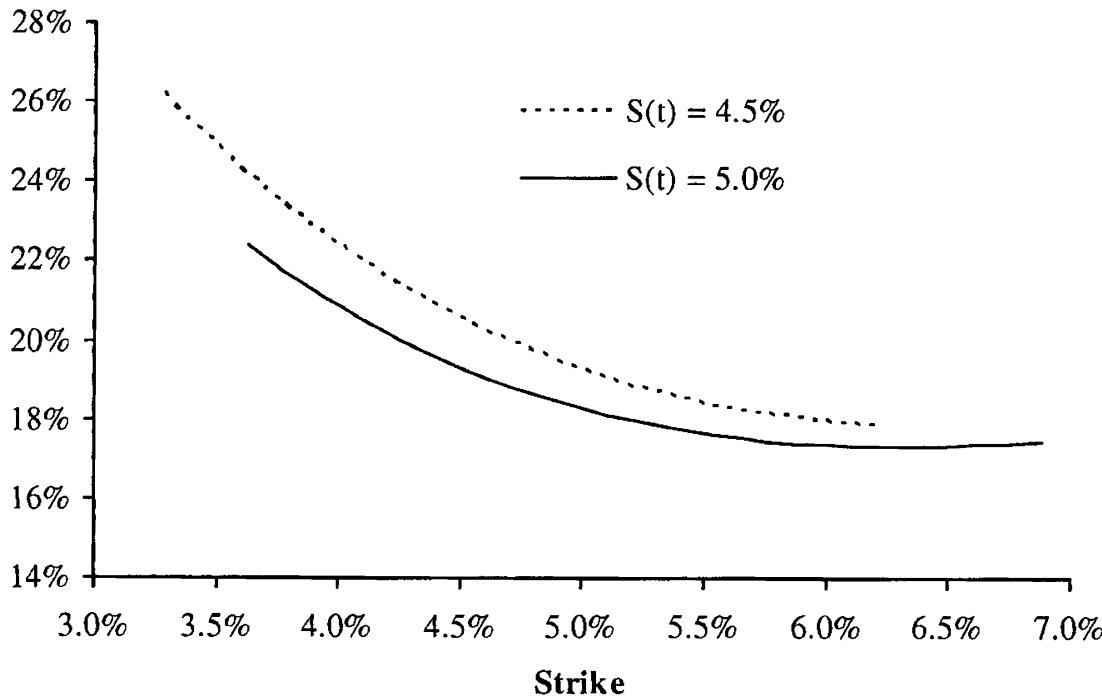
Fig. 8.6. Volatility Smile Dynamics in Heston SV Model

Notes: Time t implied volatility smile for SV model with $T = t + 1$, $S(0) = 5\%$, $z(t) = z_0 = 1$, $b = 1$, $\lambda = 20\%$, $\theta = 0.1$, $\rho = 0$, and $\eta = 1.5$. Two different values for the forward rate $S(t)$ are used, as indicated in the graph.

While Figures 8.5 and 8.6 are interesting and highlight some important differences between local and stochastic volatility models, it is more relevant in an interest rate setting to consider the case where the volatility smile has significant skew. First, we consider the local volatility case, see Figure 8.7. A shift down in $S(t)$ will increase the level of the local volatility function and raise the level of the smile; alternatively, we can interpret the move as a slide to the right. As convexity is relatively low in the graph relative to the skew, the move in $S(t)$ has little effect on the slope of the graph.

In Figure 8.8 we examine the smile dynamics of a Heston model with a significant downward skew, induced by a non-zero correlation ρ . The sticky-delta dynamics of the smile are still in effect here, causing a slide to the left when $S(t)$ is lowered, in a manner identical to that of the symmetric case in Figure 8.6.

The dynamics on display in Figures 8.7 and 8.8 appear to be diametrically opposite of each other: the smile shifts to the right in the local volatility model and to the left in the stochastic volatility model. In reality, however, differences in model dynamics are less dramatic than these graphs show. In particular, we recall that when we computed Figure 8.8, we kept $z(t)$ constant at the value z_0 . However, as $z(t)$ and $S(t)$ are negatively correlated in the model used in Figure 8.8, keeping one process constant while the other moves will clearly be wrong “on average”. A more representative

Fig. 8.7. Volatility Smile Dynamics in Quadratic LV Model

Notes: Time t implied volatility smile for quadratic LV model with $T = t + 1$, $S(0) = L = 5\%$, $b = 0.1$, $\lambda = 18\%$, and $c = 0.25$. Two different values for the forward rate $S(t)$ are used, as indicated in the graph.

characterization of the smile dynamics of the Heston process would move the variance process to its most likely outcome, given the move in the underlying. That is, we wish to set $z(t)$ equal to

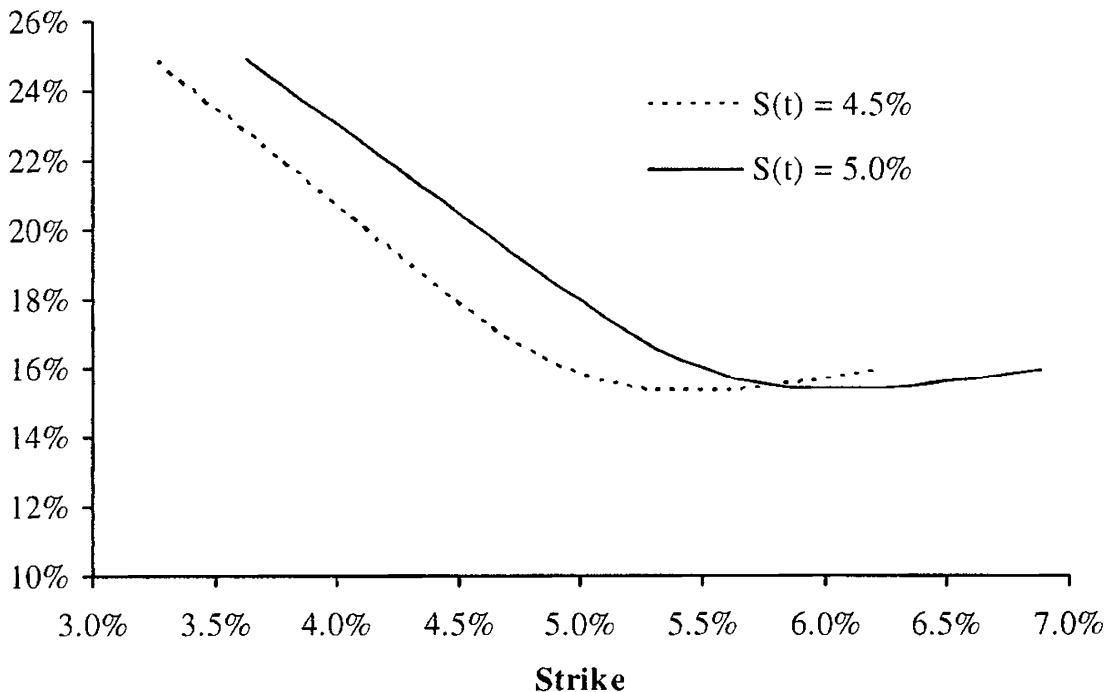
$$E(z(t)|S(t))$$

which we here compute by a simple Gaussian approximation that ignores mean reversion,

$$E(z(t)|S(t)) \approx z_0 + \frac{\eta\rho}{\lambda} \frac{S(t) - S(0)}{S(0)}. \quad (8.48)$$

Performing this modification on the data in Figure 8.8 results in the data in Figure 8.9.

With the rule in (8.48), the volatility smile shift of Figure 8.8 has reversed direction in Figure 8.9 and now looks quite similar to that of the local volatility dynamics of Figure 8.7. In other words, for volatility smiles that are “skew-dominated”, i.e. the skew is significant and the convexity is modest, smile dynamics of local and stochastic volatility models are quite similar on average. This observation is emphasized by Dupire [2006] and to some extent goes against common wisdom (see e.g. Hagan et al. [2002]) which tends to emphasize the sticky strike behavior of the stochastic volatility model. Of course, while the behavior in Figure 8.9 may be more likely

Fig. 8.8. Volatility Smile Dynamics in Heston SV Model

Notes: Time t implied volatility smile for SV model with $T = t + 1$, $z(t) = z_0 = 1$, $S(0) = 5\%$, $b = 1$, $\lambda = 20\%$, $\theta = 0.1$, $\rho = -0.6$, and $\eta = 1.5$. Two different values for the forward rate $S(t)$ are used, as indicated in the graph.

than that of Figure 8.8, both are feasible in a stochastic variance setting, depending on what value $z(t)$ happens to take. For derivatives that have convexity with respect to volatility smile moves¹⁴, what most reasonably represents “average” smile behavior is obviously less important than the fact that variance is random.

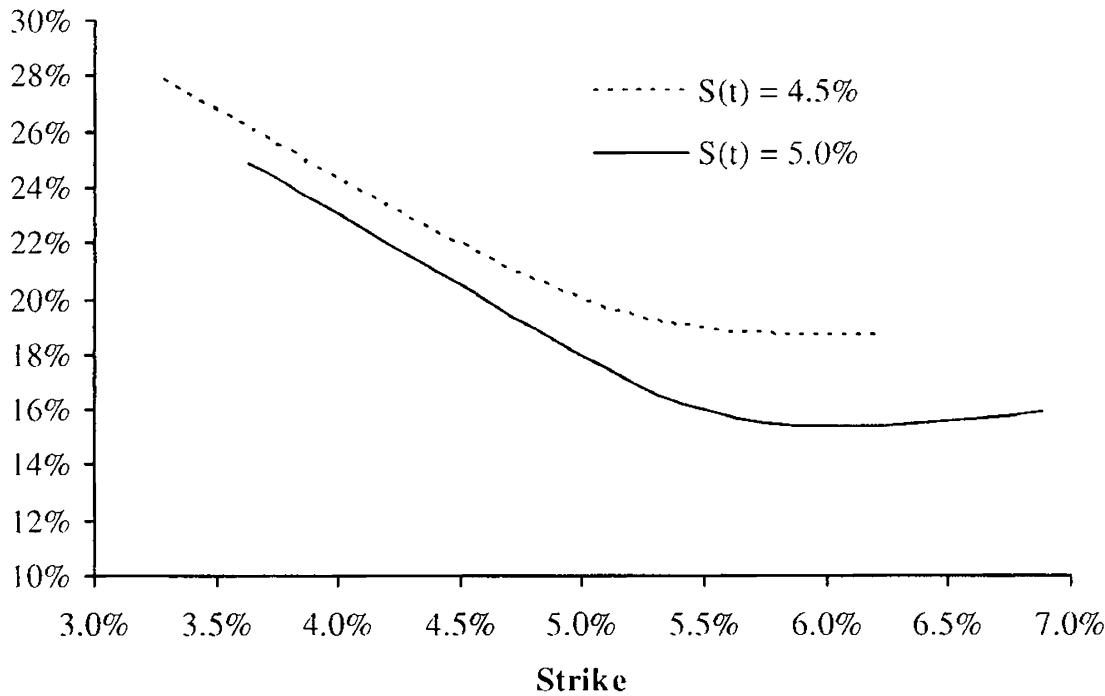
We finish this section by noting that the ideas behind (8.48) are also relevant for hedge construction in presence of stochastic volatility. We return to this topic in Section 8.9.2.

8.9 Hedging in Stochastic Volatility Models

8.9.1 Hedge Construction, Delta and Vega

Having now treated the subject of option pricing with stochastic volatility in quite some detail, let us make a foray into the topic of hedge construction. With their two generally non-collinear sources of randomness W and Z , it

¹⁴An option on implied volatility is an obvious example, although somewhat esoteric in an interest rate setting. A fairly common interest rate product with some volatility convexity is a barrier option. Many examples exist in other asset classes, such as reverse cliques and Napoleons, see Jeffery [2004].

Fig. 8.9. Volatility Smile Dynamics in Heston SV Model

Notes: Time t implied volatility smile for the SV model in Figure 8.8, but now with $z(t)$ set as computed from formula (8.48).

should be clear that stochastic volatility models of the type (8.1)–(8.2) are not complete (in the sense defined in Section 1.4) if we limit ourselves to simple delta hedging with positions only in $S(t)$ itself. However, if options with volatility sensitivity are available for trading, these can be included into the hedge portfolio to complete the market.

Assuming general dynamics (8.1)–(8.2), we proceed to consider hedging of a contingent claim $V(t)$ that depends on both $S(t)$ and $z(t)$, i.e. we write $V(t) = V(t, S(t), z(t))$. We assume existence of two traded securities $U_1(t) = U_1(t, S(t), z(t))$ and $U_2(t) = U_2(t, S(t), z(t))$. Using the framework of Section 1.7, we associate $U_1(t)$, $U_2(t)$ with the elements of the asset vector $X(t)$ from that section. Forming a hedging portfolio Π consisting of $-\pi_1(t)$ units of $U_1(t)$ and $-\pi_2(t)$ units of $U_2(t)$, we obtain from (1.26) that

$$\pi_i = \frac{\partial V}{\partial U_i}, \quad i = 1, 2.$$

A bit of calculus leads us to expressions for the hedge ratios in terms of sensitivities to the primitives S , z of the model, and the following result follows.

Lemma 8.9.1. *The portfolio $\Pi(t) = V(t) - \pi_1(t)U_1(t) - \pi_2(t)U_2(t)$ is locally riskless if*

$$\pi_1 = \left(\frac{\partial V}{\partial S} \frac{\partial U_2}{\partial z} - \frac{\partial U_2}{\partial S} \frac{\partial V}{\partial z} \right) \left(\frac{\partial U_1}{\partial S} \frac{\partial U_2}{\partial z} - \frac{\partial U_2}{\partial S} \frac{\partial U_1}{\partial z} \right)^{-1}, \quad (8.49)$$

$$\pi_2 = \left(\frac{\partial V}{\partial S} \frac{\partial U_1}{\partial z} - \frac{\partial U_1}{\partial S} \frac{\partial V}{\partial z} \right) \left(\frac{\partial U_2}{\partial S} \frac{\partial U_1}{\partial z} - \frac{\partial U_1}{\partial S} \frac{\partial U_2}{\partial z} \right)^{-1}. \quad (8.50)$$

Remark 8.9.2. In practice, the first security U_1 would often be chosen to not depend on z — for example the swap from which $S(t)$ is computed could be used as U_1 — in which case the hedge weights simplify. In particular,

$$\pi_1 = \frac{\partial V / \partial S}{\partial U_1 / \partial S} - \frac{\partial V / \partial z}{\partial U_2 / \partial z} \frac{\partial U_2 / \partial S}{\partial U_1 / \partial S}, \quad \pi_2 = \frac{\partial V / \partial z}{\partial U_2 / \partial z},$$

as one would expect.

Remark 8.9.3. The sensitivity of a given security to volatility is often called its *vega*. Even for a model with non-stochastic volatility, such as the Black model, a vega can be computed, but will not enter the hedge balance equation (1.28). In a stochastic volatility model, a vega can conveniently¹⁵ be defined to be $\partial/\partial z$ — which *will* enter the hedge balance equation. It follows that the choice (8.49)–(8.50) ensures that the hedged portfolio Π is delta-neutral, in the sense that

$$\frac{\partial \Pi(t)}{\partial S} = 0,$$

as well as *vega-neutral*,

$$\frac{\partial \Pi(t)}{\partial z} = 0.$$

8.9.2 Minimum Variance Delta Hedging

While the theoretical notion of “delta” assumes that the stochastic variance process z is kept fixed under perturbations of S , we saw earlier in Section 8.8 (see, in particular, Figure 8.9 and the discussion around it) that it sometimes might be more natural to let z float along with S , in a manner determined by the correlation between these quantities. Indeed, to the extent that our hedging strategy were to employ a position in S only, and not to separately hedge the exposure to z , the “best” hedging strategy — in the sense of locally minimizing hedging errors — is one based on such a joint move in z and S . We proceed to present this idea, using rather ad-hoc (or “deceptively simple”, to paraphrase Ewald et al. [2007]) techniques; for a full account and for a connection to the concept of the *minimal martingale measure*, see Follmer and Schweizer [1990] and Ewald et al. [2007].

First, let us return to the model (8.1)–(8.2), but now use a Cholesky decomposition to rewrite the process for $z(t)$ as

¹⁵From a theoretical viewpoint. More practical definitions of vega are covered later in the book, see Chapter 26 in particular.

$$dz(t) = O(dt) + \sigma_z(t) \left(\rho dW(t) + \sqrt{1 - \rho^2} dB(t) \right),$$

where B is a Brownian motion that is *independent* of W , and we use $\sigma_z(t) = \eta\psi(z(t))$ and $\sigma_S(t) = \lambda\varphi(S(t))\sqrt{z(t)}$ for notational clarity. Consider now a claim

$$V(t) = V(t, S(t), z(t)),$$

where, by Ito's lemma,

$$dV(t) = O(dt) + \frac{\partial V(t)}{\partial S} dS(t) + \frac{\partial V(t)}{\partial z} \sigma_z(t) \left(\rho dW(t) + \sqrt{1 - \rho^2} dB(t) \right).$$

Let us form a portfolio Π of the claim V and a position of $-\pi(t)$ in $S(t)$; that is,

$$d\Pi(t) = -\pi(t) dS(t) + dV(t). \quad (8.51)$$

We wish to set $\pi(t)$ such that $\text{Var}_t(d\Pi(t))$ is minimized.

Lemma 8.9.4. *With $d\Pi(t)$ defined in (8.51), the variance $\text{Var}_t(d\Pi(t))$ is minimized by setting $\pi(t) = \pi_{\text{mv}}(t)$, where*

$$\pi_{\text{mv}}(t) = \frac{\partial V(t)}{\partial S} + w(t), \quad w(t) = \frac{\partial V(t)}{\partial z} \frac{\rho\sigma_z(t)}{\sigma_S(t)},$$

and $\sigma_z(t) = \eta\psi(z(t))$, $\sigma_S(t) = \lambda\varphi(S(t))\sqrt{z(t)}$.

Proof. It is easily seen that

$$\begin{aligned} \text{Var}_t(d\Pi(t)) &= \left(-\pi(t)\sigma_S(t) + \frac{\partial V(t)}{\partial S}\sigma_S(t) + \frac{\partial V(t)}{\partial z}\sigma_z(t)\rho \right)^2 dt \\ &\quad + \left(\frac{\partial V(t)}{\partial z} \right)^2 \sigma_z(t)^2 (1 - \rho^2) dt. \end{aligned}$$

The first-order condition for the minimum is therefore

$$0 = -2\sigma_S(t) \left(-\pi(t)\sigma_S(t) + \frac{\partial V(t)}{\partial S}\sigma_S(t) + \frac{\partial V(t)}{\partial z}\sigma_z(t)\rho \right),$$

from which the lemma follows. \square

We notice that $w(t)$ in Lemma 8.9.4 can be written informally as

$$w(t) = \frac{\partial V(t)}{\partial z} \frac{\mathbb{E}_t(dz(t)|dS(t) = dS)}{dS}$$

which shows that the *minimum-variance* (MV) hedge ratio is obtained, in effect, by moving the z -process to its expected value, given an infinitesimal perturbation in the S -process. In other words, the hedge represents our best guess for a position in the underlying that will hedge moves in $V(t)$ caused by changes in *both* $S(t)$ and $z(t)$, as in Figure 8.9.

To further characterize the properties of the MV hedge weight, we insert the result of Lemma 8.9.4 into (8.51), which yields

$$d\Pi(t) = O(dt) + \frac{\partial V(t)}{\partial z} \sigma_z(t) \sqrt{1 - \rho^2} dB(t).$$

In other words, the MV hedge produces a portfolio that is not exposed to $W(t)$ but only to the orthogonal Brownian motion $B(t)$. If one thinks of $W(t)$ as “market” noise, we can say — in the language of the classical CAPM¹⁶ analysis — that the hedged portfolio has no *beta*. For this reason, the hedge construction in Lemma 8.9.4 is also sometimes known as a *zero-beta hedge*.

8.9.3 Minimum Variance Hedging: an Example

To better understand the practical ramifications of MV hedging, let us do a concrete example based on the SABR model from Section 8.6, which we here parameterize as

$$\begin{aligned} dS(t) &= \lambda \sqrt{z(t)} S(t)^c dW(t), \\ dz(t) &= \frac{1}{4} \eta^2 z(t) dt + \eta z(t) \left(\rho dW(t) + \sqrt{1 - \rho^2} dB(t) \right), \quad z(0) = 1. \end{aligned}$$

According to Lemma 8.9.4, the MV hedge ratio in SABR is

$$\pi_{\text{mv}}(t) = \frac{\partial V(t)}{\partial S} + \eta \sqrt{z(t)} \rho \frac{\partial V(t)}{\partial z} \frac{1}{\lambda S(t)^c}.$$

In a typical interest rate application $z(t) \approx 1$, $\lambda S(t)^c \approx 0.01$ and $\eta \approx 1$, such that, as a rule of thumb,

$$\pi_{\text{mv}}(t) \approx \frac{\partial V(t)}{\partial S} + 100 \times \rho \frac{\partial V(t)}{\partial z}.$$

For call and put options, the hedge adjustment to the “pure” delta $\partial V/\partial S$ is here typically negative, as we have $\partial V/\partial z > 0$ and, in normal market conditions, $\rho < 0$. This is consistent with Figure 8.8.

We now perform the following small experiment: we lock the correlation parameter at a pre-fixed value and then least-squares calibrate the SABR model to an actual market Black volatility smile. For a range of correlation parameters, we then compute “pure” deltas ($\partial V/\partial S$) and MV deltas (π_{mv}) for swaptions with different strikes. Using market data roughly consistent with the 5y×5y swaption volatility smile in the summer of 2005, the calibration results are in Table 8.1.

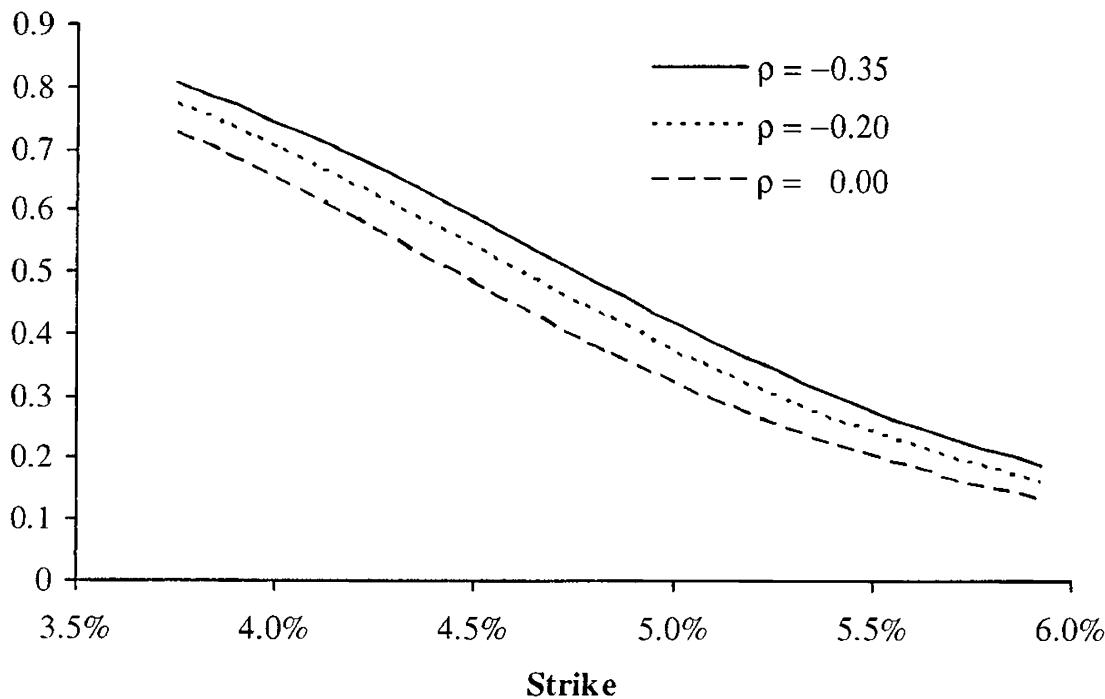
As one would expect, making correlation progressively more negative causes the skew power c to increase, from about 20% at $\rho = 0$ to nearly

¹⁶Capital Asset Pricing Model, see Sharpe [1964].

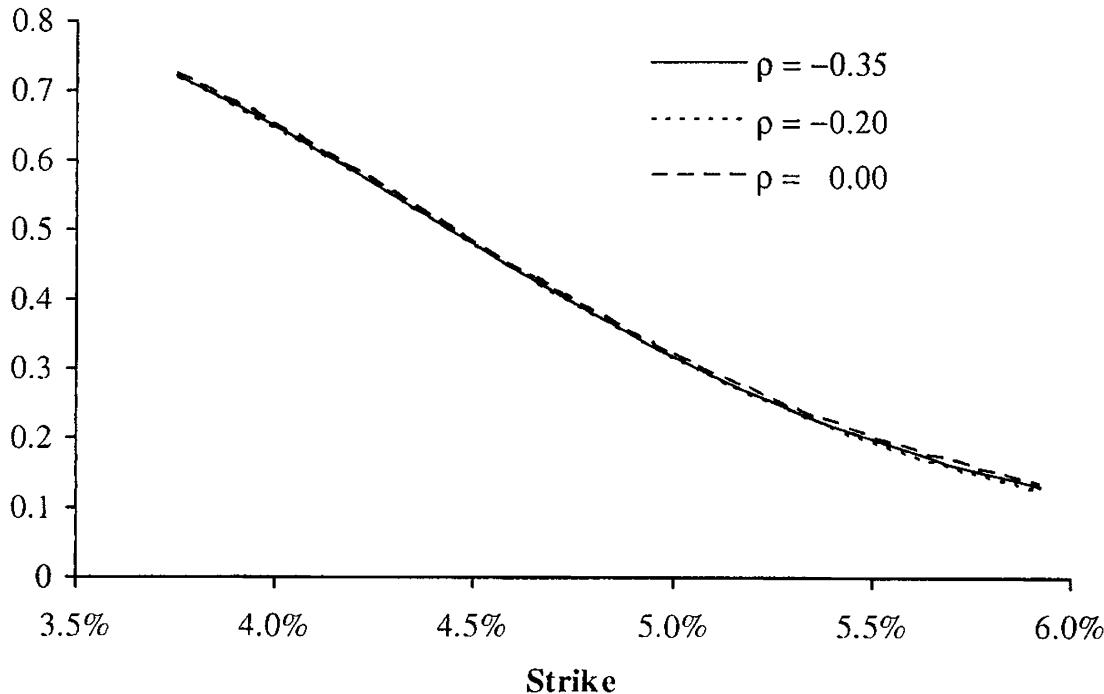
ρ	0	-0.1	-0.2	-0.3	-0.35
$\lambda S(0)^{1-c}$	0.135	0.136	0.137	0.139	0.140
c	0.223	0.432	0.648	0.877	0.999
η	0.684	0.686	0.696	0.712	0.726

Table 8.1. SABR Calibration Results

90% at $\rho = -0.3$, with other parameters being quite stable across different correlation choices. Figures 8.10 and 8.11 show the pure delta $\partial V / \partial S$ and the minimum variance delta π_{mv} for selected strikes and correlations. Clearly, the MV delta is here virtually independent of the choice of ρ , whereas the pure delta can increase quite substantially as correlation becomes more negative. It is clear from the figures that as long as hedge ratios are computed to be MV deltas, rather than pure deltas, the precise blend of local and stochastic volatility may not be critical, at least not for vanilla-like options in a skew-dominated market. This confirms a point we made earlier, in Section 8.1.

Fig. 8.10. Pure Delta

Notes: The figure shows the pure delta for the SABR models in Table 8.1.

Fig. 8.11. Minimum Variance Delta

Notes: The figure shows the minimum variance (MV) delta for the SABR models in Table 8.1.

8.A Appendix: Martingale Characterization, Moment Stability, and Other Fundamental Properties for General Variance Processes

As explained in Section 8.3, it is sometimes beneficial to consider a specification of the stochastic volatility model that is more general than (8.3)–(8.4). Let us consider a general power function for $\psi(z)$ in (8.2),

$$dS(t) = \lambda (bS(t) + (1 - b)L) \sqrt{z(t)} dW(t), \quad (8.52)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \eta z(t)^p dZ(t), \quad (8.53)$$

with $\langle dZ(t), dW(t) \rangle = \rho dt$. We assume $p > 0$. In this section we briefly outline important properties of such models. For more comprehensive treatment the reader is referred to Andersen and Piterbarg [2007]. Our first result spells out the boundary behavior of the stochastic variance process.

Proposition 8.A.1. *For the process (8.53), the following holds:*

1. 0 is always an attainable boundary for $0 < p < 1/2$.
2. 0 is an attainable boundary for $p = 1/2$, if $2z_0\theta < \eta^2$.
3. 0 is an unattainable boundary for $p > 1/2$.
4. ∞ is an unattainable boundary for all values of $p > 0$.

When $0 < p < 1/2$, the origin is always accessible and we need to impose a boundary condition at $z = 0$ to make the process unique. To ensure that the process for $z(\cdot)$ has a stationary distribution, we make the following natural choice:

Assumption 8.A.2. *For $0 < p < 1/2$, the process (8.53) for $z(\cdot)$ is reflected at the origin.*

The marginal one-dimensional distribution of $z(t)$ can in principle be computed numerically by various methods, such as PDE methods or by Fourier inversion of a characteristic function. It is often convenient, however, to have an easily-computable approximation. For that purpose, a stationary distribution, if one exists, can be useful. A stationary distribution for $z(\cdot)$ does indeed exist and can be easily computed.

Proposition 8.A.3. *Let $\pi(y)$ be the stationary distribution density for $z(\cdot)$ in (8.53). Under the assumptions listed above,*

$$\pi(y) = C(p)y^{-2p}e^{Q(y;p)}, \quad C(p)^{-1} \triangleq \int_0^\infty y^{-2p}e^{Q(y;p)}dy,$$

where the function $Q(y;p)$ is given by

1. $0 < p < 1/2$ or $1/2 < p < 1$ or $p > 1$:

$$Q(y;p) = \frac{2\theta}{\eta^2} \left(\frac{z_0 y^{1-2p}}{1-2p} - \frac{y^{2-2p}}{2-2p} \right).$$

2. $p = 1/2$:

$$Q(y;p) = \frac{2\theta}{\eta^2} (z_0 \ln y - y).$$

3. $p = 1$:

$$Q(y;p) = \frac{2\theta}{\eta^2} (-z_0/y - \ln y).$$

A-priori, $S(\cdot)$ defined by (8.52)–(8.53) is only a local martingale. In fact, under some circumstances, $S(\cdot)$ is a strict local martingale, , usually a significant technical complication. Specifically, we have the following result.

Proposition 8.A.4. *When $p \leq 1/2$ or $p > 3/2$, $S(\cdot)$ is a proper martingale. When $1/2 < p < 3/2$, $S(\cdot)$ is a martingale for $\rho \leq 0$ and a strict supermartingale for $\rho > 0$. For $p = 3/2$, $S(\cdot)$ is a martingale for $\rho \leq \frac{1}{2}\eta(\lambda b)^{-1}$ and a strict supermartingale for $\rho > \frac{1}{2}\eta(\lambda b)^{-1}$.*

What this proposition states is that the set of parameters $1/2 < p < 3/2$, $\rho > 0$, should be avoided in practical modeling. The SV model (8.3)–(8.4), as already noted, has no issues in this regard. If we use $\rho = 0$ — a typical choice in interest rate modeling as explained previously — all values of p

between 0 and $3/2$ are acceptable, at least as far as the martingale property is concerned.

In the model with $p = 1/2$, some moments of $S(\cdot)$ can become infinite, as stated in Proposition 8.3.10. With $p < 1/2$, this is no longer an issue:

Proposition 8.A.5. *In the model (8.52)–(8.53), if $p < 1/2$, moments $E(S(T)^u)$ of all orders $u \geq 1$ for all times T are finite.*

On the other hand, if $p > 1/2$ moments may be unstable. For instance:

Proposition 8.A.6. *In the model (8.52)–(8.53), if $p > 1/2$ and $\rho = 0$, all moments $E(S(T)^u)$ of all orders $u > 1$ for all times T are infinite.*

The case of non-zero correlation and $p > 1/2$ is more complicated; we refer the reader to Andersen and Piterbarg [2007].

Vanilla Models with Stochastic Volatility II

Having covered stochastic volatility models with time-homogeneous dynamics in Chapter 8, we are now ready to proceed with an analysis of the time-dependent case. As we shall see many examples of later in this book, stochastic volatility models with time-dependent parameters emerge naturally when vanilla models are used to approximate interest rate dynamics in a full term structure model.

In this chapter, we start out by modifying the Fourier analysis of Chapter 8 to cover time-dependent model parameters. We then proceed to introduce several approximation techniques that can speed up the calibration of model parameters to observable option prices. In particular, we continue our development of parameter averaging techniques, extending their scope to cover stochastic volatility and outlining in detail their usage in model calibration. Finally, the chapter gives detailed coverage of PDE and MC methods for general derivatives pricing; both of these numerical techniques are, as it turns out, rather tricky to apply to models with stochastic volatility, and an efficient implementation requires careful attention to detail.

9.1 Fourier Integration with Time-Dependent Parameters

As a start, let us consider extending the basic SV model (8.3)–(8.4) to allow for time-dependence of the volatility parameter¹ λ . That is, we now consider the P-measure dynamics

$$dS(t) = \lambda(t) (bS(t) + (1 - b)L) \sqrt{z(t)} dW(t), \quad (9.1)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \eta \sqrt{z(t)} dZ(t), \quad (9.2)$$

¹A further extension to time-dependence in η , ρ , and θ is trivial, and is covered in Remark 9.1.3.

where $\langle dZ(t), dW(t) \rangle = \rho dt$.

The model (9.1)–(9.2) still allows for call option pricing by the Fourier integration method of Section 8.4, provided that we can establish the moment-generating function (mgf) of $\ln X(t)$, with $X(t)$ being the linear function of $S(t)$ defined in Proposition 8.3.6. Let us retain the notation $\Psi_X(u; t)$ for

$$\Psi_X(u; t) = E \left(e^{u \ln X(t)} \right),$$

where the process for $X(t)$ now is modified from that of Proposition 8.3.6 to include time-dependence in λ :

$$dX(t)/X(t) = b\lambda(t)\sqrt{z(t)} dW(t), \quad X(0) = 1.$$

The following counterpart to Proposition 8.3.7 is easily proven.

Proposition 9.1.1. *In the model (9.1)–(9.2), for any $u \in \mathbb{C}$ for which the right-hand side exists, we have*

$$\Psi_X(u; t) = \Psi_{\overline{z\lambda^2}} \left(\frac{1}{2}b^2 u(u-1), u; t \right),$$

where we have defined

$$\Psi_{\overline{z\lambda^2}}(v, u; t) \triangleq E^{\widetilde{P}} \left(e^{v\overline{z\lambda^2}(t)} \right), \quad \overline{z\lambda^2}(t) \triangleq \int_0^t z(s)\lambda(s)^2 ds, \quad (9.3)$$

and under the new probability measure \widetilde{P} the process for $z(t)$ is

$$dz(t) = (\theta(z_0 - z(t)) + \rho\eta\lambda(t)buz(t)) dt + \eta\sqrt{z(t)} d\widetilde{Z}(t), \quad z(0) = z_0, \quad (9.4)$$

with $\widetilde{Z}(t)$ a \widetilde{P} -Brownian motion. If $\rho = 0$, $\widetilde{P} = P$ and $z(t)$ in (9.3) follows (9.2) rather than (9.4).

The following proposition demonstrates how to compute the moment-generating function of $\overline{z\lambda^2}(T)$.

Proposition 9.1.2. *The function $\Psi_{\overline{z\lambda^2}}(v, u; T)$ defined by (9.3) is given by*

$$\Psi_{\overline{z\lambda^2}}(v, u; T) = \exp(A(0, T) + z_0B(0, T)),$$

where $(A(t, T), B(t, T))$ solve the system of Riccati ODEs

$$\frac{d}{dt} A(t, T) + \theta z_0 B(t, T) = 0, \quad (9.5)$$

$$\frac{d}{dt} B(t, T) - (\theta - \rho\eta bu\lambda(t)) B(t, T) + \frac{\eta^2}{2} B(t, T)^2 + v\lambda(t)^2 = 0, \quad (9.6)$$

with the terminal conditions

$$B(T, T) = A(T, T) = 0.$$

Proof. Let us define

$$G(t, z) \triangleq \mathbb{E}^{\tilde{P}} \left(e^{v \int_t^T \lambda(s)^2 z(s) ds} \middle| z(t) = z \right).$$

Clearly,

$$\Psi_{z\lambda^2}(v, u; T) = G(0, z_0).$$

On the other hand, by the Feynman-Kac formula, $G(t, z)$ satisfies the following PDE,

$$\begin{aligned} \frac{\partial}{\partial t} G(t, z) + (\theta z_0 - (\theta - \rho \eta b u \lambda(t)) z) \frac{\partial}{\partial z} G(t, z) \\ + \frac{\eta^2}{2} z \frac{\partial^2}{\partial z^2} G(t, z) + v \lambda(t)^2 z G(t, z) = 0, \end{aligned} \quad (9.7)$$

with the terminal condition

$$G(T, z) = 1, \quad z \geq 0. \quad (9.8)$$

The PDE (9.7) is affine in z , i.e. all coefficients are linear functions of z . To solve it, we make the *ansatz* that the solution $G(t, z)$ is of the exponential form

$$G(t, z) = \exp(A(t, T) + zB(t, T)).$$

Substituting this conjectured solution into the PDE (9.7) and dividing by G , we get

$$\begin{aligned} \frac{d}{dt} A(t, T) + z \frac{d}{dt} B(t, T) + (\theta z_0 - (\theta - \rho \eta b u \lambda(t)) z) B(t, T) \\ + \frac{\eta^2}{2} z B(t, T)^2 + v \lambda(t)^2 z = 0. \end{aligned}$$

By collecting the coefficients on different powers of z , the two ODEs (9.5)–(9.6) emerge. Boundary conditions follow from (9.8). \square

The system of ODEs (9.5)–(9.6) can be solved numerically using the *Runge-Kutta method*, see e.g. Press et al. [1992]. In practice, it is common for the time-dependent volatility $\lambda(t)$ to be piecewise constant,

$$\lambda(t) = \lambda_i, \quad t \in (t_{i-1}, t_i],$$

for some $0 = t_0 < t_1 < \dots < t_I = T$. In this case, on each of the intervals $(t_{i-1}, t_i]$, the ODEs (9.5)–(9.6) can be solved in closed form, using the formulas from Proposition 8.3.8. By piecing these solutions together², we obtain the exact solution to the ODEs over the whole time interval $[0, T]$. However, for a given tolerance on accuracy, the Runge-Kutta method may still be faster than exact solution of the ODEs, as it avoids expensive evaluations of functions \exp , \ln , etc.

²The full procedure is described in Section 10.2.2.2.

Remark 9.1.3. So far, we assumed that η , ρ , and θ were constants. However, it follows easily from the proof of Proposition 9.1.2 that incorporation of time-dependence in η , ρ and θ is merely a matter of changing the ODEs (9.5)–(9.6) to

$$\frac{d}{dt} A(t, T) + \theta(t) z_0 B(t, T) = 0,$$

$$\frac{d}{dt} B(t, T) - (\theta(t) - \rho(t)\eta(t)b u \lambda(t)) B(t, T) + \frac{\eta(t)^2}{2} B(t, T)^2 + v \lambda(t)^2 = 0.$$

No matter which scheme is ultimately used to solve (9.5)–(9.6), combining the integration method of Theorem 8.4.4 with the integrand in Proposition 9.1.2 — possibly extended as in Remark 9.1.3 — allows for the pricing of call options by the Fourier methods in Section 8.4.

9.2 Asymptotic Expansion with Time-Dependent Volatility

As demonstrated in previous sections, the Fourier method constitutes a powerful tool for establishing a pricing algorithm for European options, provided that the underlying stochastic volatility process is of a sufficiently simple form. Should, say, the volatility function $\psi(z)$ for $z(t)$ be something other than \sqrt{z} , or should the skew function $\varphi(x)$ be more complicated than a linear form, analytic tractability (as in Proposition 9.1.2) is often lost and the Fourier method may not be feasible. However, asymptotic expansion methods can still be used in some situations and may, even for cases where Fourier methods do apply, offer a compelling (and very fast) approach to European option pricing.

To develop the asymptotic expansion approach, we return to the general skew functions $\varphi(x)$ and $\psi(z)$ in (8.1)–(8.2), under the simplifying (yet practically relevant) assumption that $\rho = 0$. As in the previous section, we will assume that the volatility $\lambda(t)$ is time-dependent. To summarize, the SDE system under consideration will be

$$dS(t) = \lambda(t)\varphi(S(t)) \sqrt{z(t)} dW(t), \quad (9.9)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \eta\psi(z(t)) dZ(t), \quad z(0) = z_0, \quad (9.10)$$

where $\langle dZ(t), dW(t) \rangle = 0$. The form of the time-dependence — as introduced here exclusively in $\lambda(t)$ — allows us to use time-change arguments similar to those in Section 7.6.1 to show that Lemma 8.5.4 as well as Proposition 8.5.5 still apply.

Lemma 9.2.1. *For the system (9.9)–(9.10) the results of Lemma 8.5.4 and Proposition 8.5.5 hold unchanged, provided we redefine (8.43) to*

$$(\mathcal{F}g)(\omega) = \int_{-\infty}^{\infty} e^{i\omega U} g(0, S(0); T^{-1}U) dU,$$

and make the substitutions

$$\lambda^2 \bar{z}(T) \rightarrow \overline{z\lambda^2}(T), \quad \lambda^2 U \rightarrow U, \quad \Psi_{\bar{z}} \rightarrow \Psi_{\overline{z\lambda^2}}.$$

For the special case $\psi(z) = \sqrt{z}$, Proposition 9.1.2 derives the expression for $\Psi_{\overline{z\lambda^2}}(u, 0; T)$. For more general choices of $\psi(z)$, we can rely on the PDE from Lemma 8.5.7, appropriately extended to time-dependent $\lambda(t)$. Specifically, $\Psi_{\overline{z\lambda^2}}(u, 0; T) = L(0, z_0; u)$, where $L(t, z; u)$ satisfies the PDE

$$\frac{\partial L}{\partial t} + \theta(z_0 - z) \frac{\partial L}{\partial z} + \frac{\eta^2}{2} \psi(z)^2 \frac{\partial^2 L}{\partial z^2} + u \lambda(t)^2 z L = 0, \quad (9.11)$$

subject to the boundary condition $L(T, z; u) = 1$. The equation can be solved numerically, or we can attempt to derive approximations. For the latter, we first introduce a centered transform

$$l(t, z; u) \triangleq L(t, z; u) e^{-u\mu_{\overline{z\lambda^2}}(t, z)}, \quad (9.12)$$

where, under mild regularity conditions on $\psi(z)$,

$$\begin{aligned} \mu_{\overline{z\lambda^2}}(t, z) &\triangleq \mathbb{E} \left(\int_t^T \lambda(s)^2 z(s) ds \middle| z(t) = z \right) \\ &= \int_t^T \lambda(s)^2 \mathbb{E}(z(s)|z(t) = z) ds \\ &= \int_t^T \lambda(s)^2 \left(z_0 + (z - z_0) e^{-\theta(s-t)} \right) ds. \end{aligned}$$

Introduction of $l(t, z; u)$ focuses attention on deviations of $\overline{z\lambda^2}(t)$ away from its mean, which can be expected to be small if η is small — a limit that we shall shortly examine. Insertion of (9.12) into (9.11) reveals that $l(t, z; u)$ satisfies

$$\frac{\partial l}{\partial t} + \theta(z_0 - z) \frac{\partial l}{\partial z} + \frac{\eta^2}{2} \psi(z)^2 \left\{ \frac{\partial^2 l}{\partial z^2} + l u^2 p(t)^2 + 2u p(t) \frac{\partial l}{\partial z} \right\} = 0, \quad (9.13)$$

where

$$p(t) = \int_t^T \lambda(s)^2 e^{-\theta(s-t)} ds \quad (9.14)$$

and $l(T, z; u) = 1$.

Lemma 9.2.2. *Let $p(t)$ be as in (9.14), and define $\tilde{\psi}(z) = \frac{1}{2}\psi(z)^2$ and $h(s, z) = z_0 + (z - z_0)e^{\theta(t-s)}$. An asymptotic expansion for the solution to (9.13) in terms of η^2 is given by*

$$l(t, z; u) = 1 + \eta^2 l_1(t, z; u) + \eta^4 l_2(t, z; u) + O(\eta^6),$$

where

$$\begin{aligned} l_1(t, z; u) &= u^2 l_{1,2}(t, z), \\ l_2(t, z; u) &= u^2 l_{2,2}(t, z) - u^3 l_{2,3}(t, z) + \frac{1}{2} u^4 (l_{1,2}(t, z))^2, \end{aligned}$$

and

$$\begin{aligned} l_{1,2}(t, z) &= \int_t^T p(s)^2 \tilde{\psi}(h(s, z)) \, ds, \\ l_{2,2}(t, z) &= \int_t^T e^{2\theta s} \tilde{\psi}(h(s, z)) \int_s^T e^{-2\theta v} p(v)^2 \tilde{\psi}''(h(v, z)) \, dv \, ds, \\ l_{2,3}(t, z) &= -2 \int_t^T e^{\theta s} p(s) \tilde{\psi}(h(s, z)) \int_s^T e^{-\theta v} p(v)^2 \tilde{\psi}'(h(v, z)) \, dv \, ds. \end{aligned}$$

Proof. Let

$$l(t, z; u) = 1 + \sum_{i \geq 1} l_i(t, z; u) \eta^{2i}.$$

Notice that odd powers of η are not used in the expansion, as only η^2 figures in the PDE (9.13). Inserting into (9.13) and collecting terms of order η^2 gives

$$\frac{\partial l_1}{\partial t} + \theta(z_0 - z) \frac{\partial l_1}{\partial z} + \frac{1}{2} u^2 p(t)^2 \psi(z)^2 = 0,$$

with terminal condition $l_1(T, z) = 0$. This simple PDE can be solved in closed form, yielding the solution listed in the lemma. The result for l_2 is established by collecting terms of order η^4 and proceeding as for l_1 . \square

While somewhat complicated in appearance, the expressions for the integrals $l_{1,2}$, $l_{2,2}$, and $l_{2,3}$ are trivial to implement on a computer. Indeed, due to the nested nature of the double integrals $l_{2,2}$ and $l_{2,3}$, all integrals can be computed in a single numerical integration loop, at negligible computational cost. In doing the integrals we start from the back, at time T , allowing us at each integration step to update the outer integral, as well as to resolve the inner integrals. In some cases of practical interest it is also possible to evaluate the integrals analytically.

Apart from potential direct application in the Fourier technique in Proposition 8.5.5, the result of Lemma 9.2.2 allows us to compute central moments as follows:

$$E \left(\left(\overline{z\lambda^2}(T) - \mu_{z\lambda^2}(0, z_0) \right)^n \right) = \left. \frac{\partial^n l(0, z_0; u)}{\partial u^n} \right|_{u=0}, \quad n = 1, 2, \dots \quad (9.15)$$

There are many ways to turn these moments into an option price expression. For instance, we could rely on a classical Gram-Charlier expansion (see

Ochi [1990]) or perhaps some parametric density family to express the full density of $z\lambda^2(T)$, to be used directly in (time-dependent generalizations of) equations (8.37) or (8.42). Alternatively, we can use Taylor expansions for a closed-form asymptotic result. Specifically, if the function g is defined as in Lemma 8.5.4, we can write

$$\begin{aligned} \mathbb{E}(f(S(T))) &= g(0, S(0); \bar{v}) \\ &\quad + \sum_{n=1}^{\infty} \frac{1}{n!T^n} \left. \frac{\partial^n g}{\partial v^n} \right|_{v=\bar{v}} \mathbb{E}\left(\left(z\lambda^2(T) - \mu_{z\lambda^2}(0, z_0)\right)^n\right), \end{aligned}$$

where the derivatives are to be evaluated at $\bar{v} \triangleq \mu_{z\lambda^2}(0, z_0)/T$.

From (9.15) and the expansion formula in Lemma 9.2.2, a few manipulations give the required result.

Lemma 9.2.3. *With $g(t, S; v)$ defined as in Lemma 8.5.4, we have to order $O(\eta^4)$*

$$\begin{aligned} \mathbb{E}(f(S(T))) &= g(0, S(0); \bar{v}) + T^{-2} (\eta^2 l_{1,2} + \eta^4 l_{2,2}) \frac{\partial^2 g}{\partial v^2} \\ &\quad - \eta^4 T^{-3} l_{2,3} \frac{\partial^3 g}{\partial v^3} + \frac{1}{2} \eta^4 T^{-4} l_{1,2}^2 \frac{\partial^4 g}{\partial v^4}, \end{aligned}$$

where all derivatives are evaluated at $\bar{v} = \mu_{z\lambda^2}(0, z_0)/T$.

To show an application of this lemma, consider the important special case of a call option $f(x) = (x - K)^+$.

Proposition 9.2.4. *Define the log-moneyness $k = \ln(K/S(0))$ and set $\tau = \int_0^T \lambda(s)^2 ds$. Also set*

$$q_1 = \mu_{z\lambda^2}(0, z_0)/T + \alpha_0 \eta^2 + \alpha_1 \eta^2 k^2 + O(\eta^4), \quad (9.16)$$

$$\begin{aligned} q_2 &= \mu_{z\lambda^2}(0, z_0)/T + (\alpha_0 \eta^2 + \beta_0 \eta^4) \\ &\quad + (\alpha_1 \eta^2 + \beta_1 \eta^4) k^2 + \beta_2 \eta^4 k^4 e^{-\Lambda \eta^2 k^2} + O(\eta^6), \end{aligned} \quad (9.17)$$

where Λ is an arbitrary positive number and the coefficients $\alpha_0, \alpha_1, \beta_0, \beta_1, \beta_2$ are given in Appendix 9.B. Then the value of a European call option in the model (9.9)–(9.10) is given by

$$c(0, S; T, K) \approx S(0)\Phi(d_+) - K\Phi(d_-), \quad (9.18)$$

$$d_{\pm} = \frac{-k \pm \sigma_{\text{imp}}^2 T/2}{\sigma_{\text{imp}} \sqrt{T}},$$

where, to order η^2 ,

$$\sigma_{\text{imp}} = \Omega_0 \sqrt{q_1} + \Omega_1 q_1^{3/2} T + O(T^2),$$

or, to order η^4 ,

$$\sigma_{\text{imp}} = \Omega_0 \sqrt{q_2} + \Omega_1 q_2^{3/2} T + O(T^2).$$

Also, we have

$$\begin{aligned}\Omega_0 &= \frac{-k}{\int_K^{S(0)} \varphi(u)^{-1} du}, \\ \Omega_1 &= -\frac{\Omega_0}{\left(\int_K^{S(0)} \varphi(u)^{-1} du\right)^2} \ln \left(\Omega_0 \left(\frac{KS(0)}{\varphi(K)\varphi(S(0))} \right)^{1/2} \right).\end{aligned}$$

Proof. (Sketch). For the case of a call option, the function g can be approximated using the small-time expansion result in Proposition 7.5.1; we here choose to expand around a log-normal model, so $\beta = 0$ in the proposition. Using the resulting expression to evaluate the terms in Lemma 9.2.3 yields, after some work, a direct expansion for the call option price. It is often more accurate to convert the price expansion into an expansion in implied “skew variance” v^* , where v^* satisfies

$$E((S(T) - K)^+) = g(0, S(0); v^*). \quad (9.19)$$

We write

$$v^* = \bar{v} + \eta^2 v_1^* + \eta^4 v_2^* + \dots, \quad (9.20)$$

insert this expression into (9.19) and Taylor-expand around \bar{v} . Matching the resulting expression against the direct expansion for the call option price yields closed-form expressions for v_1^* and v_2^* . These results are such that

$$\bar{v} + \eta^2 v_1^* = q_1, \quad \bar{v} + \eta^2 v_1^* + \eta^4 v_2^* = q_2,$$

where q_1 and q_2 are defined in (9.16) and (9.17), respectively. Another application of Proposition 7.5.1 turns the skew variance into an implied Black volatility,

$$\sigma_{\text{imp}} \sqrt{T} = \Omega_0 \sqrt{v^* T} + \Omega_1 (v^* T)^{3/2} + \dots.$$

The proposition follows. \square

Remark 9.2.5. Full details for the proof of Proposition 9.2.4 and tests of the precision of the expansion can be found in Andersen and Brotherton-Ratcliffe [2005].

9.3 Averaging Methods

The Fourier integration method from Section 9.1 involves numerical integration of a function that itself is calculated numerically by solving a coupled system of ODEs. If both the integral and the ODEs are discretized with N steps, the complexity of the scheme $O(N^2)$, which could be costly. On the other hand, the asymptotic expansion method from Section 9.2 is fast but may not be accurate enough for certain values of model parameters, especially high η . In this section we develop the parameter averaging approach to time-dependent model parameters that is both fast and accurate. We have seen applications of the method to local volatility models already, in Section 7.6.2.

9.3.1 Volatility Averaging

We initially work with the model (9.1)–(9.2) with zero correlation, $\rho = 0$. Our goal is to replace the time-dependent $\lambda(t)$ with a constant $\bar{\lambda}$ in such a way that pricing of vanilla options at a given maturity T is preserved to good approximation. For this, we first notice that a European option price can be represented as an integral of a known function against the distribution of the term stochastic variance, a representation we have already fruitfully used in Sections 8.5 and 9.2. In particular, for an at-the-money option, where $K = S(0)$,

$$\mathbb{E} \left((S(T) - S(0))^+ \right) = \mathbb{E} \left(\mathbb{E} \left((S(T) - S(0))^+ \mid \{z(t), t \in [0, T]\} \right) \right). \quad (9.21)$$

Because the Brownian motion that drives $z(t)$ is independent of the Brownian motion that drives $S(t)$, the distribution of $S(T)$ in the model (9.2) is displaced log-normal when conditioned on a particular path of $z(t)$. Hence, the inner conditional expectation in (9.21) can be evaluated easily to yield

$$\mathbb{E} \left((S(T) - S(0))^+ \right) = \mathbb{E} \left(h \left(\overline{z\lambda^2}(T) \right) \right), \quad (9.22)$$

where $\overline{z\lambda^2}(T)$ is defined by (9.3) and the function $h(x)$ is the displaced log-normal at-the-money option value as function of variance:

$$h(x) = \frac{bS(0) + (1-b)L}{b} (2\Phi(b\sqrt{x}/2) - 1). \quad (9.23)$$

Given the practical importance of correctly pricing at-the-money options, the problem of finding the effective, time-independent model volatility can be cast into the problem of finding such $\bar{\lambda}$ that

$$\mathbb{E} \left(h \left(\int_0^T \lambda(t)^2 z(t) dt \right) \right) = \mathbb{E} \left(h \left(\bar{\lambda}^2 \int_0^T z(t) dt \right) \right) \quad (9.24)$$

or, in our notations,

$$\mathbb{E} \left(h \left(\overline{z\lambda^2}(T) \right) \right) = \mathbb{E} \left(h \left(\bar{\lambda}^2 \bar{z}(T) \right) \right).$$

Neither of the expected values in (9.24) is available in closed form. However, the moment-generating functions of both $\overline{z\lambda^2}(T)$ and $\bar{z}(T)$ are available in closed form and as a solution to a system of ODEs, respectively (see Propositions 8.3.8 and 9.1.2). This observation suggests approximating $h(x)$ with a function of exponential form

$$h(x) \approx a + b e^{cx}. \quad (9.25)$$

We choose the coefficients a, b, c to get the best local second-order fit at the mean of $\overline{z\lambda^2}(T)$,

$$h(\zeta_T) = a + b e^{c\zeta_T}, \quad h'(\zeta_T) = b c e^{c\zeta_T}, \quad h''(\zeta_T) = b c^2 e^{c\zeta_T}, \quad (9.26)$$

where

$$\zeta_T = \mathbb{E} \left(\overline{z\lambda^2}(T) \right) = \mu_{\overline{z\lambda^2}}(0, z_0) = z_0 \int_0^T \lambda(t)^2 dt.$$

Clearly

$$c = \frac{h''(\zeta_T)}{h'(\zeta_T)}, \quad (9.27)$$

and the problem (9.24) can be approximated with

$$a + b \mathbb{E} \left(e^{cz\overline{\lambda^2}(T)} \right) = a + b \mathbb{E} \left(e^{c\bar{\lambda}^2 \bar{z}(T)} \right) \Rightarrow \mathbb{E} \left(e^{cz\overline{\lambda^2}(T)} \right) = \mathbb{E} \left(e^{c\bar{\lambda}^2 \bar{z}(T)} \right), \quad (9.28)$$

which gives us an *effective volatility* approximation result that we formulate as a theorem.

Theorem 9.3.1. *Values of European options with expiry T in the model (9.1)–(9.2) are well approximated by their values in the model (8.3)–(8.4) with λ set to the effective SV volatility $\bar{\lambda}$, which solves the equation*

$$\Psi_{\bar{z}} \left(\frac{h''(\zeta_T)}{h'(\zeta_T)} \bar{\lambda}^2, 0; T \right) = \Psi_{\overline{z\lambda^2}} \left(\frac{h''(\zeta_T)}{h'(\zeta_T)}, 0; T \right), \quad (9.29)$$

where

$$\zeta_T = z_0 \int_0^T \lambda(t)^2 dt,$$

the function $h(x)$ is given by (9.23), and the moment-generating functions $\Psi_{\bar{z}}$ and $\Psi_{\overline{z\lambda^2}}$ are given by Propositions 8.3.8 and 9.1.2, respectively.

Proof. Follows after replacing the problem (9.24) with (9.28), using the expression (9.27) for c . \square

Remark 9.3.2. The expression on the left-hand side of (9.29) can be computed in closed form; the right-hand side is straightforward to calculate from Proposition 9.1.2 and the accompanying remarks. Equation (9.29) can be solved for $\bar{\lambda}^2$ in just a couple of Newton-Raphson iterations, starting from an initial guess of $T^{-1} \int_0^T \lambda(t)^2 dt$.

Remark 9.3.3. The effective volatility $\bar{\lambda}$ as given by Theorem 9.3.1 is second-order accurate in the sense that the approximation (9.25) is second-order accurate with the choice of parameters in (9.26). We note that the method does not readily lend itself to higher-order approximations but this is of little relevance as the quality of the approximation is excellent as is.

9.3.2 Skew Averaging

The slope of the volatility smile in the SV model (8.3)–(8.4) is controlled by the skew parameter b . In this section we make the skew parameter time-dependent, and consider a model driven by the SDEs

$$dS(t) = \lambda(t) (b(t)S(t) + (1 - b(t)) L) \sqrt{z(t)} dW(t), \quad (9.30)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \eta \sqrt{z(t)} dZ(t), \quad (9.31)$$

with $\langle dZ(t), dW(t) \rangle = 0$. In Section 7.6.2 we derived the formula for the effective, or average, skew for local volatility models, see Proposition 7.6.2 and Corollary 7.6.3. The extension of these results to stochastic volatility models is straightforward, leading to a similar expression with somewhat more complicated averaging weights, as the following proposition demonstrates.

Proposition 9.3.4. *The effective skew \bar{b} for the equation*

$$dS(t) = \lambda(t) (b(t)S(t) + (1 - b(t)) S(0)) \sqrt{z(t)} dW(t)$$

over a time horizon $[0, T]$ is given by

$$\bar{b} = \int_0^T b(t) w_T(t) dt, \quad (9.32)$$

where the weights $w_T(t)$ are given by

$$w_T(t) = \frac{v(t)^2 \lambda(t)^2}{\int_0^T v(s)^2 \lambda(s)^2 ds}, \quad (9.33)$$

$$v(t)^2 = z_0^2 \int_0^t \lambda(s)^2 ds + z_0 \eta^2 e^{-\theta t} \int_0^t \lambda(s)^2 \frac{e^{\theta s} - e^{-\theta s}}{2\theta} ds.$$

The result in Proposition 9.3.4 can be derived by the same technique that lead to Proposition 7.6.2 and Corollary 7.6.3. Alternatively, it can be found by the small-noise expansion method in Section 7.6.3. We leave the details of these derivations to the reader and, for instructional value, instead list a third proof based on Markovian semi-groups in Appendix 9.A, see also Piterbarg [2005b]. The fact that the same result is obtained as a solution to a number of differently posed problems of skew averaging suggests robustness and general applicability.

It will be useful for the next section to derive an extension of Proposition 9.3.4 to cover the process $z(t)$ with time-dependent volatility of variance. Specifically, let us use the following dynamics for the stochastic variance process

$$dz(t) = \theta(z_0 - z(t)) dt + \eta(t)\sqrt{z(t)} dZ(t). \quad (9.34)$$

Corollary 9.3.5. *The effective skew \bar{b} for the equation*

$$dS(t) = \lambda(t)(b(t)S(t) + (1 - b(t))S(0))\sqrt{z(t)} dW(t)$$

with $z(t)$ following (9.34) over a time horizon $[0, T]$ is given by

$$\bar{b} = \int_0^T b(t)w_T(t) dt, \quad (9.35)$$

where the weights $w_T(t)$ are given by

$$w_T(t) = \frac{\widehat{v}(t)^2 \lambda(t)^2}{\int_0^T \widehat{v}(t)^2 \lambda(t)^2 dt}, \quad (9.36)$$

$$\widehat{v}(t)^2 = z_0^2 \int_0^t \lambda(s)^2 ds + z_0 e^{-\theta t} \int_0^t \lambda(s)^2 e^{-\theta s} \int_0^s \eta(u)^2 e^{2\theta u} du ds.$$

Proof. The proof or the corollary proceeds as the proof (in Appendix 9.A) of Proposition 9.3.4, but using

$$E(z(t)^2) = z_0^2 + z_0 \int_0^t \eta(u)^2 e^{-2\theta(t-u)} du \quad (9.37)$$

instead of (9.100) in (9.101) for $z(t)$ given by (9.34). \square

9.3.3 Volatility of Variance Averaging

Finally, we turn our attention to the problem of averaging the volatility of variance η in (9.1). More precisely, suppose we have a stochastic variance process with time-dependent volatility of variance (9.34). We would like to find a constant parameter $\bar{\eta}$ such that the model (9.30), (9.34) is approximated by the model (9.30), (9.31) with $\eta = \bar{\eta}$.

Before discussing our proposed solution method, we note that usage of time-dependent volatility of variance $\eta(t)$ for model calibration purposes may not be quite as necessary as for other parameters. Fundamentally, a time-dependent η will allow us to control the term structure of volatility smile convexity in the maturity direction. On the other hand, we already have control over the curvatures of volatility smiles at different times T via θ , the mean reversion of variance parameter: higher values of θ make implied volatility smiles flatten faster as option expiries increase, while lower values make them flatten slower, see Sections 8.2 and 8.7. Even though the level of control granted through θ is rather crude, it is often sufficient in practice, all the more so since the volatility smile curvatures are typically not observable to a high degree of precision.

The curvature of the volatility smile is related to the kurtosis of the distribution of $S(T)$ which, in stochastic volatility models, is controlled by the variance of the quantity

$$\overline{z\lambda^2}(T) = \int_0^T \lambda(t)^2 z(t) dt,$$

i.e. the integrated stochastic variance to expiry time T . Since the curvature of the smile is the main effect of the volatility of variance parameter η , a representative constant volatility of variance $\bar{\eta}$ should intuitively be chosen as the solution to

$$E \left(\left(\int_0^T \lambda(t)^2 \widehat{z}(t) dt \right)^2 \right) = E \left(\left(\int_0^T \lambda(t)^2 z(t) dt \right)^2 \right), \quad (9.38)$$

where $z(t)$ follows (9.34) and $\widehat{z}(t)$ follows (9.31).

Theorem 9.3.6. *For (9.34), the effective volatility of variance to maturity T , derived from the condition (9.38), is given by*

$$\bar{\eta}^2 = \frac{\int_0^T \eta(t)^2 \rho_T(t) dt}{\int_0^T \rho_T(t) dt},$$

where the weight function $\rho_T(t)$ is given by

$$\rho_T(r) = \int_r^T ds \int_s^T dt \lambda(t)^2 \lambda(s)^2 e^{-\theta(t-s)} e^{-2\theta(s-r)}.$$

Proof. While the proof is straightforward, we here provide full details in order to demonstrate some generally useful manipulations for the computations of moments in stochastic volatility models. First, we have

$$\begin{aligned}
& \mathbb{E} \left(\left(\int_0^T \lambda(t)^2 z(t) dt \right)^2 \right) \\
&= 2 \int_0^T dt \int_0^t ds \lambda(t)^2 \lambda(s)^2 \mathbb{E}(z(t)z(s)) \\
&= 2 \int_0^T dt \int_0^t ds \lambda(t)^2 \lambda(s)^2 e^{-\theta(t-s)} \mathbb{E}(z(s)^2) \\
&\quad + 2 \int_0^T dt \int_0^t ds \lambda(t)^2 \lambda(s)^2 (1 - e^{-\theta(t-s)}) z_0 \mathbb{E}(z(s)).
\end{aligned}$$

Using (9.37) for $\mathbb{E}(z(s)^2)$ we get

$$\begin{aligned}
& \mathbb{E} \left(\left(\int_0^T \lambda(t)^2 z(t) dt \right)^2 \right) = 2z_0^2 \int_0^T dt \int_0^t ds \lambda(t)^2 \lambda(s)^2 e^{-\theta(t-s)} \\
&\quad + 2z_0 \int_0^T dt \int_0^t ds \lambda(t)^2 \lambda(s)^2 e^{-\theta(t-s)} \int_0^s \eta(r)^2 e^{-2\theta(s-r)} dr \\
&\quad + 2 \int_0^T dt \int_0^t ds \lambda(t)^2 \lambda(s)^2 (1 - e^{-\theta(t-s)}) z_0 \mathbb{E}(z(s)).
\end{aligned}$$

Changing the order of integration for the second term, we obtain

$$\begin{aligned}
& \mathbb{E} \left(\left(\int_0^T \lambda(t)^2 z(t) dt \right)^2 \right) = 2z_0^2 \int_0^T dt \int_0^t ds \lambda(t)^2 \lambda(s)^2 e^{-\theta(t-s)} \\
&\quad + 2z_0 \int_0^T dr \eta(r)^2 \int_r^T ds \int_s^T dt \lambda(t)^2 \lambda(s)^2 e^{-\theta(t-s)} e^{-2\theta(s-r)} \\
&\quad + 2 \int_0^T dt \int_0^t ds \lambda(t)^2 \lambda(s)^2 (1 - e^{-\theta(t-s)}) z_0 \mathbb{E}(z(s)).
\end{aligned}$$

If we define

$$\rho_T(r) = \int_r^T ds \int_s^T dt \lambda(t)^2 \lambda(s)^2 e^{-\theta(t-s)} e^{-2\theta(s-r)},$$

the equation (9.38) can be rewritten in the form

$$\int_0^T \bar{\eta}^2 \rho_T(t) dt = \int_0^T \eta(t)^2 \rho_T(t) dt.$$

The theorem is proved. \square

Remark 9.3.7. While we used zero correlation between the underlying and its stochastic variance both in motivating our results and in deriving them,

the same approach can be applied in the non-zero correlation case. Some results, in particular Proposition 9.3.4 and Theorem 9.3.6, remain unchanged. On the other hand, the effective volatility formula in Theorem 9.3.1 is based on the representation (9.22) which, clearly, does not hold with non-zero correlation; despite that, the formula can still be used with good accuracy.

9.3.4 Calibration by Parameter Averaging

The main application of the averaging formulas developed above is in creating efficient model calibration algorithms. In this section, we discuss in some detail how such an algorithm could proceed; the principles that we outline here shall be used repeatedly later in this book. Now, suppose a collection of expiries

$$0 = T_0 < T_1 < T_2 < \dots < T_N$$

is given, as well as a collection of strikes K_1, \dots, K_M . Let the market values of European call options with expiries T_n and strikes K_m be denoted by

$$\{\hat{c}_{n,m}, \quad n = 1, \dots, N, \quad m = 1, \dots, M\}.$$

Our objective is to find time-dependent model parameters $\lambda(t)$, $b(t)$, and $\eta(t)$ such that the model

$$dS(t) = \lambda(t) (b(t)S(t) + (1 - b(t)) L) \sqrt{z(t)} dW(t), \quad (9.39)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \eta(t) \sqrt{z(t)} dZ(t), \quad (9.40)$$

values European options with expiries T_n , $n = 1, \dots, N$, and strikes K_m , $m = 1, \dots, M$, as closely as possible to their market values³ $\{\hat{c}_{n,m}\}$.

Let us denote the prices of options in the model (9.39)–(9.40) by

$$c_{n,m} = c_{n,m}(\mathcal{X}),$$

where by \mathcal{X} we denote the state of the model,

$$\mathcal{X} = \{\lambda(\cdot), b(\cdot), \eta(\cdot)\}.$$

Typically, calibration would be performed by solving the following non-linear optimization problem

$$\{\lambda(\cdot), b(\cdot), \eta(\cdot)\} = \operatorname{argmin}_{n,m} \sum (c_{n,m}(\mathcal{X}) - \hat{c}_{n,m})^2, \quad (9.41)$$

³In interest rate markets, the underlyings for options of different expiries are often different, in the sense that they represent swap rates of different tenors and fixing dates. We will deal with such complications in due time.

where⁴ $c_{n,m}(\mathcal{X})$'s are obtained in some sort of numerical procedure. With the averaging formulas, an appealing alternative is available. To describe it, let us denote triples of SV “market” parameter values by $\{\widehat{\lambda}_n, \widehat{b}_n, \widehat{\eta}_n\}$, $n = 1, \dots, N$, determined such that the market prices of European options expiring at time T_n , i.e. $\{\widehat{c}_{n,m}, m = 1, \dots, M\}$, match prices obtained in the model

$$dS(t) = \widehat{\lambda}_n \left(\widehat{b}_n S(t) + (1 - \widehat{b}_n) L \right) \sqrt{z(t)} dW(t), \quad (9.42)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \widehat{\eta}_n \sqrt{z(t)} dZ(t). \quad (9.43)$$

Sets of market parameters are routinely maintained and updated by trading desks, and instead of considering $\{\widehat{c}_{n,m}\}$ to be fundamental market inputs, we can think of $\{\widehat{\lambda}_n, \widehat{b}_n, \widehat{\eta}_n\}$, $n = 1, \dots, N$, as such. We often refer to them as “term” parameters to highlight the fact that they are constant for the whole “term”, or life, of the relevant options.

Critically, the averaging formulas link time-dependent parameters $\{\lambda(t), b(t), \eta(t)\}$ to constant parameters $\{\widehat{\lambda}_n, \widehat{b}_n, \widehat{\eta}_n\}$, $n = 1, \dots, N$, directly without referencing option values. To take advantage of this, let us denote by

$$\{\bar{\lambda}_n(\mathcal{X}), \bar{b}_n(\mathcal{X}), \bar{\eta}_n(\mathcal{X})\}$$

the averaged parameters (to time T_n) for the model (9.39)–(9.40). Then the optimization problem (9.41) can be replaced by a more convenient one,

$$\begin{aligned} \{\lambda(\cdot), b(\cdot), \eta(\cdot)\} = \operatorname{argmin} & \left(W_\lambda \sum_n (\bar{\lambda}_n(\mathcal{X}) - \widehat{\lambda}_n)^2 \right. \\ & \left. + W_b \sum_n (\bar{b}_n(\mathcal{X}) - \widehat{b}_n)^2 + W_\eta \sum_n (\bar{\eta}_n(\mathcal{X}) - \widehat{\eta}_n)^2 \right), \end{aligned} \quad (9.44)$$

where W_λ , W_b , and W_η are weights linked to relative importance of matching particular parameters. Compared to (9.41), this norm formulation is both more intuitive to traders — who often tend to think about the state of the market in terms of model parameters, rather than in terms of absolute option prices — and computationally advantageous, insofar as the norm requires no outright computation of option values.

In practice, the calibration (9.44) needs not be performed by brute-force optimization. By carefully choosing the order of calculations, calibration can be split into independent sub-calibrations: one for volatility of variance (η); one for skewness (b); and one for volatility (λ). Skew and volatility of variance calibrations can be performed by matrix manipulations, and the volatility calibration can be split into a sequence of numerically solved one-dimensional equations. To describe this calibration idea in more detail,

⁴Often different terms are weighted differently.

let us first collect all relevant averaging results for easy reference. For the volatility of variance, we have from Theorem 9.3.6,

$$\bar{\eta}_n(\mathcal{X})^2 = \frac{\int_0^{T_n} \eta(t)^2 \rho_{T_n}(t; \lambda(\cdot)) dt}{\int_0^{T_n} \rho_{T_n}(t; \lambda(\cdot)) dt}, \quad n = 1, \dots, N, \quad (9.45)$$

where we have now explicitly indicated the dependence of weights $\rho_T(t; \lambda(\cdot))$ on the volatility function $\lambda(t)$. For the skews, we have from Corollary 9.3.5,

$$\bar{b}_n(\mathcal{X}) = \int_0^{T_n} b(t) w_{T_n}(t; \lambda(\cdot), \eta(\cdot)) dt, \quad n = 1, \dots, N, \quad (9.46)$$

where again the dependence of weights $w_T(t; \lambda(\cdot), \eta(\cdot))$ on model parameters is highlighted. Finally, the equations for volatilities from Theorem 9.3.1 are

$$\bar{\lambda}_n(\mathcal{X}) = F(\lambda(\cdot); \bar{b}_n(\mathcal{X}), \bar{\eta}_n(\mathcal{X})), \quad n = 1, \dots, N, \quad (9.47)$$

where, in the notation of Theorem 9.3.1,

$$F(\lambda(\cdot); \bar{b}_n(\mathcal{X}), \bar{\eta}_n(\mathcal{X})) = \sqrt{\frac{h'(\zeta_{T_n})}{h''(\zeta_{T_n})} \times \Psi_{\bar{z}}^{-1} \left(\Psi_{\bar{z}\lambda^2} \left(\frac{h''(\zeta_{T_n})}{h'(\zeta_{T_n})}, 0; T \right), 0; T \right)},$$

$$\zeta_{T_n} = z_0 \int_0^{T_n} \lambda(t)^2 dt.$$

Note that the function F depends on \bar{b}_n through h , and on $\bar{\eta}_n$ through $\Psi_{\bar{z}\lambda^2}$ and $\Psi_{\bar{z}}$.

Equations (9.45)–(9.47) can be discretized if the model parameters are constant between option expiry dates $\{T_n\}_{n=1}^N$, a common assumption in practice. In this case, we can define λ_i , b_i and η_i by

$$\lambda(t) = \sum_{i=1}^N \lambda_i 1_{\{t \in (T_{i-1}, T_i]\}},$$

$$b(t) = \sum_{i=1}^N b_i 1_{\{t \in (T_{i-1}, T_i]\}},$$

$$\eta(t) = \sum_{i=1}^N \eta_i 1_{\{t \in (T_{i-1}, T_i]\}}.$$

In addition, we discretize the weights and define $\rho_{n,i}(\lambda(\cdot))$ and $w_{n,i}(\lambda(\cdot), \eta(\cdot))$ by

$$\rho_{T_n}(t; \lambda(\cdot)) = \sum_{i=1}^n \rho_{n,i}(\lambda(\cdot)) 1_{\{t \in (T_{i-1}, T_i]\}},$$

$$w_{T_n}(t; \lambda(\cdot), \eta(\cdot)) = \sum_{i=1}^n w_{n,i}(\lambda(\cdot), \eta(\cdot)) 1_{\{t \in (T_{i-1}, T_i]\}}.$$

Denote

$$\bar{\rho}_{n,i}(\lambda(\cdot)) = \frac{\rho_{n,i}(\lambda(\cdot))}{\int_0^{T_n} \rho_{T_n}(t; \lambda(\cdot)) dt}.$$

Our goal is to solve three systems of equations:

$$\sum_{i=1}^n \bar{\rho}_{n,i}(\lambda(\cdot)) (T_i - T_{i-1}) \eta_i^2 = (\hat{\eta}_n)^2, \quad (9.48)$$

$$\sum_{i=1}^n w_{n,i}(\lambda(\cdot), \eta(\cdot)) (T_i - T_{i-1}) b_i = \hat{b}_n, \quad (9.49)$$

$$F(\lambda(\cdot); \bar{b}_n(\mathcal{X}), \bar{\eta}_n(\mathcal{X})) = \hat{\lambda}_n, \quad (9.50)$$

for $n = 1, \dots, N$. At first glance this does not seem entirely straightforward. For example, the system (9.48) appears to be a linear system of equations in $\eta_1^2, \dots, \eta_N^2$, but the coefficients $\bar{\rho}_{n,i}(\lambda(\cdot))$ depend on $\lambda(t)$, another unknown model parameter. However, by iteratively solving these equations in the right order, we can design a very efficient algorithm, which we now proceed to describe in detail.

First, we note that the equations on volatilities (9.50) do not depend on any other model parameters. They do depend on term parameters $\bar{b}_n(\mathcal{X})$, $\bar{\eta}_n(\mathcal{X})$, which we just replace with their market values, thus solving

$$F(\lambda(\cdot); \hat{b}_n, \hat{\eta}_n) = \hat{\lambda}_n, \quad n = 1, \dots, N.$$

The n -th equation in this series only involves λ_i 's for $i = 1, \dots, n$, so the n -th equation can be rewritten as

$$F(\lambda_1, \dots, \lambda_n; \hat{b}_n, \hat{\eta}_n) = \hat{\lambda}_n.$$

The case $n = 1$ has the trivial solution

$$\lambda_1^* = \hat{\lambda}_1.$$

Proceeding iteratively in n , the n -th equation is reduced to

$$F(\lambda_1^*, \dots, \lambda_{n-1}^*, \lambda_n; \hat{b}_n, \hat{\eta}_n) = \hat{\lambda}_n, \quad (9.51)$$

where the λ_i^* , $i = 1, \dots, n-1$, are the model parameters already solved for. Thus, the first step of calibration consists of solving the system of equations (9.50) as N decoupled one-dimensional equations (9.51).

On the second step, we solve the linear system (9.48) for η_i^2 , $i = 1, \dots, N$. The coefficients of the system depend on λ_i 's which have already been computed, and we solve

$$\sum_{i=1}^n \bar{\rho}_{n,i}(\lambda^*(\cdot)) (T_i - T_{i-1}) \eta_i^2 = (\hat{\eta}_n)^2, \quad n = 1, \dots, N.$$

The solution η_i^* , $i = 1, \dots, N$, to this system can either be found by matrix methods, or by simple sequential substitution since the n -th equation involves η_i^2 for $i = 1, \dots, n$ only.

Finally, on the third step, we solve the linear system

$$\sum_{i=1}^n w_{n,i}(\lambda^*(\cdot), \eta^*(\cdot))(T_i - T_{i-1})b_i = \hat{b}_n, \quad n = 1, \dots, N, \quad (9.52)$$

for b_i , $i = 1, \dots, N$. This system is obtained from (9.49) by substituting $\lambda(\cdot)$, $\eta(\cdot)$ with their solved-for values $\lambda^*(\cdot)$, $\eta^*(\cdot)$. Again, the system can be solved sequentially.

To prevent overfitting, it is often useful to regularize the optimization problem through introduction of smoothing terms in the objective function. This can help to, for example, dampen the noise that could be present in market-observed parameters. Taking (9.52) as an example and fixing a smoothing weight $W > 0$, we can replace (9.52) with the minimization problem

$$\begin{aligned} \sum_{n=1}^N \left(\sum_{i=1}^n w_{n,i}(\lambda^*(\cdot), \eta^*(\cdot))(T_i - T_{i-1})b_i - \hat{b}_n \right)^2 \\ + W \sum_{i=2}^N (b_i - b_{i-1})^2 \rightarrow \min. \end{aligned} \quad (9.53)$$

This is a simple quadratic minimization problem with no constraints and is easily solved by linear algebra methods, see Golub and van Loan [1989]. The same regularization idea could be applied to the problem of finding $\lambda(t)$ and $b(t)$.

If the regularization weight W in (9.53) is too high then the averaged skew calculated by the model can be significantly different from the market skew, $\bar{b}_n(\mathcal{X}^*) \neq b_n$, $n = 1, \dots, N$. By itself this may not be such a bad thing as one may prefer a smoother model skew over the exact fit to market skews. However, this poses problems to the *volatility* calibration, as the equation for model volatility (9.51) used the “wrong” skew (and volatility of variance as well, were we to apply regularization to that). The exact fit to market volatilities is often much more important than the exact fit to skews or volatilities of variance. Fortunately, this problem is easy to rectify by solving the system (9.51) again, this time using the true model averaged skews $\bar{b}_n(\mathcal{X}^*)$ (and volatilities of variance) on the left-hand side of (9.51) which are available at this stage of the algorithm.

9.4 PDE Method

In the previous three sections, we discussed the development of methods for efficient model calibration and for the pricing of simple European options. In

the remainder of this chapter, we turn our attention to numerical techniques that allow a calibrated model to be used for pricing of general fixed income derivatives. We start out with the application of the PDE methods from Chapter 2.

9.4.1 PDE Formulation

The flexibility of the PDE method makes it applicable to a generalization of the specification (8.1)–(8.2) with a fully general time-dependent volatility function $\varphi(t, S)$. Let us therefore consider the following vector SDE

$$dS(t) = \varphi(t, S(t)) \sqrt{z(t)} dW(t), \quad (9.54)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \eta(t)\psi(z(t)) dZ(t), \quad (9.55)$$

where $\langle dZ(t), dW(t) \rangle = \rho dt$ and $z(0) = z_0$. Let $V(T)$ be an \mathcal{F}_T -measurable payoff and let $V(t, z, S)$ denote the numeraire-deflated value at time t , given $S(t) = S$ and $z(t) = z$, of a derivative that pays $V(T)$ at time T , $t \leq T$. By the usual arguments, $V(t, z, S)$ satisfies the following partial differential equation

$$\begin{aligned} 0 &= \frac{\partial}{\partial t} V(t, z, S) + \theta(z_0 - z) \frac{\partial}{\partial z} V(t, z, S) + \frac{\eta(t)^2}{2} \psi(z)^2 \frac{\partial^2}{\partial z^2} V(t, z, S) \\ &\quad + \frac{z}{2} \varphi(t, S)^2 \frac{\partial^2}{\partial S^2} V(t, z, S) + \rho \eta(t) \psi(z) \sqrt{z} \varphi(t, S) \frac{\partial^2}{\partial z \partial S} V(t, z, S). \end{aligned} \quad (9.56)$$

This PDE holds for $t \in [0, T]$ and $(S, z) \in \mathbb{R} \times \mathbb{R}^+$.

Fundamentally, (9.56) can be solved numerically by an application of the two-dimensional ADI scheme with a predictor-corrector step, as developed in Section 2.11.2. In an actual implementation of the ADI method, however, several issues in grid design and choice of boundary conditions must be addressed, a task to which we now turn.

9.4.2 Range for Stochastic Variance

Fixing a small probability $q_z > 0$, the range $[z_{\min}, z_{\max}]$ for z in the ADI grid can be set to cover the fraction $(1 - q_z)$ of the range of $z(T)$ in probability, i.e. from the conditions

$$P(z(T) < z_{\min}) = P(z(T) > z_{\max}) = q_z/2.$$

These probabilities are not known in closed form for $z(T)$ satisfying (9.55), so we will often have to resort to approximations. For instance, if ψ is not too different from a square root, we can replace

$$\psi(z) \rightarrow \frac{\psi(z_0)}{\sqrt{z_0}} \sqrt{z}, \quad (9.57)$$

to obtain a process of the square root type with time-dependent $\eta(t)$. From this representation, we can find an effective $\bar{\eta}$ to time horizon T by Theorem 9.3.6 and then apply the exact distribution of $z(T)$ with time-constant parameters from Proposition 8.3.2. Of course an even simpler, Gaussian, approximation is available if ψ is not too different from a constant.

A bit more crudely, but with less effort, we can also attempt to find the range for z from the stationary distribution of $z(t)$. When available, stationary distributions are a good source of approximations for tail probabilities — which is what we are interested in here — as we can often substitute large- z behavior with long-time behavior. The moments $E(z(T))$, $\text{Var}(z(T))$ of $z(T)$ that follows (9.55) are given by

$$E(z(T)) = z_0, \quad \text{Var}(z(T)) \approx \psi(z_0)^2 \int_0^T \eta(t)^2 e^{-2\theta(T-t)} dt,$$

where we have applied the approximation (9.57). Assuming that (9.57) is reasonable, the stationary distribution of $z(t)$ can be approximated with the Gamma distribution of Proposition 8.3.4; we choose the parameters of the Gamma distribution to match the mean and variance of $z(T)$,

$$\beta = \frac{E(z(T))}{\text{Var}(z(T))}, \quad \alpha = \beta E(z(T)).$$

The range of z in the ADI scheme can then be established by

$$z_{\min} = F^{-1}(q_z/2; \alpha, \beta), \quad z_{\max} = F^{-1}(1 - q_z/2; \alpha, \beta),$$

where $F(q; \alpha, \beta)$ is the Gamma CDF. Finally, we note that we can just use

$$z_{\min} = 0,$$

as long as we use one-sided discretization for boundary conditions at that point, as explained in Section 9.4.4 below.

9.4.3 Discretizing Stochastic Variance

Uniform discretization of z in the PDE (9.56) is rarely the best choice. If we look at the important case of $\psi(z) = \sqrt{z}$, assuming $z(0) = z_0 = 1$, the interval $[z_{\min}, z_{\max}]$ would be something like $[0, 10]$, with the mean of $z(t)$ being 1. Uniformly discretizing the range $[0, 10]$ would tend to put too few points in the interval $[0, 1]$, resulting in poor resolution in an important part of the range (see also Figure 9.2 in Section 9.5.3.1). To provide a remedy, we may recall the discussion in Section 7.4, which considered the transform

$$u(t) = \Psi(z(t)), \quad \Psi(z) = \int_{z_0}^z \frac{dy}{\psi(y)}. \quad (9.58)$$

Applying Ito's lemma, we get

$$\begin{aligned} du(t) &= \frac{1}{\psi(\Psi^{-1}(u(t)))} \\ &\quad \times \left(\theta(z_0 - \Psi^{-1}(u(t))) - \frac{1}{2} \frac{\psi'(\Psi^{-1}(u(t)))}{\psi(\Psi^{-1}(u(t)))} \eta(t)^2 \right) dt \\ &\quad + \eta(t) dZ(t). \end{aligned} \quad (9.59)$$

Noticing that the diffusion coefficient of $u(t)$ is not state-dependent, it appears reasonable to construct the grid in z -space from a uniform discretization in u . For this, suppose $N_z + 1$ points are used for the z -domain. We then define the grid $\{\zeta_n\}_{n=0}^{N_z}$ for z by the condition that $u_n \triangleq \Psi(\zeta_n)$ are spaced uniformly over $[\Psi(z_{\min}), \Psi(z_{\max})]$, so that

$$\begin{aligned} u_n &= \Psi(z_{\min}) + \frac{n}{N_z} (\Psi(z_{\max}) - \Psi(z_{\min})), \\ \zeta_n &= \Psi^{-1}(u_n) \\ &= \Psi^{-1}\left(\Psi(z_{\min}) + \frac{n}{N_z} (\Psi(z_{\max}) - \Psi(z_{\min}))\right), \quad n = 0, \dots, N_z. \end{aligned}$$

To give an example, consider the square root case $\psi(z) = \sqrt{z}$ where we have

$$\Psi(z) = \int_{z_0}^z \frac{dy}{\sqrt{y}} = 2(\sqrt{z} - \sqrt{z_0}), \quad \Psi^{-1}(u) = \left(\frac{u}{2} + \sqrt{z_0}\right)^2,$$

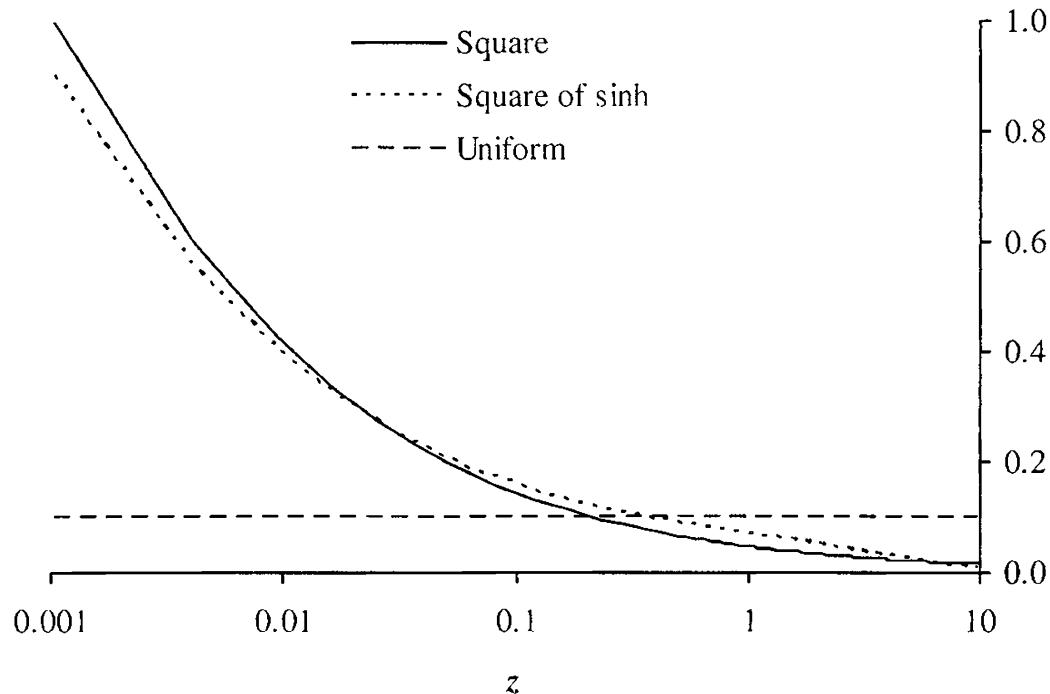
such that

$$\zeta_n = \left(\sqrt{z_{\min}} + \frac{n}{N_z} (\sqrt{z_{\max}} - \sqrt{z_{\min}}) \right)^2, \quad n = 0, \dots, N_z. \quad (9.60)$$

Empirically, it appears that concentrating points around the mean $z = z_0$ further improves numerical properties. We can achieve this effect by applying the sinh transform, see p. 167 of Tavella and Randall [2000], and then using (9.60):

$$\begin{aligned} \zeta_n &= \left(z_0 + \sinh\left(\alpha_{\min} + \frac{n}{N_z} (\alpha_{\max} - \alpha_{\min})\right) \right)^2, \\ \alpha_{\min, \max} &= \sinh^{-1}(\sqrt{z_{\min, \max}} - z_0). \end{aligned} \quad (9.61)$$

To illustrate the discretization strategies above, Figure 9.1 shows the density of grid points over $[z_{\min}, z_{\max}]$ using uniform discretization, quadratic discretization (9.60), and the sinh-quadratic discretization (9.61). As discussed, the quadratic and sinh-quadratic discretizations both increase the density of points in $(0, z_0]$, relative to a uniform discretization. In addition,

Fig. 9.1. Grid Density

Notes: Density of grid points (number of grid points per unit length) as a function of z for three different discretization schemes for z -domain: uniform, quadratic (9.60), and sinh-quadratic (9.61). We assume $z_{\min} = 0$, $z_0 = 1$, $z_{\max} = 10$. The abscissa axis is in logarithmic scale.

the sinh-quadratic scheme places more points around z_0 than does the quadratic scheme.

Let us finally note that instead of drawing on (9.58) as an inspiration for grid discretization in z , we could in principle use the variable u directly in the ADI scheme. Indeed, all that would be required is to rewrite (9.56) in terms of u , S and apply a uniform discretization to u . However, the drift of $u(t)$ is rather complicated and, importantly, grows to infinity as $u \rightarrow 0$ in the special case of $\psi(z) = \sqrt{z}$, see (9.59). A scheme that can handle large values of the drift robustly, such as the upwinding scheme from Section 2.6.1, would therefore be a necessity.

9.4.4 Boundary Conditions for Stochastic Variance

Practical experience shows that numerical schemes for solving the PDE (9.56) are quite robust with respect to the specifications of boundary conditions for z . Any reasonable choice from Chapter 2 appears to work well, including the standard $\partial^2 V / \partial z^2 = 0$ for $z = z_{\min}$, $z = z_{\max}$. In the case of $\psi(z) = \sqrt{z}$, if $z_{\min} = 0$, i.e. if we use $z = 0$ as the lower bound on the grid, for best results we should derive the boundary conditions for z_{\min} from the PDE itself, see Section 2.2.2. Setting $z = 0$ in (9.56) we obtain

$$0 = \frac{\partial}{\partial t} V(t, 0, S) + \theta z_0 \frac{\partial}{\partial z} V(t, 0, S), \quad (9.62)$$

a boundary condition of Neumann type. The validity of this boundary condition is intuitively justified by the fact that the solution to the SDE for $z(t)$ is unique, i.e. the behavior of $z(t)$ at the boundary $z = 0$ is determined by the SDE itself — and hence the boundary condition is determined by setting $z = 0$ in the PDE⁵. Incorporation of (9.62) into the finite difference solver would generally require one to discretize the z -derivative by one-sided differences; see Section 10.1.5.2 for details in a slightly more general setting.

Another, also reasonable, specification for the boundary $z = 0$ is obtained from the fact that the square-root process for $z(t)$ is strongly reflecting at $z = 0$, see Proposition 8.3.1. A reflection at the boundary translates into the boundary condition

$$\frac{\partial}{\partial z} V(t, 0, S) = 0$$

(see Karatzas and Shreve [1997]), which is quite similar to (9.62) and is another reasonable choice.

Interestingly, using the correct boundary conditions for the *forward* PDE, i.e. the forward Kolmogorov equation that the density of the process satisfies, is crucial, especially when the Feller condition (Proposition 8.3.1) is violated. As we have no use for forward PDEs for stochastic volatility processes in this book, we refer the reader to Lucic [2008] for the details.

9.4.5 Range for Underlying

To obtain the range

$$[S_{\min}, S_{\max}]$$

for the underlying S , we need to compute the approximate distribution of $S(T)$. Replacing the stochastic variance process with its expected value $E(z(t)) = z_0$, we obtain

$$dS(t) \approx \varphi(t, S(t)) \sqrt{z_0} dW(t).$$

To proceed, we can for example use the connection between option prices and the probability density, and apply various asymptotic results for local volatility models from Section 7.5. In the important special case of a time-dependent linear local volatility

$$\varphi(t, S) = \lambda(t) (b(t)S + (1 - b(t))L), \quad (9.63)$$

a reasonable approach is to replace time-dependent $b(t)$ with the effective time-independent skew \bar{b} via Proposition 9.3.4, and then apply a Gaussian approximation:

⁵A formal proof that (9.62) is theoretically correct, at least for payoffs that depend on z only (and not on S), is given in Ekström and Tysk [2008].

$$S(T) \approx [(\bar{b}S(0) + (1 - \bar{b})L) e^{\xi} - (1 - \bar{b})L] / \bar{b}, \quad (9.64)$$

$$\xi \sim \mathcal{N}\left(-\frac{z_0 \bar{b}^2}{2} \int_0^T \lambda(t)^2 dt, z_0 \bar{b}^2 \int_0^T \lambda(t)^2 dt\right).$$

As ξ is Gaussian, it is easy to find $[\xi_{\min}, \xi_{\max}]$ so that

$$P(\xi < \xi_{\min}) = P(\xi > \xi_{\max}) = q_S/2$$

for a given small probability $q_S > 0$. This trivially translates into the range for $S(T)$.

9.4.6 Discretizing the Underlying

The representation (9.64) proves useful for discretizing S as well. One approach is to discretize S so that the grid is uniform in ξ ,

$$S_n = [(\bar{b}S(0) + (1 - \bar{b})L) e^{\xi_n} - (1 - \bar{b})L] / \bar{b},$$

$$\xi_n = \xi_{\min} + \frac{n}{N_S} (\xi_{\max} - \xi_{\min}),$$

$$n = 0, \dots, N_S,$$

where N_S is the grid size. Alternatively, we can apply a transformation

$$y(S) = \ln \left(\frac{\bar{b}S + (1 - \bar{b})L}{\bar{b}S(0) + (1 - \bar{b})L} \right),$$

rewrite the PDE (9.56) in y instead of S , and discretize in y uniformly.

To conclude we note that even if $\varphi(t, S)$ is not of the form (9.63), we can always approximate it as such in order to compute the effective \bar{b} that is then used in discretization for S or in the mapping $S \rightarrow y$. Alternatively, we can always employ the same strategy (integral variable transform) that was advocated in Section 9.4.3 for z — which is what we used in Section 7.4 for discretizing local volatility models as well.

9.5 Monte Carlo Method

For generic stochastic volatility models such as (9.54)–(9.55), little can be said about Monte Carlo simulation that has not already been covered in Chapter 3. For any particular model parameterization, however, special-purpose discretization schemes can be constructed that have significant computational advantages over, say, the general-purpose Ito-Taylor schemes in Section 3.2.6. To demonstrate, we shall here specialize to the standard SV model, i.e. we consider the system

$$dS(t) = \lambda (bS(t) + (1 - b)L) \sqrt{z(t)} dW(t), \quad (9.65)$$

$$dz(t) = \theta(z_0 - z(t)) dt + \eta \sqrt{z(t)} dZ(t), \quad (9.66)$$

with $\langle dZ(t), dW(t) \rangle = \rho dt$ and $z(0) = z_0$. Our primary objective is to establish a scheme that allows us to time-discretize the SV model dynamics in an efficient manner; as it turns out, this is a surprisingly challenging, particularly for the z -process. We shall consequently deal with the Monte Carlo simulation of the SV model in a fairly careful manner, listing a number of schemes with different efficiency/bias trade-offs.

Remark 9.5.1. While we have assumed that parameters in the SV process are constants, all that is ultimately required is that parameters are piecewise constant on the simulation time line. As such, the schemes we suggest will also apply to time-dependent dynamics.

9.5.1 Exact Simulation of Variance Process

According to Proposition 8.3.2, the distribution of $z(t + \Delta)$ given $z(t)$ is known in closed form, and generation of a random sample of $z(t + \Delta)$ given $z(t)$ can be done entirely bias-free by sampling from a non-central chi-square distribution. Using the fact that a non-central chi-square distribution can be seen as a regular chi-square distribution with Poisson-distributed degrees of freedom (see Section 3.1.1.3), the following algorithm can be used.

1. Draw a Poisson random variable N , with mean $\frac{1}{2}z(t)n(t, t + \Delta)$ (here $n(t, T)$ is defined in (8.6)).
2. Given N , draw a regular chi-square random variable χ_v^2 , with $v = d + 2N$ degrees of freedom (d is defined in (8.6)).
3. Set $z(t + \Delta) = \chi_v^2 \cdot \exp(-\theta\Delta)/n(t, t + \Delta)$.

Steps 1 and 3 of this algorithm are straightforward, and Step 2 can be accomplished using the acceptance-rejection technique discussed in Section 3.1.1.2.

As mentioned in Section 3.1.1.3, if $d > 1$ it may be numerically advantageous to use a different algorithm, based on the relation

$$\chi_d'^2(\gamma) \stackrel{d}{=} (Z + \sqrt{\gamma})^2 + \chi_{d-1}^2, \quad d > 1, \quad (9.67)$$

where $\stackrel{d}{=}$ denotes equality in distribution, $\chi_d'^2(\gamma)$ is a non-central chi-square variable with d degrees of freedom and non-centrality parameter γ , and Z is an ordinary $\mathcal{N}(0, 1)$ Gaussian variable. We trust that the reader can complete the details on application of (9.67) in a simulation algorithm for $z(t + \Delta)$.

One might think that the existence of an exact simulation scheme for $z(t + \Delta)$ would settle once and for all the question of how to generate paths

of the square-root process. In practice, however, several complications may arise with the application of the algorithm above. Indeed, the scheme is quite complex compared with many standard SDE discretization schemes and may not fit smoothly into existing software architecture for SDE simulation routines. Also, computational speed may be an issue, and the application of acceptance-rejection sampling will potentially cause a “scrambling effect” when process parameters are perturbed⁶, resulting in poor convergence of numerically computed sensitivities, see Section 3.3. While caching techniques can be designed to overcome some of these issues, storage, look-up, and interpolation of such a cache pose their own challenges. Further, the basic scheme above provides no explicit link between the paths of the Brownian motion $Z(t)$ and that of $z(t)$, complicating applications in which, say, multiple correlated Brownian motions need to be advanced through time.

In light of the discussion above, it seems reasonable to also investigate the application of simpler simulation algorithms. These will typically exhibit a bias — in the sense discussed in Section 3.2.8 — for finite values of Δ , but convenience and speed may more than compensate for this, especially if the bias is small and easy to control by reduction of step size. We proceed to discuss several classes of such schemes.

9.5.2 Biased Taylor-Type Schemes for Variance Process

9.5.2.1 Euler Schemes

Going forward, let us use \hat{z} to denote a discrete-time (biased) approximation to z . A classical approach to simulating a path \hat{z} involves the application of Ito-Taylor expansions, suitably truncated, see Sections 3.2.3 and 3.2.6 for details. The simplest such scheme is the Euler scheme, a direct application of which would here give

$$\hat{z}(t + \Delta) = \hat{z}(t) + \theta(z_0 - \hat{z}(t))\Delta + \eta\sqrt{\hat{z}(t)}Z\sqrt{\Delta}, \quad (9.68)$$

where Z is a $\mathcal{N}(0, 1)$ Gaussian variable. One immediate (and fatal) problem with (9.68) is that the discrete process \hat{z} can become negative with non-zero probability. The first time this happens on a path, computation of $\sqrt{\hat{z}(t)}$ will be impossible and the time-stepping scheme will fail. To get around this problem, several remedies have been proposed in the literature, starting with the suggestion in Kloeden and Platen [2000] that one simply replace $\sqrt{\hat{z}(t)}$ in (9.68) with $\sqrt{|\hat{z}(t)|}$. Lord et al. [2006] review a number of similar “fixes” and conclude that the following works best:

$$\hat{z}(t + \Delta) = \hat{z}(t) + \theta(z_0 - \hat{z}(t)^+)\Delta + \eta\sqrt{\hat{z}(t)^+}Z\sqrt{\Delta}. \quad (9.69)$$

⁶After a perturbation of parameters, the number of rejected samples in the Monte Carlo trial will likely change.

In Lord et al. [2006] this scheme is denoted “full truncation”; its main characteristic is that the process for \hat{z} is allowed to go below zero, at which point \hat{z} becomes deterministic with an upward drift of θz_0 .

9.5.2.2 Higher-Order Schemes

The scheme (9.69) has first-order weak convergence, i.e. expectations of functions of \hat{z} will approach their true values as $O(\Delta)$. To improve convergence, it is tempting to apply a Milstein scheme (see Section 3.2.6.3), the most basic of which is

$$\hat{z}(t + \Delta) = \hat{z}(t) + \theta(z_0 - \hat{z}(t))\Delta + \eta\sqrt{\hat{z}(t)}Z\sqrt{\Delta} + \frac{1}{4}\eta^2\Delta(Z^2 - 1).$$

As was the case for (9.68), this scheme has a positive probability of generating negative values of \hat{z} and therefore cannot be used without suitable modifications. Kahl and Jäckel [2006] list several other Milstein-type schemes, some of which allow for a certain degree of control over the likelihood of generating negative values. One interesting variation is the *implicit Milstein scheme*, defined as

$$\hat{z}(t + \Delta) = \frac{\hat{z}(t) + \theta z_0\Delta + \eta\sqrt{\hat{z}(t)}Z\sqrt{\Delta} + \frac{1}{4}\eta^2\Delta(Z^2 - 1)}{1 + \theta\Delta}. \quad (9.70)$$

It is easy to verify that this discretization scheme will result in strictly positive paths for the z process if $4\theta z_0 > \eta^2$. For cases where this bound does not hold, it will be necessary to modify (9.70) to prevent problems with the computation of $\sqrt{\hat{z}(t)}$. For instance, whenever $\hat{z}(t)$ drops below zero, we could use (9.69) rather than (9.70).

Under certain sufficient regularity conditions, we have seen in Chapter 3 that Milstein schemes have second-order weak convergence. Due to the presence of a square root in (9.66), these sufficient conditions are violated here, and one should not expect (9.70) to have second-order convergence for all parameter values, even the ones that satisfy $4\theta z_0 > \eta^2$. Numerical tests of Milstein schemes for square-root processes can be found in Kahl and Jäckel [2006] and Glasserman [2004]; overall these schemes perform fairly well in benign parameter regimes, but are typically less robust than the Euler scheme.

9.5.3 Moment Matching Schemes for Variance Process

9.5.3.1 Log-normal Approximation

The simulation schemes introduced in Section 9.5.2 all suffer to various degrees from an inability to keep the path of z non-negative. One, rather obvious, way around this is to draw $\hat{z}(t + \Delta)$ from a user-selected probability

distribution that i) is reasonably close to the true distribution of $z(t + \Delta)$; and ii) is certain not to produce negative values⁷. To ensure that i) is satisfied, it is natural to select the parameters of the chosen distribution to match one or more of the true moments for $z(t + \Delta)$, conditional upon $z(t) = \hat{z}(t)$. For instance, if we assume that the true distribution of $z(t + \Delta)$ is well approximated by a log-normal distribution with parameters μ and σ^2 , we write (see Andersen and Brotherton-Ratcliffe [2005])

$$\hat{z}(t + \Delta) = e^{\mu + \sigma Z}, \quad (9.71)$$

where Z is a standard Gaussian random variable, and μ, σ are chosen to satisfy

$$e^{\mu + \frac{1}{2}\sigma^2} = E(z(t + \Delta)|z(t) = \hat{z}(t)), \quad (9.72)$$

$$e^{2(\mu + \frac{1}{2}\sigma^2)} (e^{\sigma^2} - 1) = \text{Var}(z(t + \Delta)|z(t) = \hat{z}(t)). \quad (9.73)$$

The results in Corollary 8.3.3 can be used to compute the right-hand sides of this system of equations, which can then easily be solved analytically for μ and σ .

As is the case for many other schemes, (9.71) works best if the Feller condition, as defined in Proposition 8.3.1, is satisfied. If not, the lower tail of the log-normal distribution is often too thin to capture the true distribution shape of $\hat{z}(t + \Delta)$ — see Figure 9.2 for an example.

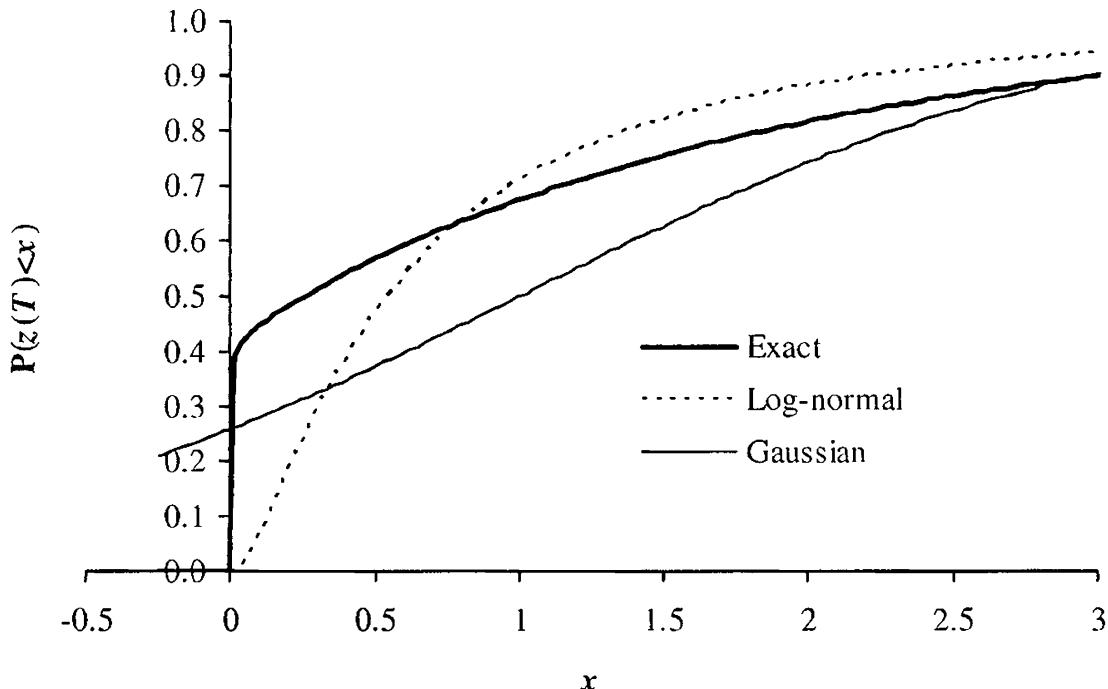
9.5.3.2 Truncated Gaussian

Figure 9.2 demonstrates that the density of $z(t + \Delta)|z(t)$ may sometimes be nearly singular at the origin. To accommodate this, one could contemplate inserting an actual singularity through outright truncation at the origin of a distribution that may otherwise go negative. Using a Gaussian distribution for this, say, one could write

$$\hat{z}(t + \Delta) = (\mu + \sigma Z)^+, \quad (9.74)$$

where μ and σ are determined by moment-matching, along the same lines as in Section 9.5.3.1 above. While this moment-matching exercise cannot be done in entirely analytical fashion, a number of caching tricks outlined in Andersen [2008] can be used to make the determination of μ and σ essentially instantaneous. As documented in Andersen [2008], the scheme

⁷As pointed out in Section 3.2.2, weak consistency — convergence of the first and second moments in the discretization scheme to those of the original SDE — is sufficient (together with some regularity conditions) for weak convergence. Hence, the actual distribution used for time-stepping can be chosen almost arbitrarily. Of course, matching other characteristics of the actual distribution may substantially improve the performance of the scheme.

Fig. 9.2. Cumulative Distribution of z 

Notes: The figure shows the cumulative distribution function for $z(T)$ given $z(0)$, with $T = 0.1$. Model parameters were $z(0) = z_0 = 1$, $\theta = 50\%$, and $\eta = 100\%$. The log-normal and Gaussian distributions in the graph were parameterized by matching mean and variances to the exact distribution of $z(T)$.

(9.74) is robust and generally has attractive convergence properties when applied to standard option pricing problems. Being fundamentally Gaussian when $\hat{z}(t)$ is far from the origin, (9.74) is qualitatively similar to the Euler scheme (9.69), although performance of (9.74) is typically somewhat better than (9.69). Unlike (9.69), the truncated Gaussian scheme (9.74) also ensures, by construction, that negative values of $\hat{z}(t + \Delta)$ cannot be attained.

9.5.3.3 Quadratic-Exponential

We finish our discussion of biased schemes for (9.66) with a more elaborate moment-matched scheme, based on a combination of a squared Gaussian and an exponential distribution. In this scheme, for large values of $\hat{z}(t)$, we write

$$\hat{z}(t + \Delta) = a(b + Z)^2, \quad (9.75)$$

where Z is a standard Gaussian random variable, and a and b are certain constants, to be determined by moment-matching. The constants a and b will depend on the time step Δ and $\hat{z}(t)$, as well as the parameters of the SDE for $z(t)$. While based on well-established asymptotics for the non-central chi-square distribution (see Andersen [2008]), formula (9.75) does not work well for low values of $\hat{z}(t)$ — in fact, the moment-matching exercise fails to

work — so we supplement it with a scheme to be used when $\hat{z}(t)$ is small. Examination of the true conditional density for $z(t + \Delta)|z(t)$ shows that the upper density tail decays exponentially, so a good choice is to approximate the distribution of $\hat{z}(t + \Delta)$ with

$$P(\hat{z}(t + \Delta) \in [x, x + dx]) = (p\delta(x) + \beta(1 - p)e^{-\beta x}) dx, \quad x \geq 0, \quad (9.76)$$

where δ is the Dirac delta function, and p and β are non-negative constants to be determined. As in the scheme in Section 9.5.3.2, we have a probability mass at the origin, but now the strength of this mass (p) is explicitly specified, rather than implied from other parameters. It can be verified that if $p \in [0, 1]$ and $\beta \geq 0$, then (9.76) constitutes a valid density function.

Assuming that we have determined a and b , Monte Carlo sampling from (9.75) is trivial. To draw samples in accordance with (9.76), we can generate a cumulative distribution function

$$\Psi(x) = P(\hat{z}(t + \Delta) \leq x) = p + (1 - p)(1 - e^{-\beta x}), \quad x \geq 0. \quad (9.77)$$

Here, the inverse of Ψ is readily computable:

$$\Psi^{-1}(u) = \Psi^{-1}(u; p, \beta) = \begin{cases} 0, & 0 \leq u \leq p, \\ \beta^{-1} \ln\left(\frac{1-p}{1-u}\right), & p < u < 1. \end{cases} \quad (9.78)$$

By the standard inverse distribution function method from Section 3.1.1.1, we thus get the simple sampling scheme

$$\hat{z}(t + \Delta) = \Psi^{-1}(U_z; p, \beta) \quad (9.79)$$

where U_z is a draw from a uniform distribution. Note that this scheme is extremely fast to execute.

Equations (9.75) and (9.79) together define the QE (for Quadratic-Exponential) discretization scheme. What remains is the determination of the constants a , b , p , and β , as well as a rule for when to switch from (9.75) to (9.79). The first problem is easily settled by moment-matching techniques, as shown in the following two propositions. We omit their straightforward proofs, which can be found in Andersen [2008].

Proposition 9.5.2. *Let*

$$m \triangleq E(z(t + \Delta)|z(t) = \hat{z}(t)), \quad s^2 \triangleq \text{Var}(z(t + \Delta)|z(t) = \hat{z}(t)),$$

and set $\psi = s^2/m^2$. Provided that $\psi \leq 2$, set

$$b^2 = 2\psi^{-1} - 1 + \sqrt{2\psi^{-1}}\sqrt{2\psi^{-1} - 1} \geq 0 \quad (9.80)$$

and

$$a = \frac{m}{1 + b^2}. \quad (9.81)$$

Let $\hat{z}(t + \Delta)$ be as defined in (9.75); then $E(\hat{z}(t + \Delta)) = m$ and $\text{Var}(\hat{z}(t + \Delta)) = s^2$.

Proposition 9.5.3. Let m , s , and ψ be as defined in Proposition 9.5.2. Assume that $\psi \geq 1$ and set

$$p = \frac{\psi - 1}{\psi + 1} \in [0, 1), \quad (9.82)$$

and

$$\beta = \frac{1 - p}{m} = \frac{2}{m(\psi + 1)} > 0. \quad (9.83)$$

Let $\hat{z}(t + \Delta)$ be sampled from (9.79); then $E(\hat{z}(t + \Delta)) = m$ and $\text{Var}(\hat{z}(t + \Delta)) = s^2$.

The terms m, s, ψ defined in the two propositions above are explicitly computable from the result in Corollary 8.3.3. For any ψ_c in $[1, 2]$, a valid *switching rule* is to use (9.75) if $\psi \leq \psi_c$ and to sample (9.77) otherwise. The exact value selected for ψ_c is non-critical; $\psi_c = 1.5$ is a natural choice.

9.5.3.4 Summary of QE Algorithm

As the QE algorithm is fairly complex, let us for convenience summarize the entire sampling algorithm step-by-step.

Assume that some arbitrary level $\psi_c \in [1, 2]$ has been selected. The detailed algorithm for the QE simulation step from $\hat{z}(t)$ to $\hat{z}(t + \Delta)$ is then:

1. Given $z(t) = \hat{z}(t)$, compute $m = E(z(t + \Delta)|z(t) = \hat{z}(t))$ and $s^2 = \text{Var}(z(t + \Delta)|z(t) = \hat{z}(t))$ from Corollary 8.3.3.
2. Compute $\psi = s^2/m^2$.
3. Draw a uniform random number U_z .
4. **If** $\psi \leq \psi_c$:
 - a) Compute a and b from equations (9.81) and (9.80).
 - b) Compute $Z = \Phi^{-1}(U_z)$.
 - c) Use (9.75), i.e. set $\hat{z}(t + \Delta) = a(b + Z)^2$.
5. **Otherwise**, if $\psi > \psi_c$:
 - a) Compute p and β according to equations (9.82) and (9.83).
 - b) Use (9.79), i.e. set $\hat{z}(t + \Delta) = \Psi^{-1}(U_z; p, \beta)$, where Ψ^{-1} is given in (9.78).

For efficiency, exponentials used in computation of m and s^2 should be pre-cached. The inversion of the Gaussian CDF in Step 4 can be done using the techniques described in Section 3.1.1.1.

The quadratic-exponential (QE) scheme outlined above is typically the most accurate of the biased schemes discussed here. Indeed, in most practical application the bias introduced by the scheme is statistically undetectable at the levels of Monte Carlo noise typically encountered in practical applications; see Andersen [2008] for numerical tests under a range of challenging conditions. Variations on the QE scheme without an explicit singularity in zero can also be found in Andersen [2008].

9.5.4 Broadie-Kaya Scheme for the Underlying

At this point, we are done discussing simulation schemes for the z -process, and now turn to the underlying process (9.65) itself.

For numerical work, it is useful to work with a logarithmic transformation of $S(t)$, rather than $S(t)$ itself. Specifically, we set

$$X(t) = \frac{bS(t) + (1-b)L}{bS(0) + (1-b)L},$$

the logarithm of which, from Proposition 8.3.6, satisfies the SDE

$$d \ln X(t) = -\frac{1}{2} \lambda^2 b^2 z(t) dt + \lambda b \sqrt{z(t)} dW(t). \quad (9.84)$$

As demonstrated in Broadie and Kaya [2006], it is possible to simulate (9.84) bias-free. To show this, first integrate the SDE for $z(t)$ in (9.66) and rearrange:

$$\int_t^{t+\Delta} \sqrt{z(u)} dZ(u) = \frac{1}{\eta} \left(z(t+\Delta) - z(t) - \theta z_0 \Delta + \theta \int_t^{t+\Delta} z(u) du \right). \quad (9.85)$$

Performing a Cholesky decomposition we can also write

$$d \ln X(t) = -\frac{1}{2} \lambda^2 b^2 z(t) dt + \lambda b \left(\rho \sqrt{z(t)} dZ(t) + \sqrt{1-\rho^2} \sqrt{z(t)} dB(t) \right),$$

where B is a Brownian motion independent of Z . An integration then yields

$$\begin{aligned} \ln X(t+\Delta) &= \ln X(t) + \frac{\rho \lambda b}{\eta} (z(t+\Delta) - z(t) - \theta z_0 \Delta) \\ &+ \left(\frac{\theta \rho \lambda b}{\eta} - \frac{\lambda^2 b^2}{2} \right) \int_t^{t+\Delta} z(u) du + \lambda b \sqrt{1-\rho^2} \int_t^{t+\Delta} \sqrt{z(u)} dB(u), \end{aligned} \quad (9.86)$$

where we have used (9.85). Conditional on $z(t+\Delta)$ and $\int_t^{t+\Delta} z(u) du$, it is clear that the distribution of $\ln X(t+\Delta)$ is Gaussian with easily computable moments. After first sampling $z(t+\Delta)$ bias-free from the non-central chi-square distribution (as described in Section 9.5.1), one then performs the following steps:

1. Conditional on $z(t+\Delta)$ (and $z(t)$) draw a bias-free sample of $I = \int_t^{t+\Delta} z(u) du$.
2. Conditional on $z(t+\Delta)$ and I , use (9.86) to draw a sample of $\ln X(t+\Delta)$ from a Gaussian distribution.

While execution of the second step is straightforward, the first one is decidedly not, as the conditional distribution of the integral I is not known in closed form. In Broadie and Kaya [2006], the authors instead derive a characteristic function, which they numerically Fourier-invert to generate the cumulative distribution function for I , given $z(t + \Delta)$ and $z(t)$. Numerical inversion of this distribution function over a uniform random variable finally allows for generation of a sample of I . The total algorithm requires great care in numerical discretization to prevent introduction of noticeable biases and is further complicated by the fact that the characteristic function for I contains two modified Bessel functions.

The Broadie-Kaya algorithm is bias-free by construction, but its complexity and lack of speed is problematic in many applications. Smith [2007] and Glasserman and Kim [2008] discuss various techniques to improve computational efficiency of the basic algorithm, but even with such improvements it is safe to say that the method is competitive only for applications that involve long time steps and require very high accuracy (and neither are the norm for fixed income applications).

9.5.5 Other Schemes for the Underlying

9.5.5.1 Taylor-Type Schemes

In their examination of “fixed” Euler-schemes, Lord et al. [2006] suggest simulation of the Heston model by combining (9.69) with the following scheme for $\ln X$:

$$\ln \widehat{X}(t + \Delta) = \ln \widehat{X}(t) - \frac{1}{2} \lambda^2 b^2 \widehat{z}(t)^+ \Delta + \lambda b \sqrt{\widehat{z}(t)^+} W \sqrt{\Delta}, \quad (9.87)$$

where W is a Gaussian $\mathcal{N}(0, 1)$ draw, correlated to Z in (9.69) with correlation coefficient ρ . For the periods where \widehat{z} drops below zero in (9.69), the process for \widehat{X} comes to a standstill.

Kahl and Jäckel [2006] examine the usage of Ito-Taylor expansions for joint simulation of $X(t)$ and $z(t)$, proposing several concrete schemes. As these schemes are rather complex, we simply refer the reader to Kahl and Jäckel [2006] for the details. Andersen [2008] tests the most prominent of the schemes in Kahl and Jäckel [2006] (the “IJK” scheme) and concludes that the scheme works well in benign parameter ranges, but has a tendency to deteriorate when parameters are made more extreme.

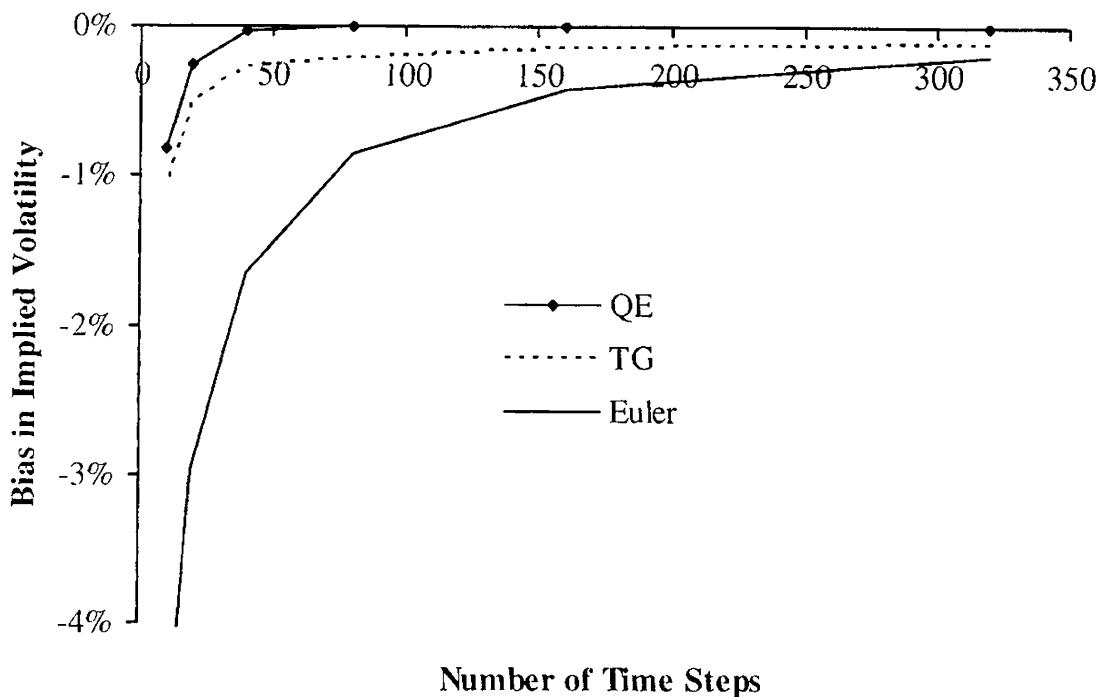
9.5.5.2 Simplified Broadie-Kaya

We recall from the discussion earlier that the complicated part of the Broadie-Kaya algorithm was the computation of $\int_t^{t+\Delta} z(u) du$, conditional on $z(t)$ and $z(t + \Delta)$. Andersen [2008] suggests a naive, but effective, approximation, based on the idea that

$$\int_t^{t+\Delta} z(u) du \approx \Delta [\gamma_1 z(t) + \gamma_2 z(t + \Delta)] , \quad (9.88)$$

for certain constants γ_1 and γ_2 . The constants γ_1 and γ_2 can be found by moment-matching techniques (using calculations similar to those from the proof of Theorem 9.3.6, or results from Dufresne [2001], p. 16), but Andersen [2008] presents evidence that it will often be sufficient to use either an Euler-like setting ($\gamma_1 = 1, \gamma_2 = 0$) or a central discretization ($\gamma_1 = \gamma_2 = \frac{1}{2}$). In any case, (9.88) combined with (9.86) gives rise to a scheme for Y -simulation that can be combined with any basic algorithm that can produce $\hat{z}(t)$ and $\hat{z}(t + \Delta)$. Andersen [2008] contains numerical results for the case where $\hat{z}(t)$ and $\hat{z}(t + \Delta)$ are simulated by the algorithms in Sections 9.5.3.2 and 9.5.3.3; results are excellent, particularly when the QE algorithm in Section 9.5.3.3 is used to sample \hat{z} . Figure 9.3 reproduces some sample convergence results from Andersen [2008].

Fig. 9.3. Convergence of Bias



Notes: The figure shows the convergence of the call option price bias in implied volatility terms, as a function of the number of time steps per path ($=T/\Delta$). The Euler scheme graph was computed using the full truncation scheme in (9.69), and the QE scheme used $\gamma_1 = \gamma_2 = 0.5$ and $\psi_c = 1.5$. Model parameters: $S(0) = L = 100$, $b = 1$, $z(0) = z_0 = 1$, $\theta = 0.5$, $\rho = -0.9$, $\eta = 1$, $\lambda = 20\%$. The option maturity is $T = 10$ and the strike is $K = 100$. The bias was estimated from 1,000,000 simulation paths, using the Fourier technique to establish exact prices.

9.5.5.3 Martingale Correction

Finally, let us note that some of the schemes outlined above, including the one in Section 9.5.5.2, will generally not lead to martingale behavior of \widehat{X} ; that is, $E(\widehat{X}(t + \Delta)|\widehat{X}(t)) \neq \widehat{X}(t)$. For the cases where the error $e = E(\widehat{X}(t + \Delta)|\widehat{X}(t)) - \widehat{X}(t)$ is analytically computable, it is, however, straightforward to remove the bias by simply adding $-e$ to the sample value for $\widehat{X}(t + \Delta)$. Andersen [2008] gives several examples of this idea and shows that, for the QE scheme at least, the improvements from martingale correction are minor.

9.A Appendix: Proof of Proposition 9.3.4

Let us fix a time horizon $T > 0$. Let $f(t, x)$ be a local volatility function,

$$f(t, x) \in C^1([0, T] \times \mathbb{R}),$$

satisfying the usual growth requirements. Let $\lambda(t)$, $t \in [0, T]$, be a function of time only. Fix $x_0 \in \mathbb{R}$. For any $\epsilon \geq 0$, define a rescaled local volatility function

$$f_\epsilon(t, x) = f(t\epsilon^2, x_0 + (x - x_0)\epsilon). \quad (9.89)$$

Without loss of generality we can assume that

$$f(t, x_0) \equiv 1, \quad t \in [0, T],$$

which implies

$$f_\epsilon(t, x_0) \equiv 1, \quad t \in [0, T]. \quad (9.90)$$

Let $w(t)$, $t \in [0, T]$, be a weight function such that

$$\int_0^T w(t) dt = 1, \quad (9.91)$$

and let us define an averaged local volatility function

$$\overline{f}_\epsilon(x)^2 = \int_0^T f_\epsilon(t, x)^2 w(t) dt. \quad (9.92)$$

Define two families of diffusions indexed by ϵ ,

$$\begin{aligned} dX_\epsilon(t) &= f_\epsilon(t, X_\epsilon(t)) \sqrt{z(t)} \lambda(t) dW(t), \quad X_\epsilon(0) = x_0, \\ dY_\epsilon(t) &= \overline{f}_\epsilon(Y_\epsilon(t)) \sqrt{z(t)} \lambda(t) dW(t), \quad Y_\epsilon(0) = x_0, \end{aligned}$$

for $t \in [0, T]$, where $z(t)$ is defined by (9.31). The following theorem can be found in Piterbarg [2005b].

Theorem 9.A.1. *If the weight function $w(t)$ is set to equal $w_T(t)$, where*

$$w_T(t) \triangleq \frac{v(t)^2 \lambda(t)^2}{\int_0^T v(t)^2 \lambda(t)^2 dt}, \quad (9.93)$$

$$v(t)^2 = \mathbb{E} \left(z(t) (X_0(t) - x_0)^2 \right),$$

then, as $\epsilon \rightarrow 0$,

$$\mathbb{E} \left((X_\epsilon(T) - x_0)^2 \right) - \mathbb{E} \left((Y_\epsilon(T) - x_0)^2 \right) = o(\epsilon^2), \quad (9.94)$$

$$\mathbb{E} \left((X_\epsilon(T) - x_0)^3 \right) - \mathbb{E} \left((Y_\epsilon(T) - x_0)^3 \right) = o(\epsilon^2). \quad (9.95)$$

Proof. The stochastic variance process $z(t)$ is Markovian. We denote its infinitesimal generator by L^z ,

$$L^z : \phi \mapsto \theta(z_0 - z) \frac{\partial \phi}{\partial z} + \frac{1}{2} \eta^2 z \frac{\partial^2 \phi}{\partial z^2}.$$

We note that the process $X_0(t)$ ($\equiv Y_0(t)$) satisfies the following SDE,

$$dX_0(t) = \sqrt{z(t)} \lambda(t) dW(t), \quad X_0(0) = x_0.$$

Let us denote the Markov semi-group of operators that corresponds to the process $(X_0(t), z(t))$ by $P_0(s, t)$, and the time-dependent infinitesimal generator by $L_0(t)$,

$$[P_0(s, t) \phi](x, z) = \mathbb{E}_s (\phi(X_0(t), z(t)) | X_0(s) = x, z(s) = z),$$

$$L_0(t) : \phi \mapsto \frac{1}{2} \lambda(t)^2 z \frac{\partial^2 \phi}{\partial x^2} + L^z.$$

Let us denote the same for $(X_\epsilon(t), z(t))$ and for $(Y_\epsilon(t), z(t))$ by $P_\epsilon^X(s, t)$, $L_\epsilon^X(t)$ and $P_\epsilon^Y(s, t)$, $L_\epsilon^Y(t)$, respectively.

From the general operator semigroup theory (see Ethier and Kurtz [1986]) it follows that

$$P_\epsilon^Y(0, T) = P_\epsilon^X(0, T) + \int_0^T P_\epsilon^Y(0, t) (L_\epsilon^Y(t) - L_\epsilon^X(t)) P_\epsilon^X(t, T) dt. \quad (9.96)$$

By Proposition 8.4.13 applied to $f(x) = (x - x_0)^2/2$ and $f(x) = (x - x_0)^3/6$,

$$\begin{aligned} \frac{1}{2} \mathbb{E} (X_\epsilon(T) - x_0)^2 &= \int_{-\infty}^{x_0} \mathbb{E} (K - X_\epsilon(T))^+ dK + \int_{x_0}^{\infty} \mathbb{E} (X_\epsilon(T) - K)^+ dK, \\ \frac{1}{6} \mathbb{E} (X_\epsilon(T) - x_0)^3 &= \int_{-\infty}^{x_0} (K - x_0) \mathbb{E} (K - X_\epsilon(T))^+ dK \\ &\quad + \int_{x_0}^{\infty} (K - x_0) \mathbb{E} (X_\epsilon(T) - K)^+ dK, \end{aligned}$$

and the same for Y_ϵ . Expressed in terms of the Markovian semigroup,

$$\frac{1}{(i+2)!} \mathbb{E} (X_\epsilon(T) - x_0)^{i+2} = \int_{-\infty}^{\infty} \langle \delta_{x_0, z_0}, (K - x_0)^i P_\epsilon^X(0, T) \pi_K \rangle dK,$$

(and the same for Y_ϵ) for $i = 0, 1$, where we have defined the payoff π_K by

$$\pi_K(x, z) = \begin{cases} (x - K)^+, & K \geq x_0, \\ (K - x)^+, & K < x_0. \end{cases}$$

Let us denote

$$\Delta(i) = \frac{1}{(i+2)!} \left(\mathbb{E} (Y_\epsilon(T) - x_0)^{i+2} - \mathbb{E} (X_\epsilon(T) - x_0)^{i+2} \right), \quad i = 0, 1.$$

To prove the theorem, we need to show that with the appropriate choice of weights $w_T(t)$,

$$\Delta(i) = o(\epsilon^2), \quad \epsilon \rightarrow 0, \quad i = 0, 1. \quad (9.97)$$

Clearly,

$$\Delta(i) = \int_{-\infty}^{\infty} (K - x_0)^i \langle \delta_{x_0, z_0}, (P_\epsilon^Y(0, T) - P_\epsilon^X(0, T)) \pi_K \rangle dK.$$

By (9.96) we have,

$$\begin{aligned} \Delta(i) &= \int_{-\infty}^{\infty} (K - x_0)^i \\ &\times \left(\int_0^T \langle \delta_{x_0, z_0}, P_\epsilon^Y(0, t) (L_\epsilon^Y(t) - L_\epsilon^X(t)) P_\epsilon^X(t, T) \pi_K \rangle dt \right) dK. \end{aligned}$$

After a series of manipulations (see Piterbarg [2005b] for details) we obtain, to order $o(\epsilon^2)$,

$$\Delta(i) = \frac{1}{2} \int_0^T \int \widehat{p}(x, z) (x - x_0)^i \left(\bar{f}_\epsilon(x)^2 - f_\epsilon(t, x)^2 \right) \lambda(t)^2 dx dt, \quad (9.98)$$

$$\widehat{p}(t, x) \triangleq \mathbb{E} (z(t) \delta(X_0(t) - x_0)).$$

Expanding f, \bar{f} to the first order around (s, x_0) , we obtain

$$\begin{aligned}
\delta(t; i) &\triangleq \int \widehat{p}(t, x) (x - x_0)^i \left(\overline{f}_\epsilon(x)^2 - f_\epsilon(t, x)^2 \right) dx \\
&= 2\epsilon \left(\frac{\partial f(s\epsilon^2, x_0)}{\partial x} - \int_0^T \frac{\partial f(s\epsilon^2, x_0)}{\partial x} w(s) ds \right) \\
&\quad \times \int \widehat{p}(t, x) (x - x_0)^{i+1} dx \\
&\quad + \epsilon^2 \left(\left[\frac{\partial f(s\epsilon^2, x_0)}{\partial x} \right]^2 - \int_0^T \left[\frac{\partial f(s\epsilon^2, x_0)}{\partial x} \right]^2 w(s) ds \right) \\
&\quad \times \int \widehat{p}(t, x) (x - x_0)^{i+2} dx \\
&\quad + o(\epsilon^2).
\end{aligned}$$

Calculating the integrals, we obtain to order $o(\epsilon^2)$,

$$\begin{aligned}
\delta(t; i) &= 2\epsilon v(t)^2 \left(\frac{\partial f(s\epsilon^2, x_0)}{\partial x} - \int_0^T \frac{\partial f(s\epsilon^2, x_0)}{\partial x} w(s) ds \right), \\
\Delta(i) &= \frac{1}{2} \int_0^T \delta(t; i) \lambda(t)^2 dt.
\end{aligned}$$

For $w(t) = w_T(t)$, we obtain $\Delta(i) = 0$, $i = 0, 1$, and the theorem follows. \square

Proposition 9.3.4 is proved by applying Theorem 9.A.1 to the equation (9.30). To compute $v(t)^2$, conditioning on $z(t)$ and using conditional independence of $X_0(t)$ and $z(t)$ we obtain,

$$\begin{aligned}
\mathbb{E}((X_0(t) - x_0)^2 z(t)) &= \mathbb{E}\left(z(t) \mathbb{E}\left((X_0(t) - x_0)^2 \mid z(\cdot)\right)\right) \quad (9.99) \\
&= \mathbb{E}\left(z(t) \int_0^t z(s) \lambda(s)^2 ds\right) \\
&= \int_0^t \lambda(s)^2 \mathbb{E}(z(t) z(s)) ds.
\end{aligned}$$

Clearly

$$z(t) - z_0 = e^{-\theta(t-s)} (z(s) - z_0) + O(dW),$$

so that

$$\begin{aligned}
\mathbb{E}(z(t) z(s)) &= z_0^2 + \mathbb{E}\left(e^{-\theta(t-s)} (z(s) - z_0) z(s)\right) \\
&= e^{-\theta(t-s)} \mathbb{E}(z(s)^2) + \left(1 - e^{-\theta(t-s)}\right) z_0^2.
\end{aligned}$$

We also have that

$$\mathbb{E}(z(s)^2) = z_0^2 + z_0 \eta^2 \frac{1 - e^{-2\theta s}}{2\theta}. \quad (9.100)$$

Substituting into (9.99) yields

$$v(t)^2 = \mathbb{E} \left((S(t) - x_0)^2 z(t) \right) \quad (9.101)$$

$$\begin{aligned} &= \int_0^t \lambda(s)^2 \left(e^{-\theta(t-s)} \mathbb{E}(z(s)^2) + \left(1 - e^{-\theta(t-s)}\right) z_0^2 \right) ds \\ &= \int_0^t \lambda(s)^2 \left(e^{-\theta(t-s)} z_0^2 + e^{-\theta(t-s)} z_0 \eta^2 \frac{1 - e^{-2\theta s}}{2\theta} \right. \\ &\quad \left. + z_0^2 \left(1 - e^{-\theta(t-s)}\right) \right) ds \\ &= \int_0^t \lambda(s)^2 \left(z_0^2 + z_0 \eta^2 e^{-\theta(t-s)} \frac{1 - e^{-2\theta s}}{2\theta} \right) ds \\ &= z_0^2 \int_0^t \lambda(s)^2 ds + z_0 \eta^2 e^{-\theta t} \int_0^t \lambda(s)^2 \frac{e^{\theta s} - e^{-\theta s}}{2\theta} ds. \end{aligned} \quad (9.102)$$

9.B Appendix: Coefficients for Asymptotic Expansion

Set $\Omega = \Omega_0 \bar{v}^{1/2} \tau^{1/2} + \Omega_1 \bar{v}^{3/2} \tau^{3/2}$ where $\bar{v} = \mu_{z \lambda^2}(0, z_0)/T$. Also define the easily computed quantities

$$\Omega_{mn} = \frac{\partial^m \Omega / \partial \bar{v}^m}{\partial^n \Omega / \partial \bar{v}^n}.$$

Then the expansion coefficients in Proposition 9.2.4 are given by

$$\alpha_0 = \tau^{-2} l_{1,2} \left(\Omega_{21} - \frac{1}{4} \Omega^2 \Omega_{10} \right), \quad \alpha_1 = \tau^{-2} l_{1,2} \Omega^{-2} \Omega_{10},$$

and

$$\begin{aligned} \beta_0 = & \tau^{-2} l_{2,2} \left(\Omega_{21} - \frac{1}{4} \Omega^2 \Omega_{10} \right) \\ & - \tau^{-3} l_{2,3} \left(\Omega_{31} - \Omega_{21}^2 - \frac{1}{4} \Omega^2 (\Omega_{20} + \Omega_{10}^2) + \left(\Omega_{21} - \frac{1}{4} \Omega^2 \Omega_{10} \right)^2 \right) \\ & + \frac{1}{2} \tau^{-4} l_{1,2}^2 \left(\Omega_{41} - 3\Omega_{31}\Omega_{21} + 2\Omega_{21}^2 - \frac{1}{4} \Omega^2 \Omega_{30} - \frac{3}{4} \Omega^2 \Omega_{10} \Omega_{20} \right) \\ & + \frac{3}{2} \tau^{-4} l_{1,2}^2 \left(\Omega_{21} - \frac{1}{4} \Omega^2 \Omega_{10} \right) \left(\Omega_{31} - \Omega_{21}^2 - \frac{1}{4} \Omega^2 (\Omega_{20} + \Omega_{10}^2) \right), \end{aligned}$$

$$\begin{aligned}
\beta_1 = & \Omega^{-2} \tau^{-2} l_{2,2} \Omega_{10} \\
& - \Omega^{-2} \tau^{-3} l_{2,3} \left(\Omega_{20} - 3\Omega_{10}^2 + 2\Omega_{10} \left(\Omega_{21} - \frac{1}{4}\Omega^2 \Omega_{10} \right)^2 \right) \\
& + \Omega^{-2} \frac{1}{2} \tau^{-4} l_{1,2}^2 \left(\Omega_{30} - 9\Omega_{10} \Omega_{20} + 12\Omega_{10}^3 \right. \\
& \quad \left. + 3\Omega_{10} \left(\Omega_{31} - \Omega_{21}^2 - \frac{1}{4}\Omega^2 (\Omega_{20} + \Omega_{10}^2) \right) \right) \\
& + \Omega^{-2} \frac{3}{2} \tau^{-4} l_{1,2}^2 \left(\Omega_{21} - \frac{1}{4}\Omega^2 \Omega_{10} \right) (\Omega_{20} - 3\Omega_{10}^2),
\end{aligned}$$

$$\beta_2 = -\Omega^{-4} \tau^{-3} l_{2,3} \Omega_{10}^2 + \Omega^{-4} \frac{3}{2} \tau^{-4} l_{1,2}^2 \Omega_{10} (\Omega_{20} - 3\Omega_{10}^2).$$

Part III

Term Structure Models

One-Factor Short Rate Models I

So far, our focus has been on vanilla models suitable for simple securities for which a change of measure allows the price to be expressed as an expectation of (a function of) a single random variable, typically a forward swap or Libor rate. However, many practically important securities, such as those that are callable or path-dependent, depend on interest rates in a substantially more complex manner, necessitating the construction of models for the dynamics of the *entire* discount curve — and not just a select few points on it. We have already, in Chapter 4, outlined the HJM theory that governs all dynamic discount curve models driven by vector-valued Brownian motions. The general HJM class with its infinite-dimensional Markovian dynamics is, however, too unwieldy to work with in practice, so it is of considerable interest to identify HJM model sub-classes that involve a finite number of Markov state variables only. We shall devote several chapters to this task, covering first the “classical” approach of writing down an explicit SDE for the short rate $r(t)$.

In our treatment of short rate models, we start out in this chapter with an in-depth analysis of the one-factor mean-reverting Gaussian model, providing a classical perspective on a model that we encountered in a modern HJM setting in Chapter 4. The chapter also covers the affine one-factor model, of which the Gaussian model is a special case. In Chapter 11, we generalize our discussion to arbitrary one-factor SDEs for the short rate, and finally, in Chapter 12, we introduce the class of multi-factor short rate models.

For derivatives pricing purposes, the short rate modeling approach has largely been superseded by newer approaches. Still, short rate models remain quite popular in empirical work, and a good understanding of these models provides a strong foundation for work with more sophisticated models.

10.1 The One-Factor Gaussian Short Rate Model

We recall that discount bond prices are given by the risk-neutral expectation

$$P(t, T) = \mathbb{E}_t^Q \left(e^{-\int_t^T r(u) du} \right), \quad (10.1)$$

so knowledge of the risk-neutral dynamics for $r(t)$ is in principle sufficient to compute time t discount bond prices to all maturities $T > t$. In practice, the expectation in (10.1) may, of course, not be computable in closed form, so to make short rate models operational in practice we must look for the sub-class of models where (10.1) is either analytically tractable or, at the very least, amenable to fast numerical methods.

One approach for which (10.1) becomes particularly tractable is to model the short rate as a Gaussian random variable. The resulting *Gaussian short rate* (GSR) model has a long and distinguished history in the financial literature. While our applications focus leaves us little room for historical ruminations, we shall make a slight concession here, by developing the GSR model progressively from the historically important — yet ultimately impractical — special case in Ho and Lee [1986]. Our development of the model will also initially progress by classical “bottoms-up” means, developing the dynamics of the forward curve from an SDE for the short rate, rather than the other way around. Besides providing some historical perspective, our style of presentation involves several generally applicable techniques and should give the reader additional intuition about the mechanics of the models involved.

10.1.1 The Ho-Lee Model

10.1.1.1 Notations and First Steps

Starting from the fundamental assumption that the short rate $r(t)$ is adapted to a single Brownian motion $W(t)$, the simplest possible dynamics we can imagine is the martingale process $r(t) = r(0) + \sigma_r W(t)$, or

$$dr(t) = \sigma_r dW(t), \quad (10.2)$$

where $\sigma_r > 0$ is a constant and $W(t)$ is a Brownian motion in the risk-neutral measure Q . From the basic risk-neutral pricing relationship (10.1), the time t discount bond maturing at time T then must have the price

$$P(t, T) = \mathbb{E}_t \left(e^{-\int_t^T r(u) du} \right) = \mathbb{E}_t \left(e^{-r(0)(T-t) - \sigma_r \int_t^T W(u) du} \right), \quad (10.3)$$

where $\mathbb{E}_t = \mathbb{E}_t^Q$ is the time t risk-neutral expectation operator.

Lemma 10.1.1. *If $r(t)$ follows (10.2) in the risk-neutral measure, then*

$$\mathbb{E}_t \left(e^{-\int_t^T r(u) du} \right) = \exp \left(-r(t)(T-t) + \frac{1}{6}\sigma_r^2(T-t)^3 \right).$$

Proof. We notice that

$$r(u) = r(t) + \int_t^u \sigma_r dW(s), \quad u > t,$$

so that

$$-\int_t^T r(u) du = -r(t)(T-t) - \sigma_r \int_t^T \int_t^u dW(s) du.$$

The order of integration can be changed by Fubini's theorem (see Duffie [2001]), such that

$$\int_t^T \int_t^u dW(s) du = \int_t^T \int_s^T du dW(s) = \int_t^T (T-s) dW(s).$$

By the Ito isometry, it then follows that $-\int_t^T r(u) du$ is Gaussian with mean $-r(t)(T-t)$ and variance

$$\text{Var}_t \left(\int_t^T r(u) du \right) = \sigma_r^2 \int_t^T (T-s)^2 ds = \frac{1}{3} \sigma_r^2 (T-t)^3.$$

The result of the lemma then follows from basic moment properties of log-normal variables, see e.g. (1.22). \square

Let us define a yield $y(t, T) = -\ln P(t, T)/(T-t)$, such that

$$y(t, T) = r(t) - \frac{1}{6} \sigma_r^2 (T-t)^2.$$

The yield curve shapes that can be produced by the simple model in (10.2) are rather primitive, as is evident from this expression. In particular, the yield curve is always downward-sloping in $T-t$ and $y_\infty = \lim_{T \rightarrow \infty} y(t, T) = -\infty$.

10.1.1.2 Fitting the Term Structure of Discount Bonds

The model presented above effectively has only two parameters — $r(0)$ and σ_r — with which one can attempt to fit the initial yield curve. It should be clear that this is insufficient to properly match observable discount bond prices, which effectively disqualifies the model from practical pricing applications. Fortunately, as realized in the paper Ho and Lee [1986], a remedy is quite straightforward¹: simply introduce a deterministic function $a(t)$ and alter the model to be

$$r(t) = r(0) + a(t) + \sigma_r W(t), \quad a(0) = 0, \quad (10.4)$$

¹The original paper by Ho and Lee was set exclusively in discrete time. The continuous-time version of the model developed here is, we feel, significantly more transparent.

such that

$$dr(t) = a'(t) dt + \sigma_r dW(t), \quad (10.5)$$

where $a'(t)$ is the first-order derivative of $a(t)$. To match the discount bond curve at time 0, $a(t)$ cannot be freely stipulated, but must be set as specified in Lemma 10.1.2 below.

Lemma 10.1.2. *Let $r(t)$ be given as in (10.4), and assume that discount bond prices at time 0, $P(0, T)$, are known for all $T > 0$. Set*

$$a(t) = f(0, t) - r(0) + \frac{1}{2}\sigma_r^2 t^2, \quad f(0, t) = -\frac{\partial \ln P(0, t)}{\partial t}.$$

Then, for any $T > 0$,

$$\mathbb{E} \left(e^{-\int_0^T r(u) du} \right) = P(0, T). \quad (10.6)$$

Proof. Applying Lemma 10.1.1, we get

$$\mathbb{E} \left(e^{-\int_0^t r(u) du} \right) = \exp \left(-r(0)t + \frac{1}{6}\sigma_r^2 t^3 \right) \times \exp \left(-\int_0^t a(u) du \right),$$

from which it follows that (10.6) is satisfied if

$$-\int_0^t a(u) du = \ln P(0, t) + r(0)t - \frac{1}{6}\sigma_r^2 t^3.$$

Taking derivatives with respect to t yields

$$a(t) = -\frac{\partial \ln P(0, t)}{\partial t} - r(0) + \frac{1}{2}\sigma_r^2 t^2 = f(0, t) - r(0) + \frac{1}{2}\sigma_r^2 t^2.$$

□

The model (10.4) with $a(t)$ set as in Lemma 10.1.2 is known as the *Ho-Lee* model. We characterize the model further in the following proposition.

Proposition 10.1.3. *In the Ho-Lee model, the risk-neutral process for $r(t)$ is*

$$dr(t) = \left(\frac{\partial f(0, t)}{\partial t} + \sigma_r^2 t \right) dt + \sigma_r dW(t), \quad (10.7)$$

and bond prices at time t can be reconstituted from $r(t)$ through the expression

$$P(t, T) = \frac{P(0, T)}{P(0, t)} \exp \left(- (r(t) - f(0, t)) (T - t) - \frac{1}{2}\sigma_r^2 t (T - t)^2 \right).$$

Proof. Equation (10.7) follows directly from (10.5) when $a(t)$ satisfies Lemma 10.1.2. To show the second part of the proposition, applying Lemma 10.1.1 to $r(t) - a(t)$ yields

$$\begin{aligned}
P(t, T) &= \exp \left(-(r(t) - a(t))(T - t) + \frac{1}{6} \sigma_r^2 (T - t)^3 \right) \\
&\quad \times \exp \left(- \int_t^T a(u) du \right) \\
&= \exp \left(-(r(t) - a(t))(T - t) + \frac{1}{6} \sigma_r^2 (T - t)^3 \right) \\
&\quad \times \exp \left(\ln P(0, T) + r(0)T - \frac{1}{6} \sigma_r^2 T^3 - \ln P(0, t) - r(0)t + \frac{1}{6} \sigma_r^2 t^3 \right) \\
&= \frac{P(0, T)}{P(0, t)} \exp \left(-(r(t) - a(t) - r(0))(T - t) + \frac{1}{2} \sigma_r^2 (Tt^2 - T^2t) \right).
\end{aligned}$$

In this expression $-a(t) - r(0) = -f(0, t) - \sigma_r^2 t^2 / 2$ from the definition of $a(t)$. The result follows. \square

10.1.1.3 Analysis and Comparison with HJM Approach

To gain a better understanding of the Ho-Lee model, let us examine the dynamics for bonds and forward rates implied by the model. From Proposition 10.1.3, we get

$$f(t, T) = -\frac{\partial \ln P(t, T)}{\partial T} = f(0, T) + r(t) - f(0, t) + \sigma_r^2 t (T - t) \quad (10.8)$$

and

$$df(t, T) = dr(t) + \left(\sigma_r^2 (T - 2t) - \frac{\partial f(0, t)}{\partial t} \right) dt = \sigma_r^2 (T - t) dt + \sigma_r dW(t). \quad (10.9)$$

In similar fashion, we get

$$dP(t, T)/P(t, T) = r(t) dt - \sigma_r (T - t) dW(t).$$

In the notations of Section 4.4, we have thus established that forward rate volatilities in the Ho-Lee model are $\sigma_f(t, T) = \sigma_r$ and discount bond volatilities are $\sigma_P(t, T) = \sigma_r(T - t)$. Due to the constancy of $\sigma_f(t, T)$, random perturbations of the forward curve from movement in the dW term will thus be parallel², in the sense that all points on the forward curve will move by identical amounts. Discount bond volatilities, on the other hand, approach

²Due to the presence of the T -dependent “convexity adjustment” term $\sigma_r^2(T - t)$ in the drift of the forward rate process, net forward curve movements are not perfectly parallel. Were this the case, it is well-known that the model would be arbitrageable.

zero in linear fashion as $t \rightarrow T$, reflecting the pull to par phenomenon discussed earlier in Chapter 4.

Setting aside for a moment the question about whether the Ho-Lee model is a reasonable representation of the real world, let us make a brief interlude to point out that we could, in fact, have specified the model directly as an HJM model with $\sigma_f(t, T) = \sigma_r$ and a single Brownian motion. The HJM result, Lemma 4.4.1, then immediately establishes the drift in the SDE for $f(t, T)$ to be

$$\mu_f(t, T) = \sigma_r \int_t^T \sigma_r du = \sigma_r^2 (T - t),$$

consistent with (10.9) above. Integrating this equation establishes (10.8), from which the discount bond reconstitution formula in Proposition (10.1.3) follows. To establish (10.7), we simply write $r(t) = f(t, t)$ and differentiate:

$$dr(t) = df(t, T)|_{T=t} + \left. \frac{\partial f(t, T)}{\partial T} \right|_{T=t} dt = \sigma_r dW(t) + \left(\frac{\partial f(0, t)}{\partial t} + \sigma_r^2 t \right) dt,$$

where the second equality uses (10.8)–(10.9). Notice that arriving at Proposition 10.1.3 in this manner did not involve evaluation of any expectations.

The Ho-Lee model has several drawbacks that disqualifies it for most, if not all, pricing applications. We list some of them below.

- The constancy of forward rate volatilities as a function of forward rate maturity $(T - t)$ is unrealistic: long-dated forward rates are less volatile than short-dated ones.
- The constancy of forward rate volatilities as a function of calendar time t gives the model time-stationary dynamics, but also results in the model having far too few degrees of freedom to allow for calibration to quoted option prices.
- Spot and forward interest rates are Gaussian and can therefore become negative, which is unrealistic.
- The model has only one driving Brownian motion and instantaneous moves of all forward rates are therefore perfectly correlated, contrary to empirical evidence.

The last objection is common for all one-factor short rate models and will disqualify these models for the pricing of options that have strong payoff dependency on non-parallel moves of the yield curve, e.g. spread options (see Chapter 17). The possibility of generating negative rates also cannot be helped unless we abandon the Gaussian setting (which we shall do later in this chapter), but we *can* address the problems associated with using constant forward rate volatility. We turn to this problem next.

10.1.2 The Mean-Reverting GSR Model

10.1.2.1 The Vasicek Model

Many empirical studies find that interest rates exhibit *mean reversion*, in the sense that if an interest rate is high by historical standards, it will most likely fall in the future (and vice versa if the interest rate is low). To model this phenomenon, Vasicek [1977] assumed that the short rate follows a one-factor *Ornstein-Uhlenbeck* process in the risk-neutral measure:

$$dr(t) = \kappa(\vartheta - r(t)) dt + \sigma_r dW(t), \quad (10.10)$$

where $\kappa, \vartheta, \sigma_r$ are positive constants. From results for the linear SDE in Section 1.6, it follows that the short rate can be written

$$r(t) = \vartheta + (r(0) - \vartheta)e^{-\kappa t} + \sigma_r \int_0^t e^{-\kappa(t-s)} dW(s). \quad (10.11)$$

It follows that $r(t)$ is a Gaussian random variable with moments

$$\mathbb{E}(r(t)) = \vartheta + (r(0) - \vartheta)e^{-\kappa t}, \quad (10.12)$$

$$\text{Var}(r(t)) = \frac{\sigma_r^2}{2\kappa} (1 - e^{-2\kappa t}). \quad (10.13)$$

As $t \rightarrow \infty$, the mean of the short rate approaches ϑ and the variance goes to $\sigma_r^2/(2\kappa)$. Accordingly, ϑ is often known as the *long-term level* (or sometimes the *mean reversion level*) of the short rate. The speed at which the short rate can be expected to revert to its long-term level is determined by κ , known as the *mean reversion speed*.

To establish a discount bond pricing formula in the Vasicek model, we use (10.11) to write

$$\begin{aligned} - \int_0^t r(u) du &= -\vartheta t - (r(0) - \vartheta)(1 - e^{-\kappa t})/\kappa \\ &\quad - \sigma_r \int_0^t \int_0^u e^{-\kappa(u-s)} dW(s) du. \end{aligned}$$

Clearly $-\int_0^t r(u) du$ is Gaussian, with mean

$$-\vartheta t - (r(0) - \vartheta)(1 - e^{-\kappa t})/\kappa.$$

To establish the variance, we follow the approach in Lemma 10.1.1 and reverse the order of integration in the stochastic integral, followed by an application of the Ito isometry. The result is

$$\begin{aligned} \text{Var}\left(\sigma_r \int_0^t \int_0^u e^{-\kappa(u-s)} dW(s) du\right) &= \sigma_r^2 \int_0^t e^{2\kappa s} \left(\int_s^t e^{-\kappa u} du\right)^2 ds \\ &= \frac{\sigma_r^2}{2\kappa^3} (-e^{-2\kappa t} + 4e^{-\kappa t} + 2t\kappa - 3). \end{aligned}$$

From the usual result for log-normal variables, it follows that discount bond prices in the Vasicek model can be computed as

$$\begin{aligned} P(0, t) &= \exp \left(E \left(- \int_0^t r(u) du \right) + \frac{1}{2} \text{Var} \left(- \int_0^t r(u) du \right) \right) \\ &= \exp \left(- \frac{1 - e^{-\kappa t}}{\kappa} r(0) - \vartheta t + \frac{\vartheta}{\kappa} (1 - e^{-\kappa t}) \right) \\ &\quad \times \exp \left(\frac{\sigma_r^2}{4\kappa^3} (-e^{-2\kappa t} + 4e^{-\kappa t} + 2t\kappa - 3) \right). \end{aligned}$$

More generally, we have the following proposition, the proof of which is straightforward.

Proposition 10.1.4. Define

$$\begin{aligned} B(t, T) &= \frac{1 - e^{-\kappa(T-t)}}{\kappa}, \\ A(t, T) &= \left(\vartheta - \frac{\sigma_r^2}{2\kappa^2} \right) (B(t, T) - (T-t)) - \frac{\sigma_r^2 B(t, T)^2}{4\kappa}. \end{aligned}$$

Then, in the Vasicek model (10.10),

$$P(t, T) = \exp(A(t, T) - B(t, T)r(t)).$$

As we did for the model in Section 10.1.1.1, define $y(t, T) = -\ln P(t, T)/(T-t)$, and notice that now a finite limit exists,

$$y_\infty = \lim_{T \rightarrow \infty} y(t, T) = \vartheta - \sigma_r^2 / (2\kappa^2).$$

In the Vasicek model, three different yield curve shapes are possible.

Lemma 10.1.5. Let $y(t, T) = -\ln P(t, T)/(T-t)$. Then

- If $r(t) > \vartheta$, then $y(t, T)$ decreases in $T-t$.
- If $r(t) < y_\infty - \sigma_r^2/(4\kappa^2)$, then $y(t, T)$ increases in $T-t$.
- Otherwise, $y(t, T)$ first increases in $T-t$ and then decreases (i.e. $y(t, T)$ is humped).

Proof. By straightforward, but tedious, calculus. \square

While this is certainly an improvement over the martingale model we encountered in Section 10.1.1.1, the Vasicek model is still not capable of fitting the observable yield curve accurately enough for pricing applications. It should be obvious that the way to solve this problem is to mimic the step that lead to the Ho-Lee model in Section 10.1.1.2: introduce a deterministic function of time into the definition (10.11). That is, we write

$$r(t) = a(t) + \vartheta + (r(0) - \vartheta)e^{-\kappa t} + \sigma_r \int_0^t e^{-\kappa(t-s)} dW(s) = a(t) + r_{OU}(t),$$

where $a(t)$ is a deterministic function and $r_{OU}(t)$ is the short rate in the Vasicek model. The function $a(t)$ is determined from the condition that

$$\mathbb{E} \left(e^{-\int_0^t r_{OU}(u) du} \right) e^{-\int_0^t a(u) du} = P(0, t),$$

where the right-hand side is assumed given. Further development of this model proceeds as in Section 10.1.1.2, and results are easily imagined; we skip the analysis as the resulting model is a special case of the more general setup in Section 10.1.2.2 below. We do note, however, that the Vasicek model — both with and without adjustment to fit the initial yield curve — is easily shown to have forward rate and discount bond volatilities of

$$\sigma_f(t, T) = \sigma_r e^{-\kappa(T-t)}, \quad \sigma_P(t, T) = \sigma_r \left(\frac{1 - e^{-\kappa(T-t)}}{\kappa} \right).$$

Introduction of mean reversion into the model will thus introduce exponential decay in the term structure of forward rate volatilities. From an empirical standpoint this is considerably more appealing than the maturity-independent forward rate volatilities in the Ho-Lee model, and also in qualitative agreement with the fact that short- and medium-maturity interest rate options trade at higher implied volatilities³ than do long-dated options. While this is a step up from the Ho-Lee model, the model still has too few degrees of freedom for many derivatives pricing applications, as the model will rarely calibrate well to observed prices of vanilla options (e.g. European swaptions and caps). We improve on this in the next section.

10.1.2.2 The General One-Factor GSR Model

The most general form of the one-factor GSR model is given by the SDE

$$dr(t) = \kappa(t) (\vartheta(t) - r(t)) dt + \sigma_r(t) dW(t), \quad (10.14)$$

i.e. we have now allowed all parameters in the Vasicek model to depend on time. While this model can be developed by classical means (see e.g. Hull and White [1994a] for, often laborious, details), it is significantly easier to work within an HJM setting. In fact, we already showed in Section 4.5.2 that short rate dynamics of the form in (10.14) must originate from a “separable” HJM model of the form

$$df(t, T) = \sigma_f(t, T) \left(\int_t^T \sigma_f(t, u) du \right) dt + \sigma_f(t, T) dW(t), \quad (10.15)$$

$$\sigma_f(t, T) = \sigma_r(t) \exp \left(- \int_t^T \kappa(u) du \right).$$

Chapter 4 also proved the following result for the function $\vartheta(t)$.

³An exception to this observation is the humped volatility term structure that can often be observed in caplet markets. We return to this issue in Section 10.1.2.3.

Proposition 10.1.6. *For the general one-factor GSR model (10.14) to match the initial yield curve, we must have*

$$\vartheta(t) = \frac{1}{\varkappa(t)} \frac{\partial f(0, t)}{\partial t} + f(0, t) + \frac{1}{\varkappa(t)} \int_0^t e^{-2 \int_u^t \varkappa(s) ds} \sigma_r(u)^2 du.$$

Proof. Follows from Proposition 4.5.4, when $d = 1$. \square

We notice the presence of $\partial f(0, t)/\partial t$ in the expression for $\vartheta(t)$ (a similar term was, of course, present in the Ho-Lee model) which can be a nuisance in applications where the initial forward curve is not smooth, as when we have used simple bootstrapping to construct the curve. To get rid of the term, we now switch variables, from $r(t)$ itself to $x(t) = r(t) - f(0, t)$. Dynamics for $x(t)$, as well as the bond reconstitution formula for (10.14) in terms of $x(t)$ are listed next.

Proposition 10.1.7. *Define*

$$x(t) \triangleq r(t) - f(0, t).$$

Then, for the model (10.14)–(10.15),

$$dx(t) = (y(t) - \varkappa(t)x(t)) dt + \sigma_r(t) dW(t), \quad x(0) = 0, \quad (10.16)$$

where

$$y(t) = \int_0^t e^{-2 \int_u^t \varkappa(s) ds} \sigma_r(u)^2 du. \quad (10.17)$$

The bond reconstitution formula is

$$P(t, T) = \frac{P(0, T)}{P(0, t)} \exp \left(-x(t)G(t, T) - \frac{1}{2}y(t)G(t, T)^2 \right), \quad (10.18)$$

$$G(t, T) = \int_t^T e^{-\int_t^u \varkappa(s) ds} du.$$

Proof. To simplify notation, define $K(t) = \int_0^t \varkappa(u) du$, and set $g(t) = \sigma_r(t)e^{K(t)}$, $h(t) = e^{-K(t)}$. Then $\sigma_f(t, T) = g(t)h(T)$ and, by integration of (10.15),

$$f(t, T) = f(0, T) + h(T) \int_0^t g(u)^2 \int_u^T h(s) ds du + h(T) \int_0^t g(u) dW(u). \quad (10.19)$$

Set

$$x(t) = h(t) \int_0^t g(u)^2 \int_u^t h(s) ds du + h(t) \int_0^t g(u) dW(u),$$

and note that, by the Leibniz rule for differentiation of an integral,

$$\begin{aligned}
dx(t) &= h'(t) \left(\int_0^t g(u)^2 \int_u^t h(s) ds du \right) dt + h(t)^2 \left(\int_0^t g(u)^2 du \right) dt \\
&\quad + h(t)g(t) dW(t) + h'(t) \int_0^t g(u) dW(u) dt \\
&= \left(\frac{h'(t)}{h(t)} x(t) + y(t) \right) dt + h(t)g(t) dW(t) \\
&= (y(t) - \varkappa(t)x(t)) dt + \sigma_r(t) dW(t),
\end{aligned}$$

where

$$y(t) = h(t)^2 \int_0^t g(u)^2 du$$

was defined in (10.17). From (10.19) we have

$$\begin{aligned}
f(t, T) &= f(0, T) + \frac{h(T)}{h(t)} x(t) + h(T) \int_0^t g(u)^2 \int_u^T h(s) ds du \\
&\quad - h(T) \int_0^t g(u)^2 \int_u^t h(s) ds du \\
&= f(0, T) + \frac{h(T)}{h(t)} x(t) + h(T) \int_t^T h(s) ds \int_0^t g(u)^2 du \\
&= f(0, T) + \frac{h(T)}{h(t)} \left(x(t) + \frac{y(t)}{h(t)} \int_t^T h(s) ds \right),
\end{aligned}$$

such that in particular $r(t) = f(t, t) = f(0, t) + x(t)$, as claimed earlier. Inserting the expression for $f(t, T)$ into the basic relation

$$P(t, T) = \exp \left(- \int_t^T f(t, u) du \right)$$

produces (10.18) after a few rearrangements. \square

Remark 10.1.8. The discount bond dynamics for $P(t, T)$ are

$$dP(t, T)/P(t, T) = r(t) dt - \sigma_P(t, T) dW(t), \quad \sigma_P(t, T) = \sigma_r(t) G(t, T).$$

Remark 10.1.9. In the reconstitution formula (10.18), notice that

$$G(t, T) = (G(0, T) - G(0, t)) e^{\int_0^t \varkappa(s) ds},$$

a result that is often useful in grid-based numerical work (see Section 10.1.5).

Proposition 10.1.7 is an important result and shall serve as the foundation for most of the remaining discussion of Gaussian short rate models.

10.1.2.3 Time-Stationarity and Caplet Hump

A Gaussian HJM model is said to be *time-stationary* if the instantaneous volatility $\sigma_f(t, T)$ is only a function of $T - t$, i.e. the time *to* maturity rather than the time *of* maturity T . Time stationarity is an appealing feature, as it implies that the volatility term structure of forward rates will look the same in the future as it does today; in the absence of other information, this prediction is often very reasonable and in good agreement with empirical observation. In the setting of the one-factor GSR model, imposing time-stationarity will require us to set both $\sigma_r(t)$ and $\kappa(t)$ to constants, such that

$$\sigma_f(t, T) = \sigma_r e^{-\kappa(T-t)}. \quad (10.20)$$

In other words, the only time-stationary forward rate volatility term structure that can be constructed in the GSR model is an exponentially decaying one. In practice, however, it is quite common to observe (from the caplet market, say) forward rate volatility structures that have a marked “hump”, with short-dated options trading at very low volatilities. This effect can largely be attributed to central bank activity, as the extreme short end of the forward curve tends to move primarily in response to central bank changes to funding rates. As such changes are relatively infrequent and normally quite predictable⁴, short-dated forward rates are typically associated with relatively little uncertainty and, consequently, have low volatilities.

If we attempt to match a GSR model to a humped forward volatility structure, it follows from (10.20) that this cannot be done in a stationary manner and we are forced to let κ become a function of time. To see this, suppose that we at time 0 observe forward volatilities $\sigma_f(0, T) = b(T)$, where $b(T)$ is a humped function of T , i.e. $b(T)$ initially increases in T but ultimately decreases in T . Ideally, we would like to set $\sigma_f(t, T) = b(T - t)$, but this is not possible in the GSR setting, as explained above. To make the GSR model match $b(T)$ at time 0, we instead are forced to make κ a function of time, determined from the relation

$$\sigma_f(0, T) = \sigma_r e^{-\int_0^T \kappa(u) du} = b(T), \quad \sigma_r = b(0).$$

Taking logarithms and differentiating gives

$$\kappa(t) = -\frac{d(\ln b(t))}{dt} = -\frac{b'(t)}{b(t)}.$$

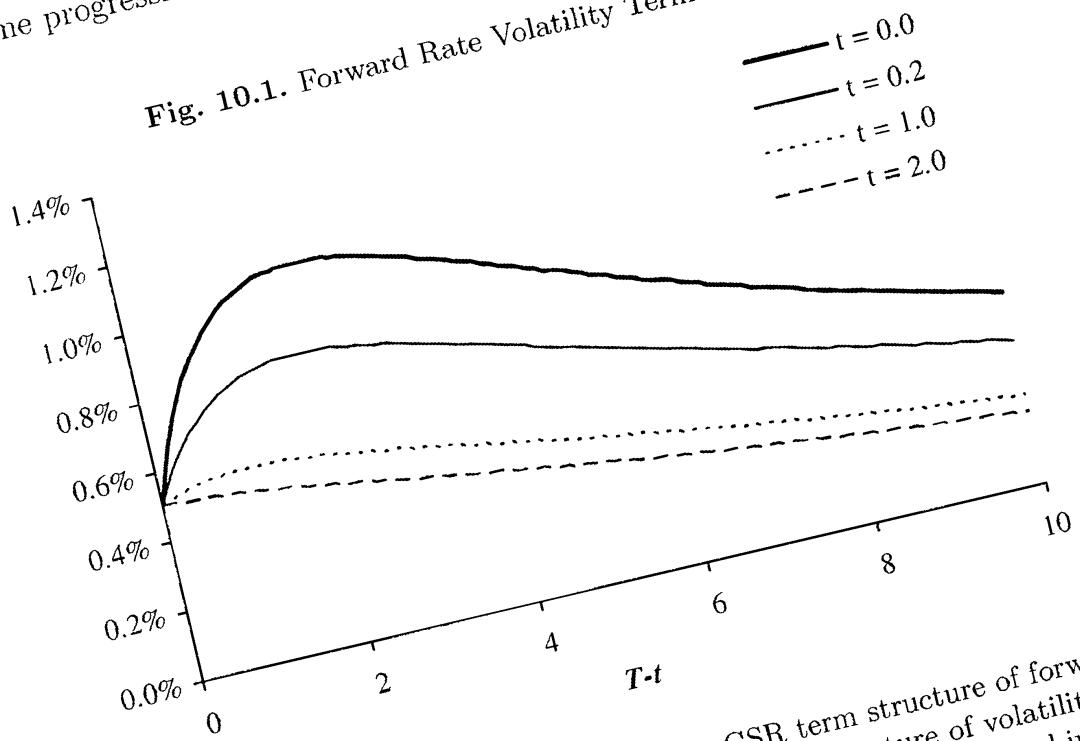
If t_p is the time t at which $b(t)$ reaches its peak (i.e. $b'(t_p) = 0$), it follows that $\kappa(t)$ will be negative for all $t < t_p$ and positive for all $t > t_p$. At time $t > 0$, our so-calibrated GSR model will produce instantaneous forward volatilities of

⁴On occasion there is significant uncertainty in the market about the intentions of monetary authorities, in which case the caplet hump may disappear temporarily.

$$\sigma_f(t, T) = \sigma_r e^{-\int_t^T \kappa(u) du} = \sigma_r b(T)/b(t).$$

Clearly $\sigma_f(t, T)$ is not stationary. In fact, once $t > t_p$ the model no longer produces a hump at all, as $b(T)/b(t)$ is monotonically decaying in $T - t$ for $t > t_p$. Figure 10.1 demonstrates this effect; notice also how volatilities become progressively lower as time t moves forward.

Fig. 10.1. Forward Rate Volatility Term Structure



Notes: The figure shows the evolution of the GSR term structure of forward rate volatility with t . The model fit to the original term structure of volatility ($t = 0$) was done solely through the mean reversion function $\kappa(t)$, as described in the text.

The lesson to be learned from the examination above is essentially this: in the GSR model and, in fact, in all short rate models, one should resist making κ a function of time lest one is willing to accept strongly non-stationary evolution of forward rate volatilities. Yet, working with perfectly time-stationary GSR models is often too constraining in practical applications, as the resulting model has too few degrees of freedom in its volatility characteristics to calibrate against observed vanilla option prices in a meaningful way. Our recommendation for most applications is to freeze κ at a constant value, but to allow σ_r to be a function of time (see Section 13.1.8 for much more on mean reversion calibration). That is, we would set

$$\sigma_f(t, T) = \sigma_r(t) e^{-\kappa(T-t)}.$$

While the resulting model is not time-stationary, it retains through time a persistent exponential shape of its instantaneous forward volatility structure.

The reader may at this point reasonably ask whether models exist that can produce a time-stationary hump in instantaneous forward volatilities. The answer is yes, but such models would generally need more than a single Markov variable to characterize moves in the yield curve. We return to this issue in Chapter 12 and, indeed, in many later chapters on multi-factor models.

10.1.3 European Option Pricing

In the general one-factor GSR model (10.14), suppose that we fix the mean reversion function $\kappa(t)$ exogenously, e.g. based on empirical observations or from observation of typical decay speed of implied volatilities with option maturity⁵. The function $\vartheta(t)$ in (10.14) is then uniquely fixed by the initial forward curve, so to complete the specification of the model it remains to determine the function $\sigma_r(t)$. In pricing applications, this function is normally found by calibration of the GSR model to observed prices of liquid European options, such as caps and swaptions. While we shall postpone most of the intricacies of volatility calibration to later chapters, it should be clear that for a calibration to caps and swaptions to be efficient, we need computationally efficient methods for the valuation of these instruments.

In Section 4.5.1, we showed that for any Gaussian HJM model — whether the short rate is Markov or not — caplets can be priced by simple Black-Scholes formulas; see Proposition 4.5.2 for the details. Consequently, we here focus our attention on the pricing of swaptions. For concreteness, consider a payer swaption expiring at time T_0 , with the underlying swap paying an annualized coupon c at times $T_1 < T_2 < \dots < T_N$, with $T_1 > T_0$. We recall from Chapter 5 that the swaption payout at time T_0 is

$$V_{\text{swaption}}(T_0) = \left(1 - P(T_0, T_N) - c \sum_{i=0}^{N-1} \tau_i P(T_0, T_{i+1}) \right)^+, \quad \tau_i = T_{i+1} - T_i. \quad (10.21)$$

10.1.3.1 The Jamshidian Decomposition

Our first approach is exact, and is based on a method developed by Jamshidian [1989]. The basic idea is to rewrite the swaption payout from an option on a sum of discount bonds to a sum of options on discount bonds. To develop the idea in detail, let us write $P(T_0, T_N) = P(T_0, T_N, x(T_0))$ to recognize the dependence of $P(T_0, T_N)$ on $x(T_0) = r(T_0) - f(0, T_0)$ through the reconstitution formula (10.18). We also define a “critical” value x^* for

⁵As argued above, normally we would pick $\kappa(t)$ to be a constant. A more detailed examination of the estimation of mean reversions — and the role it plays in Bermudan swaption pricing — can be found in Chapter 13.

which the swap at time T_0 is exactly zero; x^* can be found by numerical root search on the equation

$$P(T_0, T_N, x^*) + c \sum_{i=0}^{N-1} \tau_i P(T_0, T_{i+1}, x^*) = 1. \quad (10.22)$$

Finally, define “strikes”

$$K_i = P(T_0, T_i, x^*), \quad i = 1, \dots, N;$$

it follows that

$$K_N + c \sum_{i=0}^{N-1} \tau_i K_{i+1} = 1. \quad (10.23)$$

We are now ready to apply the Jamshidian “trick”. Inspection of (10.18) shows that all zero-coupon bonds $P(T_0, T_i, x(T_0))$ are monotonically decreasing in $x(T_0)$, whereby the swaption only pays out a positive amount if $x(T_0) > x^*$. That is,

$$\begin{aligned} V_{\text{swaption}}(T_0) &= \\ &= \left(1 - P(T_0, T_N, x(T_0)) - c \sum_{i=0}^{N-1} \tau_i P(T_0, T_{i+1}, x(T_0)) \right) \mathbf{1}_{\{x(T_0) > x^*\}} \\ &= \left(K_N + c \sum_{i=0}^{N-1} \tau_i K_{i+1} - P(T_0, T_N, x(T_0)) \right. \\ &\quad \left. - c \sum_{i=0}^{N-1} \tau_i P(T_0, T_{i+1}, x(T_0)) \right) \mathbf{1}_{\{x(T_0) > x^*\}}, \end{aligned}$$

where the second equality follows from (10.23). Thus,

$$\begin{aligned} V_{\text{swaption}}(T_0) &= (K_N - P(T_0, T_N, x(T_0))) \mathbf{1}_{\{x(T_0) > x^*\}} \\ &\quad + c \sum_{i=0}^{N-1} \tau_i (K_{i+1} - P(T_0, T_{i+1}, x(T_0))) \mathbf{1}_{\{x(T_0) > x^*\}} \\ &= (K_N - P(T_0, T_N, x(T_0)))^+ \\ &\quad + c \sum_{i=0}^{N-1} \tau_i (K_{i+1} - P(T_0, T_{i+1}, x(T_0)))^+, \end{aligned} \quad (10.24)$$

where we used monotonicity of $P(T_0, T_i, x(T_0))$ on the last step. With this result, we have decomposed the swaption payout into $N + 1$ put options on zero-coupon bonds. Such options are easily valued using the formula from Proposition 4.5.1, allowing us to price the swaption in closed form.

10.1.3.2 Gaussian Swap Rate Approximation

While the Jamshidian approach above is perfectly adequate for many applications, its use of numerical root search and the need to price a potentially large amount of zero-coupon options can be cumbersome. One may wonder, then, whether perhaps a simpler approach is possible, given the simplicity of the dynamics of rates in the GSR framework. One obvious idea is to examine the SDE of the forward swap rate in an appropriate annuity measure, introducing approximations as needed to make the dynamics tractable. This idea shall be used many times in this book, often in combination with sophisticated techniques for simplification of the swap rate SDEs. Here, we have more modest aspirations and will be content with a simpler — yet still functional — approach. The reader shall consider this section a warm-up exercise for more accurate approximations to come, in particular in Sections 13.1.4 and 13.1.5 that also cover the GSR case.

We start by rewriting the swaption payout as

$$V_{\text{swaption}}(T_0) = A(T_0)(S(T_0) - c)^+,$$

where $A(t)$ and $S(t)$ are the swap annuity and forward swap rate, respectively, see (5.13)–(5.14):

$$A(t) \triangleq A_{0,N}(t) = \sum_{i=0}^{N-1} \tau_i P(t, T_{i+1}), \quad S(t) \triangleq S_{0,N}(t) = \frac{P(t, T_0) - P(t, T_N)}{A(t)}.$$

Let Q^A be the measure induced by using $A(t)$ as the numeraire, such that

$$V_{\text{swaption}}(0) = A(0)E^A((S(T_0) - c)^+), \quad (10.25)$$

where E^A denotes expectation in measure Q^A . To evaluate (10.25), we need to determine the dynamics of $S(t)$ in Q^A . Lemma 4.2.4 establishes that $S(t)$ is a martingale under Q^A . From the reconstitution formula (10.18) we also know that $S(t)$ and $A(t)$ must be deterministic functions of $x(t)$:

$$S(t) = S(t, x(t)), \quad A(t) = A(t, x(t)),$$

so from Ito's lemma

$$dS(t) = q(t, x(t)) \sigma_r(t) dW^A(t), \quad q(t, x) = \frac{\partial}{\partial x} \frac{P(t, T_0, x) - P(t, T_N, x)}{\sum_{i=0}^{N-1} \tau_i P(t, T_{i+1}, x)},$$

where W^A is a Q^A -Brownian motion and where we use (10.18) to express the $P(t, T_i)$'s as functions of x . Evaluating the partial derivatives yields

$$\begin{aligned} q(t, x) &= - \frac{P(t, T_0, x)G(t, T_0) - P(t, T_N, x)G(t, T_N)}{A(t, x)} \\ &\quad + \frac{S(t, x)}{A(t, x)} \sum_{i=0}^{N-1} \tau_i P(t, T_{i+1}, x)G(t, T_{i+1}), \end{aligned} \quad (10.26)$$

where we recall that

$$G(t, T) = \int_t^T e^{-\int_t^u \kappa(s) ds} du.$$

The function $q(t, x)$ can be experimentally verified to be close to a constant in x -direction so, as a good approximation, we can write

$$q(t, x(t)) \approx q(t, \bar{x}(t)), \quad (10.27)$$

where $\bar{x}(t)$ is some *deterministic* proxy for $x(t)$. With this, the option formula in the Normal model, see Remark 7.2.9, immediately leads to the following lemma:

Lemma 10.1.10. *Let $\bar{x}(t)$ be a deterministic function of time, and assume that (10.27) holds. Then*

$$V_{\text{swaption}}(0) \approx A(0) [(S(0) - c) \Phi(d) + \sqrt{v} \varphi(d)],$$

where

$$d = \frac{S(0) - c}{\sqrt{v}}, \quad v = \int_0^{T_0} q(t, \bar{x}(t))^2 \sigma_r(t)^2 dt. \quad (10.28)$$

It remains to choose $\bar{x}(t)$. An easy choice is to set $\bar{x}(t) = 0$, which will yield good precision if $\sigma_r(t)$ is not too high. What also works reasonably well is to simply evaluate $q(t, \bar{x})$ at the forward discount bond curve, i.e. replace $P(t, T_i, \bar{x})$ with $P(0, T_i)/P(0, t)$ in (10.26). More accurate choices for $\bar{x}(t)$, as well as refinements to the approximation (10.27), are developed in Sections 13.1.4 and 13.1.5.

10.1.4 Swaption Calibration

In a typical application of the model, the European option pricing formulas from Section 10.1.3 are used to calibrate the model, i.e. to find the volatility curve $\sigma_r(t)$ so as to match the market prices of one or more calibration targets, most often European swaptions.

Let us assume that we are given a collection of $N - 1$ swaptions defined on a maturity grid $0 = T_0 < T_1 < \dots < T_N$ such that the i -th swaption expires at times T_i , $i = 1, \dots, N - 1$. Such a collection⁶ is often called a *swaption strip*. A common choice of swaption strip (used, for instance, for Bermudan swaptions) would have the underlying swaps for all swaptions

⁶Note that it is common to set $T_0 = 0$ when defining swaption strips, a convention that slightly clashes with the notation used above when deriving swaption formulas (where the swaption maturity $T_0 > 0$). We shall later, in Chapter 14, develop more formal notation for indexation, but for now trust the reader's ability to adapt generic swaption formulas to the swaption strip convention.

mature on the same date T_N . If this is the case, the strip is called the *coterminal swaption strip*.

With the mean reversion $\kappa(t)$ fixed, we can make the important observation that the value of the swaption expiring at time T_i depends on the volatility curve $\sigma_r(s)$ for $s \in [0, T_i]$ *only*. This can be seen most clearly from the formula for v in Lemma 10.1.10, but is also evident from the pricing formula (10.24) and the fact that the discount bond reconstitution formula (10.18) for $P(t, T)$ does not depend on $\sigma_r(s)$ for $s \geq t$.

The special structure of volatility dependence allows us to perform calibration for one swaption at a time, replacing a potentially multi-dimensional optimization problem with a series of one-dimensional root searches. Assume that $\sigma_r(t)$ is piecewise flat on the maturity grid, with σ_i denoting the flat value on $[T_i, T_{i+1}]$. A possible algorithm based on the formula (10.24) would then work as follows.

1. Assume $\sigma_0, \dots, \sigma_{i-1}$ have been found.
2. Set the value σ_i such that the model price of the $(i+1)$ -th swaption, i.e. a swaption that expires at T_{i+1} , is equal its market price, by numerically inverting (10.24) for σ_i , while $\sigma_0, \dots, \sigma_{i-1}$ are kept constant.
3. Repeat Step 2 for $i = 0, \dots, N-2$.

At first glance, it may appear that the pricing formula from Lemma 10.1.10 will give rise to a linear system on $\sigma_0^2, \dots, \sigma_{N-2}^2$, allowing us to execute Step 2 above by simple matrix inversion. The reality, however, is slightly more complex as the weight functions $q(\cdot, \cdot)$ also depend on the volatility curve $\sigma_r(\cdot)$ through $P(t, T)$'s dependence on $y(t)$ in (10.18). Nevertheless, even with the proper update of $y(t)$ through (10.17) in each step, the inversion in Step 2 above is simple fare for any one-dimensional root solver. Further details can be gleaned from Section 13.1.7 that discusses volatility calibration for the closely related quasi-Gaussian models.

We should note that the volatility calibration scheme above is not guaranteed to always work: a condition sometimes called a “volatility squeeze” may cause the inversion in Step 2 to fail if the market value of the T_{i+1} -expiry swaption is significantly below that of the swaption expiring at T_i . In practice, market data is rarely extreme enough for this to happen, and sometimes the problem can be cured by increasing the mean reversion speed $\kappa(t)$. Some care must be exercised here, though, as the usage of unrealistically high mean reversions will impact the inter-temporal correlations of the model (see Chapter 13), which may lead to unrealistic prices for exotics options, as discussed in Chapter 18.

10.1.5 Finite Difference Methods

We round off our discussion of GSR models with some brief comments on numerical implementation. We start with finite difference methods here, and

turn to Monte Carlo applications in Section 10.1.6. Our discussion of both techniques is rather brief; for further analysis and alternatives we simply refer to Chapters 2 and 3.

10.1.5.1 PDE and Spatial Boundary Conditions

Our treatment of finite difference methods for the GSR model — and for short rate models in general, see Section 11.3.1 — essentially involves little outside of straightforward applications of schemes from Chapter 2. Still, let us start by noting that the algorithms we describe here nevertheless deviate quite markedly from the somewhat old-fashioned (and often suboptimal) tree-based schemes that abound in the short rate literature, even in recent work.

Consider a claim V with the terminal payout $V(T)$ that depends on the discount curve at time T . As the discount curve at time T can be reconstituted solely from knowledge of $x(T)$, we write $V(T) = V(T, x(T))$. By standard results (see Section 1.8), we write $V(t) = V(t, x(t))$, where $V(t, x)$ satisfies the PDE

$$\frac{\partial V}{\partial t} + (y(t) - \varkappa(t)x) \frac{\partial V}{\partial x} + \frac{1}{2}\sigma_r(t)^2 \frac{\partial^2 V}{\partial x^2} = (x + f(0, t))V, \quad (10.29)$$

subject to a known terminal (payout) condition for $V(T, x)$. This PDE can be solved numerically using finite difference methods, e.g. the Crank-Nicolson method in Section 2.2.

In setting up the finite difference scheme, we require knowledge of spatial boundary conditions in the x -domain. In the absence of contractually agreed-upon boundary conditions (as would be the case for e.g. barrier options) one possibility is to set

$$\left. \frac{\partial^2 V}{\partial x^2} \right|_{x=x_{\min}} = \left. \frac{\partial^2 V}{\partial x^2} \right|_{x=x_{\max}} = 0, \quad (10.30)$$

as recommended in Section 2.2.2, where x_{\max} and x_{\min} are the grid boundaries. The boundaries are typically determined by probabilistic means, e.g.

$$x_{\max} = E(x(T)) + \alpha \sqrt{\text{Var}(x(T))}, \quad x_{\min} = E(x(T)) - \alpha \sqrt{\text{Var}(x(T))}, \quad (10.31)$$

for some confidence multiplier α . The moments required in this computation can be found from equations (10.12)–(10.13); see also (10.40)–(10.41).

While workable in practice, the specification (10.30) is not particularly accurate for many actual payout types. As a consequence, one often finds that α needs to be set quite large⁷ (e.g. at values of 5–6, or larger) to

⁷An alternative approach that is advocated by some is to set the mean reversion \varkappa to zero when determining x_{\max} and x_{\min} , in which case α can be reduced.

prevent mis-specification errors at the boundary from affecting the solution at $(t, x) = (0, 0)$. This, in turn, implies that significant computational effort is spent in areas of the x -domain that are probabilistically insignificant. One way to improve on this situation is to rely on the PDE itself to generate boundary conditions, as described earlier in Section 2.2.2 (see also Section 9.4.4). We present the details of this idea in the next section.

10.1.5.2 Determining Spatial Boundary Conditions from PDE

We assume that the PDE (10.29) has been discretized on a spatial grid $\{x_j\}_{j=0}^{m+1}$, so that $V_j(t) = V(t, x_j)$, etc. Let us focus on establishing the boundary condition at $x_0 = x_{\min}$, say. Using a θ -method discretization scheme, as in Section 2.2, with an *upward* discretization of the x -derivatives we get, at some time step $[t, t + \delta]$,

$$\frac{V_0(t + \delta) - V_0(t)}{\delta} + \theta \mu(t, x_0) \frac{V_1(t) - V_0(t)}{x_1 - x_0} \quad (10.32)$$

$$\begin{aligned} & + (1 - \theta) \mu(t + \delta, x_0) \frac{V_1(t + \delta) - V_0(t + \delta)}{x_1 - x_0} \\ & + \frac{\theta}{2} \sigma_r(t)^2 \left\{ \frac{V_2(t) - V_1(t)}{x_2 - x_1} - \frac{V_1(t) - V_0(t)}{x_1 - x_0} \right\} \frac{1}{\frac{1}{2}(x_2 - x_0)} \\ & + \frac{1 - \theta}{2} \sigma_r(t + \delta)^2 \end{aligned} \quad (10.33)$$

$$\times \left\{ \frac{V_2(t + \delta) - V_1(t + \delta)}{x_2 - x_1} - \frac{V_1(t + \delta) - V_0(t + \delta)}{x_1 - x_0} \right\} \frac{1}{\frac{1}{2}(x_2 - x_0)}$$

$$= \theta (x_0 + f(0, t)) V_0(t) + (1 - \theta) (x_0 + f(0, t + \delta)) V_0(t + \delta), \quad (10.34)$$

where $\mu(t, x) \triangleq y(t) - \varkappa(t)x$. This equation can be rearranged to write $V_0(t)$ as

$$V_0(t) = k_1(t)V_1(t) + k_2(t)V_2(t) + g_0(t + \delta), \quad (10.35)$$

where $k_1(t)$ and $k_2(t)$ are easily computed functions of the process parameters, and where $g_0(t + \delta)$ is a function of $V_0(t + \delta)$, $V_1(t + \delta)$, and $V_2(t + \delta)$. We leave it to the reader to write out k_1 , k_2 , and g in detail. Applying similar principles, we get

$$V_{m+1}(t) = k_{m-1}(t)V_{m-1}(t) + k_m(t)V_m(t) + g_{m+1}(t + \delta). \quad (10.36)$$

Comparing (10.35)–(10.36) with the equations (2.12)–(2.13), we see that the boundary conditions (10.35)–(10.36) can be incorporated into our usual tri-diagonal roll-back scheme by simply interpreting $f(t, x_0) = g_0(t + \delta)$ and $\bar{f}(t, x_{m+1}) = g_{m+1}(t + \delta)$ in the scheme of Section 2.2. As we are rolling back in time (from $t + \delta$ to t) when using the finite difference equations, both $g_0(t + \delta)$ and $g_{m+1}(t + \delta)$ are known at time t , so this interpretation involves no difficulties.

10.1.5.3 Upwinding

For the PDE (10.29), notice that the condition (2.34) states that convection domination can cause spurious oscillations to creep into the finite difference scheme unless

$$|y(t) - \kappa(t)x| \Delta_x \leq \sigma_r(t)^2, \quad (10.37)$$

for all x spanned by the finite difference grid. Since $\sigma_r(t)^2$ is typically a small number (around 0.001), it is not uncommon for this inequality to be violated at the edges of the finite difference grid (i.e. in the neighborhoods around x_0 and x_{m+1}) where the mean reversion pushes or pulls strongly at $x(t)$. To avoid numerical difficulties with the finite difference scheme, it is therefore recommended to apply the upwinding scheme in Section 2.6.1. While in principle this may reduce the spatial convergence order of the scheme, in practice the effect of upwinding on convergence is often minimal provided that the finite difference grid is dimensioned in such a way that (10.37) is only violated in a fairly small portion of the grid.

10.1.6 Monte Carlo Simulation

10.1.6.1 Exact Discretization

Consider the problem of pricing a derivative security that pays an amount $V(T)$ at time T , where $V(T)$ may be a function of the entire path of the discount curve over time interval $[0, T]$. Working in the risk-neutral measure, we are thus interested in computing

$$\begin{aligned} V(0) &= \mathbb{E} \left(V(T) e^{-\int_0^T r(u) du} \right) \\ &= P(0, T) \mathbb{E} \left(V(T; \{x(t) : 0 \leq t \leq T\}) e^{-\int_0^T x(u) du} \right), \end{aligned} \quad (10.38)$$

where the second equality shifts variables to $x(t) = r(t) - f(0, t)$ and emphasizes the dependence of $V(T)$ on the entire path of $x(t)$. Recall from the discussion in connection with Proposition 10.1.7 that there are distinct advantages to working with the variable $x(t) = r(t) - f(0, t)$ rather than $r(t)$. In the GSR model, the dynamics for $x(t)$ are given by (10.16), i.e.

$$dx(t) = (y(t) - \kappa(t)x(t)) dt + \sigma_r(t) dW(t), \quad y(t) = \int_0^t e^{-2 \int_u^t \kappa(u) du} \sigma_r(u)^2 du.$$

For the purpose of Monte Carlo pricing of (10.38), we discretize the time-interval into a schedule $t_0 < t_1 < \dots < t_N$, with $t_0 = 0$ and $t_N = T$. The exact choice of the schedule depends on the particulars of the payout $V(T)$; if, say, $V(T)$ only depends on the yield curve at time T , it suffices to set $N = 1$. Now, we can solve the Gaussian SDE for $x(t)$ (see Section 1.6) to write

$$\begin{aligned} x(t_{i+1}) &= e^{-\int_{t_i}^{t_{i+1}} \kappa(u) du} x(t_i) + \int_{t_i}^{t_{i+1}} e^{-\int_s^{t_{i+1}} \kappa(u) du} y(s) ds \\ &\quad + \int_{t_i}^{t_{i+1}} e^{-\int_s^{t_{i+1}} \kappa(u) du} \sigma_r(s) dW(s), \end{aligned} \quad (10.39)$$

which we recognize as being a Gaussian random variable with moments

$$\begin{aligned} \mathbb{E}(x(t_{i+1})|x(t_i)) &= e^{-\int_{t_i}^{t_{i+1}} \kappa(u) du} x(t_i) + \int_{t_i}^{t_{i+1}} e^{-\int_s^{t_{i+1}} \kappa(u) du} y(s) ds, \\ \text{Var}(x(t_{i+1})|x(t_i)) &= \int_{t_i}^{t_{i+1}} \left(e^{-\int_s^{t_{i+1}} \kappa(u) du} \sigma_r(s) \right)^2 ds. \end{aligned} \quad (10.40)$$

Advancement of $x(t)$ on the schedule can thus be done in bias-free manner, by writing

$$x(t_{i+1}) = \mathbb{E}(x(t_{i+1})|x(t_i)) + \sqrt{\text{Var}(x(t_{i+1})|x(t_i))} Z_i, \quad i = 0, \dots, N-1,$$

where Z_0, \dots, Z_{N-1} is a sequence of independent standard Gaussian random variables.

For every date on the simulated path $x(t_0), \dots, x(t_N)$, we can use the reconstitution formula in Proposition 10.1.7 to reconstitute the entire discount curve, in turn allowing us to compute $V(T)$ on the path. To evaluate (10.38), it remains to simulate the quantity

$$I(T) = - \int_0^T x(u) du$$

on the path. Given $x(t_0), \dots, x(t_N)$, an obvious choice would be to compute $I(T)$ by quadrature (e.g. trapezoidal integration, or similar). As this inevitably introduces a discretization bias (see Andersen and Boyle [2000] for more analysis), it is preferable to use the following result.

Lemma 10.1.11. *Let $G(t, T)$ be as in Proposition 10.1.7. Given $I(t_i)$ and $x(t_i)$, $I(t_{i+1})$ is Gaussian with moments*

$$\begin{aligned} \mathbb{E}(I(t_{i+1})|I(t_i), x(t_i)) &= I(t_i) - x(t_i)G(t_i, t_{i+1}) - \int_{t_i}^{t_{i+1}} \int_{t_i}^u e^{-\int_s^u \kappa(v) dv} y(s) ds du, \end{aligned} \quad (10.42)$$

and

$$\begin{aligned} \text{Var}(I(t_{i+1})|I(t_i), x(t_i)) &= 2 \int_{t_i}^{t_{i+1}} \int_{t_i}^u e^{-\int_s^u \kappa(v) dv} y(s) ds du - y(t_i)G(t_i, t_{i+1})^2. \end{aligned} \quad (10.43)$$

Also, we have

$$\begin{aligned}\text{Cov}(x(t_{i+1}), I(t_{i+1})|I(t_i), x(t_i)) \\ = \int_{t_i}^{t_{i+1}} \int_{t_i}^u \sigma_r(s)^2 e^{-\int_s^u \kappa(v)dv} e^{-\int_s^{t_{i+1}} \kappa(v)dv} ds du.\end{aligned}$$

Proof. Straightforward but tedious calculations for Gaussian random variables. \square

Over the time step $[t_i, t_{i+1}]$ we advance $I(t)$ according to the formula

$$I(t_{i+1}) = \mathbb{E}(I(t_{i+1})|I(t_i), x(t_i)) + \sqrt{\text{Var}(I(t_{i+1})|I(t_i), x(t_i))} \tilde{Z}_i,$$

where $\tilde{Z}_0, \dots, \tilde{Z}_{N-1}$ is a sequence of independent standard Gaussian random variables, and where the required moments of $I(t_{i+1})$ can be found in Lemma 10.1.11. To honor the covariance between the $x(t)$ and $I(t)$ processes, we require that the variables Z_i and \tilde{Z}_i are correlated:

$$\text{Corr}(Z_i, \tilde{Z}_i) = \frac{\text{Cov}(I(t_{i+1}), x(t_{i+1})|I(t_i), x(t_i))}{\sqrt{\text{Var}(I(t_{i+1})|I(t_i), x(t_i))} \sqrt{\text{Var}(x(t_{i+1})|I(t_i), x(t_i))}}.$$

As explained in Chapter 3, correlated Gaussian samples can be generated from uncorrelated samples through the Cholesky decomposition.

The scheme outlined above allows us to simulate bias-free paths of the variables $x(t)$ and $I(t)$, which in turn allows us to compute independent, unbiased samples of $V(T)e^{I(T)}$. Monte Carlo estimation of the expectation for $V(0)$ can then be performed in standard Monte Carlo fashion, by forming sample averages. The discretization scheme involves several time-integrals over dates in the observation schedule, many of them nested; it goes without saying that these integrals should be pre-computed before actual path simulations commence.

10.1.6.2 Approximate Discretization

For a quick-and-dirty implementation of the Gaussian model, we may elect to skip the algorithm in the previous section and instead apply one of the approximate discretization schemes in Section 3.2. As a starting point, we have the vector SDE

$$d \begin{pmatrix} x(t) \\ I(t) \end{pmatrix} = \begin{pmatrix} y(t) - \kappa(t)x(t) \\ -x(t) \end{pmatrix} dt + \begin{pmatrix} \sigma_r(t) \\ 0 \end{pmatrix} dW(t).$$

A plain-vanilla Euler scheme from Section 3.2.3 would write

$$\begin{pmatrix} \hat{x}(t_{i+1}) \\ \hat{I}(t_{i+1}) \end{pmatrix} = \begin{pmatrix} \hat{x}(t_i) \\ \hat{I}(t_i) \end{pmatrix} + \begin{pmatrix} y(t_i) - \kappa(t_i)\hat{x}(t_i) \\ -\hat{x}(t_i) \end{pmatrix} \Delta_i + \begin{pmatrix} \sigma_r(t_i) \\ 0 \end{pmatrix} Z_i \sqrt{\Delta_i},$$

where $\Delta_i = t_{i+1} - t_i$ and Z_i is a standard Gaussian random variable. Unless κ is small, this scheme *cannot* be recommended due to the stability issues discussed in Section 3.2.3. As explained in Section 3.2.3.1, it is preferable to incorporate the fact that

$$\begin{aligned}\mathbb{E}(x(t_{i+1})|x(t_i)) &= e^{-\int_{t_i}^{t_{i+1}} \kappa(u)du} x(t_i) + \int_{t_i}^{t_{i+1}} e^{-\int_s^{t_{i+1}} \kappa(u)du} y(s) ds \\ &\approx e^{-\kappa(t_i)(t_{i+1}-t_i)} x(t_i) + \frac{1 - e^{\kappa(t_i)(t_{i+1}-t_i)}}{\kappa(t_i)} y(t_i).\end{aligned}$$

That is, we write

$$\begin{pmatrix} \widehat{x}(t_{i+1}) \\ \widehat{I}(t_{i+1}) \end{pmatrix} = \begin{pmatrix} e^{-\kappa(t_i)\Delta_i} x(t_i) + \frac{1 - e^{\kappa(t_i)\Delta_i}}{\kappa(t_i)} y(t_i) \\ \widehat{I}(t_i) - \widehat{x}(t_i)\Delta_i \end{pmatrix} + \begin{pmatrix} \sigma_r(t_i) \\ 0 \end{pmatrix} Z_i \sqrt{\Delta_i}.$$

This scheme has first-order (weak) convergence. Higher-order schemes can be found in Section 3.2, but are essentially obsolete here: if truly low bias is critically important, we should use the unbiased scheme from Section 10.1.6.1.

10.1.6.3 Using other Measures for Simulation

The need to simulate $I(t)$ can be avoided entirely by a suitable change of the probability measure. Switching to the terminal measure Q^T (see Section 4.2.4), we rewrite (10.38) as

$$V(0) = P(0, T) \mathbb{E}^T(V(T)),$$

and observe that we now need to simulate $x(t)$ only in order to calculate the payoff. The dynamics of $x(t)$ under the terminal measure Q^T follow from (4.34),

$$\begin{aligned}dx(t) &= (y(t) - \kappa(t)x(t)) dt + \sigma_r(t) (dW^T(t) - \sigma_P(t, T) dt) \\ &= (y(t) - \sigma_r(t)^2 G(t, T) - \kappa(t)x(t)) dt + \sigma_r(t) dW^T(t),\end{aligned}$$

with $W^T(t)$ being a Q^T -Brownian motion. The dynamics remain Gaussian and Markov under the terminal measure, and hence $x(t)$ can be simulated bias-free on the time grid $\{t_i\}_{i=1}^N$.

An alternative to the terminal measure that is “closer” in some ways to the risk-neutral measure, yet still allows one to avoid the simulation of $I(t)$, is the spot measure from Section 4.2.3. We recall that this measure is associated with the discretely compounded money market account $B(t)$,

$$B(t) = P(t, t_{i+1}) \prod_{n=0}^i \frac{1}{P(t_n, t_{n+1})}, \quad t \in (t_i, t_{i+1}].$$

Under the spot measure Q^B , we obtain

$$V(0) = \mathbb{E}^B \left(\prod_{n=0}^{N-1} P(t_n, t_{n+1}, x(t_n)) V(T) \right),$$

where we explicitly indicated the dependence of the discount bond on the state process $x(t)$. Notice that the random variable under the expectation operator is a function of $x(t)$ on the grid $\{t_i\}_{i=1}^N$ only. Moreover, over the interval $[t_n, t_{n+1}]$, the measure Q^B coincides with the T_{n+1} -forward measure, which gives us the dynamics of $x(t)$,

$$\begin{aligned} dx(t) &= (y(t) - \sigma_r(t)^2 G(t, t_{n+1}) - \kappa(t)x(t)) dt \\ &\quad + \sigma_r(t) dW^B(t), \quad t \in (t_n, t_{n+1}], \end{aligned}$$

with $W^B(t)$ a Q^B -Brownian motion. Again, we can generate a sample of $x(t_{n+1})$ from $x(t_n)$ in a bias-free manner. We refer the reader to Chapter 14 for more details on numeraire simulation strategies.

10.2 The Affine One-Factor Model

Earlier (in Section 10.1.1.3) we identified the non-zero probability of negative interest rates as one of the drawbacks of the one-factor GSR model. Another problem is the lack of interest rate dependence in the GSR short rate volatility, leaving the user with no means of controlling the volatility skew implied by the model. While there are different ways of addressing these issues, one type of model that can, in part at least, address both of these shortcomings of the GSR model is the *affine short rate model*. This model — or, rather, model *class* — constitutes a significant extension of the GSR model (which in fact is a member of the affine class), yet retains a high degree of analytical tractability. Originally introduced by Duffie and Kan [1996], the affine class of short rate models enjoys high popularity among practitioners and academics alike, particularly for econometric work. The affine models are also quite useful for derivatives pricing, although ultimately the constraints one need to impose on diffusion dynamics can be too strong for some applications.

10.2.1 Basic Definitions

10.2.1.1 SDE

Consider a time-homogeneous one-factor short rate process of the form

$$dr(t) = \kappa(\vartheta - r(t)) dt + \sigma v(r(t)) dW(t), \quad (10.44)$$

where $W(t)$ is a Brownian motion in the risk-neutral measure, $\kappa > 0$, $\sigma > 0$, and ϑ are constants, and $v(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is a deterministic function of the level of the short rate. We notice that the drift of (10.44) is affine, i.e. linear, in $r(t)$. If the square of the diffusion term in (10.44) is also affine, we say that (10.44) is a time-homogeneous *affine one-factor short rate model*. Evidently, the function $v(r)$ is thus limited to the form

$$v(r) = \sqrt{\alpha + \beta r}, \quad (10.45)$$

for constants α and β . We notice that the special case of $\beta = 0$ produces the GSR model of Section 10.1.2.2, whereas the case $\alpha = 0$ produces a square-root type model similar to those encountered, for stochastic volatility modeling, in Chapter 8. The case $\alpha = 0$, $\beta = 1$ was first studied by Cox et al. [1985] and is known as the *Cox-Ingersol-Ross* (CIR) model.

10.2.1.2 Regularity Issues

Not all combinations of parameters in (10.44) and (10.45) produce a well-defined SDE. If $\beta = 0$ for all t , we must require that $\alpha \geq 0$ for all t to ensure that $v(r(t))$ is defined. If $\beta \neq 0$, for the square root in (10.45) to exist, we must ensure that the drift term in (10.44) has the same sign as β whenever $\alpha + \beta r(t) = 0$. That is,

$$\kappa\beta(\vartheta + \alpha/\beta) \geq 0, \quad \beta \neq 0, \quad (10.46)$$

for all $t \geq 0$. Notice that if we wish for the volatility term in (10.45) to be strictly positive ($\alpha + \beta r(t) > 0$), we need to replace this condition with the stronger Feller condition (recall Proposition 8.3.1)

$$\kappa\beta(\vartheta + \alpha/\beta) \geq \frac{1}{2}\beta^2\sigma^2.$$

For the CIR model the requirement that $r(t)$ stays strictly positive can be seen to translate into the classical condition $2\kappa\vartheta \geq \sigma^2$.

For the purposes of modeling interest rates, it is most reasonable to assume that $\kappa > 0$ to ensure that rates are mean-reverting rather than mean-fleeing, and that $\beta \geq 0$. In this case, the domain of the short rate becomes

$$r(t) \in [-\alpha/\beta, \infty), \quad \beta > 0, \quad (10.47)$$

and $r(t) \in (-\infty, \infty)$ for the case $\beta = 0$, $\alpha > 0$ (Gaussian model). Evidently, to keep $r(t)$ non-negative for all t , we need to set $\alpha \leq 0$, subject to the restriction that $-\alpha/\beta \leq r(0)$.

10.2.1.3 Volatility Skew

The parameters α and β in (10.44) effectively determine the volatility skew behavior of the affine model. If both parameters are non-negative, the affine

model can generate skews ranging from a Gaussian process ($\alpha \gg \beta$) to a square-root process ($\alpha \ll \beta$). In the usual language, for non-negative α, β , the skew “power” of the affine model thus ranges from 0 to 0.5. By allowing α to be negative, effective skew powers above 0.5 are possible, although the allowed range of the underlying process $r(t)$ will then be floored at some positive level, which may have undesirable side effects if α/β is not close to zero.

10.2.1.4 Time-Dependent Parameters

The SDE (10.44) does not depend on time and hence will generally not match the initial yield curve at time 0. As we did for the Gaussian model, we may extend the SDE to have time-dependent parameters, e.g.

$$dr(t) = \kappa(t)(\vartheta(t) - r(t))dt + \sigma(t)\sqrt{\alpha + \beta r(t)}dW(t). \quad (10.48)$$

Notice that we have not introduced time-dependence in α and β , leaving the domain (10.47) unchanged⁸. Not all functions $\kappa, \vartheta, \sigma$ produce a well-defined SDE; for instance, if $\kappa(t)$ is positive (which is always the case in practice) and $\beta > 0$, then (10.46) shows that we need

$$\kappa(t)\beta\vartheta(t) \geq -\alpha \quad (10.49)$$

in order for (10.48) to be well-defined.

Remark 10.2.1. For generality we allow $\kappa(t)$ to be a function of time throughout. As argued in Section 10.1.2.3, however, it is often most reasonable to let $\kappa(t)$ be a constant.

10.2.2 Discount Bond Pricing and Extended Transform

Starting from the time-dependent SDE (10.48), we now turn to the search for a discount bond reconstitution formula, i.e. a formula that allows us to compute the risk-neutral expectation

$$P(t, T) = E_t \left(e^{-\int_t^T r(u)du} \right), \quad (10.50)$$

as an explicit function of $r(t)$. Rather than directly attacking (10.50), we turn to the more general problem of establishing the so-called *extended transform* $g(t, T; c_1, c_2)$ defined by

$$g(t, T; c_1, c_2) = E_t \left(e^{-c_1 r(T) - c_2 \int_t^T r(u)du} \right), \quad c_1, c_2 \in \mathbb{C}. \quad (10.51)$$

⁸The results of the next sections do, however, often generalize to time-dependence in α and β . See Remark 10.2.3, for example.

Notice how this generalizes the idea of the moment-generating function from Chapters 8 and 9. Also note that the knowledge of g allows us to find discount bond prices as a special case,

$$P(t, T) = g(t, T; 0, 1).$$

For the values of c_1 and c_2 for which g exists, we can use the following result, which is an extension of Proposition 9.1.2.

Proposition 10.2.2. *For the model (10.48), whenever the extended transform g in (10.51) is defined, it is given by*

$$g(t, T; c_1, c_2) = \exp(A(t, T; c_1, c_2) - B(t, T; c_1, c_2)r(t)), \quad (10.52)$$

where A and B satisfy a system of Riccati ODEs

$$\frac{dA}{dt} - \varkappa(t)\vartheta(t)B + \frac{1}{2}\sigma(t)^2\alpha B^2 = 0, \quad (10.53)$$

$$-\frac{dB}{dt} + \varkappa(t)B + \frac{1}{2}\sigma(t)^2\beta B^2 = c_2, \quad (10.54)$$

subject to the terminal conditions $A(T; T, c_1, c_2) = 0$, $B(T; T, c_1, c_2) = c_1$.

Proof. Follows that of Proposition 9.1.2 closely. \square

Remark 10.2.3. If the parameters α and β are functions of time, Proposition 10.2.2 continues to hold if we simply replace α , β with $\alpha(t)$, $\beta(t)$ in the Riccati equations for A and B .

Proposition 10.2.2 establishes that the joint characteristic function of $r(T)$ and $\int_0^T r(u) du$ is known analytically for the affine model, a result that accounts for much of its popularity in the financial literature. Solution of the Riccati equations (10.53)–(10.54) can be done quickly and robustly by any number of standard ODE schemes, such as the Runge-Kutta method (see Press et al. [1992]). For the case where parameters are piecewise constant in time, establishing A and B in Proposition 10.2.2 can also be done analytically; see Section 10.2.2.1 below.

10.2.2.1 Constant Parameters

We now turn to establishing the extended transform g for the special case where all parameters in (10.48) are constants. As a warm-up case, we first list a result for the CIR case.

Proposition 10.2.4. *Consider the CIR model*

$$dr(t) = \varkappa(\vartheta - r(t)) dt + \sigma\sqrt{r(t)} dW(t), \quad (10.55)$$

and let $g(t, T; c_1, c_2)$ be defined as in (10.51). Set

$$\gamma = \gamma(c_2, \sigma) = \sqrt{\kappa^2 + 2\sigma^2 c_2}.$$

Then

$$g(t, T; c_1, c_2) = \exp(A_{\text{CIR}}(t, T; \vartheta, \sigma, c_1, c_2) - B_{\text{CIR}}(t, T; \vartheta, \sigma, c_1, c_2)r(t)),$$

where

$$\begin{aligned} A_{\text{CIR}}(t, T; \vartheta, \sigma, c_1, c_2) &= \kappa \vartheta \sigma^{-2} (\kappa + \gamma(c_2, \sigma)) (T - t) \\ &- 2\kappa \vartheta \sigma^{-2} \ln \left(1 + \frac{(\kappa + \gamma(c_2, \sigma) - c_1 \sigma^2) (e^{\gamma(c_2, \sigma)(T-t)} - 1)}{2\gamma(c_2, \sigma)} \right), \end{aligned}$$

and

$$\begin{aligned} B_{\text{CIR}}(t, T; \vartheta, \sigma, c_1, c_2) &= \\ &\frac{(2c_2 - \kappa c_1) (1 - e^{-\gamma(c_2, \sigma)(T-t)}) + \gamma(c_2, \sigma) c_1 (1 + e^{-\gamma(c_2, \sigma)(T-t)})}{(\kappa + \gamma(c_2, \sigma) + c_1 \sigma^2) (1 - e^{-\gamma(c_2, \sigma)(T-t)}) + 2\gamma(c_2, \sigma) e^{-\gamma(c_2, \sigma)(T-t)}}. \end{aligned}$$

Proof. The result is a small extension of Proposition 8.3.8, and follows by direct solution of the ODEs (10.53)–(10.54). \square

Armed with this result, it is straightforward to extend it to the general constant-parameter case

$$dr(t) = \kappa(\vartheta - r(t)) dt + \sigma \sqrt{\alpha + \beta r(t)} dW(t), \quad \beta > 0. \quad (10.56)$$

In particular, we notice that if $y(t) = \alpha + \beta r(t)$, then $y(t)$ follows the SDE

$$dy(t) = \beta dr(t) = \kappa(\beta \vartheta + \alpha - y(t)) dt + \beta \sigma \sqrt{y(t)} dW(t),$$

which is of the form (10.55). We also have

$$\begin{aligned} g(t, T; c_1, c_2) &= E_t \left(\exp \left(-c_1 r(T) - c_2 \int_t^T r(u) du \right) \right) \\ &= E_t \left(\exp \left(-c_1 \left(\frac{y(T) - \alpha}{\beta} \right) - c_2 \int_t^T \left(\frac{y(u) - \alpha}{\beta} \right) du \right) \right) \\ &= e^{c_1 \alpha / \beta} e^{c_2 \alpha (T-t) / \beta} E_t \left(\exp \left(-\frac{c_1}{\beta} y(T) - \frac{c_2}{\beta} \int_t^T y(u) du \right) \right). \end{aligned}$$

The expectation involved in the last equality can here be evaluated directly from Proposition 10.2.4, leading to the following lemma.

Lemma 10.2.5. *The extended transform for the constant parameter affine model (10.56) is*

$$g(t, T; c_1, c_2) = \exp(A(t, T; c_1, c_2) - B(t, T; c_1, c_2)r(t)),$$

where

$$\begin{aligned} A(t, T; c_1, c_2) &= c_1\alpha/\beta + c_2\alpha(T-t)/\beta + A_{\text{CIR}}\left(t, T; \beta\vartheta + \alpha, \beta\sigma, \frac{c_1}{\beta}, \frac{c_2}{\beta}\right) \\ &\quad - \alpha B_{\text{CIR}}\left(t, T; \beta\vartheta + \alpha, \beta\sigma, \frac{c_1}{\beta}, \frac{c_2}{\beta}\right), \\ B(t, T; c_1, c_2) &= \beta B_{\text{CIR}}\left(t, T; \beta\vartheta + \alpha, \beta\sigma, \frac{c_1}{\beta}, \frac{c_2}{\beta}\right), \end{aligned}$$

and the functions A_{CIR} and B_{CIR} are given in Proposition 10.2.4.

10.2.2.2 Piecewise Constant Parameters

We can use the results established in Section 10.2.2.1 to compute extended transforms for the case where we are given a time grid $0 = t_0 < t_1 < t_2 < \dots$, on which all model parameters \varkappa and σ can be assumed piecewise constant. The resulting recursive routine is a robust and efficient⁹ alternative to Runge-Kutta solvers.

For simplicity of notation, let us define $g(t_i, t_j; c_1, c_2) = g_{i,j}(c_1, c_2)$, $A(t_i, t_j; c_1, c_2) = A_{i,j}(c_1, c_2)$, and so on. Then, from Proposition 10.2.2,

$$g_{i,j}(c_1, c_2) = e^{A_{i,j}(c_1, c_2) - r(t_i)B_{i,j}(c_1, c_2)}, \quad j > i, \quad (10.57)$$

and, using the law of iterated conditional expectations,

$$\begin{aligned} g_{i-1,j}(c_1, c_2) &= \mathbb{E}_{t_{i-1}}\left(e^{-c_1r(t_i)-c_2\int_{t_{i-1}}^{t_i} r(u)du}\right) \\ &= \mathbb{E}_{t_{i-1}}\left(\mathbb{E}_{t_i}\left(e^{-c_1r(t_i)-c_2\int_{t_{i-1}}^{t_i} r(u)du}\right)\right) \\ &= \mathbb{E}_{t_{i-1}}\left(e^{-c_2\int_{t_{i-1}}^{t_i} r(u)du} \mathbb{E}_{t_i}\left(e^{-c_1r(t_i)-c_2\int_{t_i}^{t_j} r(u)du}\right)\right) \\ &= \mathbb{E}_{t_{i-1}}\left(e^{-c_2\int_{t_{i-1}}^{t_i} r(u)du} g_{i,j}(c_1, c_2)\right). \end{aligned}$$

Inserting (10.57) into the last equation then yields

⁹As pointed out in Section 9.1, depending on the level of accuracy required, the Runge-Kutta numerical solution of the ODEs can sometimes have higher computational efficiency.

$$\begin{aligned} g_{i-1,j}(c_1, c_2) &= \mathbb{E}_{t_{i-1}} \left(e^{-c_2 \int_{t_{i-1}}^{t_i} r(u) du} e^{A_{i,j}(c_1, c_2) - r(t_i) B_{i,j}(c_1, c_2)} \right) \\ &= e^{A_{i,j}(c_1, c_2)} g_{i-1,i}(B_{i,j}(c_1, c_2), c_2). \end{aligned}$$

Applying (10.57) to the right-hand side of this equation leads to

$$\begin{aligned} &e^{A_{i-1,j}(c_1, c_2) - r(t_{i-1}) B_{i-1,j}(c_1, c_2)} \\ &= e^{A_{i,j}(c_1, c_2)} e^{A_{i-1,i}(B_{i,j}(c_1, c_2), c_2) - r(t_{i-1}) B_{i-1,i}(B_{i,j}(c_1, c_2), c_2)}, \end{aligned}$$

or, finally,

$$A_{i-1,j}(c_1, c_2) = A_{i,j}(c_1, c_2) + A_{i-1,i}(B_{i,j}(c_1, c_2), c_2), \quad (10.58)$$

$$B_{i-1,j}(c_1, c_2) = B_{i-1,i}(B_{i,j}(c_1, c_2), c_2). \quad (10.59)$$

As parameters are constant on the time grid, the functions $A_{i-1,i}$ and $B_{i-1,i}$ can be computed in closed form from the results of Lemma 10.2.5. For a fixed j , (10.58)–(10.59) can be used in backward fashion to establish $A_{i,j}$ and $B_{i,j}$ for $i = j-1, j-2, \dots, 0$; the recursion starts with an application of Lemma 10.2.5 to compute $A_{j-1,j}$ and $B_{j-1,j}$.

10.2.3 Discount Bond Calibration

10.2.3.1 Change of Variables

In the affine SDE (10.48), the role of the mean reversion level $\vartheta(t)$ is to calibrate the model to the initial term structure of discount bonds. As we discussed in the context of the GSR model, $\vartheta(t)$ will depend on the derivative $\partial f(0, t)/\partial t$ which may, for many curve construction algorithms, be irregular. For practical applications of affine models, it is therefore strongly recommended to follow the advice of Section 10.1.2.2 and rewrite the model in terms of a variable that measures the difference between $r(t)$ and $f(0, t)$. Let this variable be $x(t)$, defined as

$$x(t) = r(t) - f(0, t).$$

The SDE for $x(t)$ becomes

$$\begin{aligned} dx(t) &= dr(t) - \frac{\partial f(0, t)}{\partial t} dt \\ &= (\omega(t) - \varkappa(t)x(t)) dt + \sigma(t)\sqrt{\xi(t) + \beta x(t)} dW(t), \end{aligned} \quad (10.60)$$

where $x(0) = 0$, $\xi(t) = \alpha + \beta f(0, t)$, and

$$\omega(t) = \varkappa(t)\vartheta(t) - \partial f(0, t)/\partial t - \varkappa(t)f(0, t).$$

The deterministic function $\omega(t)$ is likely to be smooth even if the forward curve is not.

Written in terms of $x(t)$, the extended transform in Proposition 10.2.2 becomes

$$\begin{aligned} g(t, T; c_1, c_2) &= e^{-c_1 f(0, T) - c_2 \int_t^T f(0, u) du} E_t \left(e^{-c_1 x(T) - c_2 \int_t^T x(u) du} \right) \\ &= e^{-c_1 f(0, T)} \frac{P(0, T)^{c_2}}{P(0, t)^{c_2}} \exp(C(t, T; c_1, c_2) - x(t)B(t, T; c_1, c_2)), \end{aligned} \quad (10.61)$$

where B solves (10.54) and C can, after suitable translation of the results in Proposition 10.2.2 to the process (10.60), be written as the solution to the Riccati ODE:

$$\frac{dC}{dt} - \omega(t)B + \frac{1}{2}\sigma(t)^2\xi(t)B^2 = 0. \quad (10.62)$$

10.2.3.2 Algorithm for $\omega(t)$

We now assume (but see Section 10.2.5) that α and β have been fixed, and that $\varkappa(t)$ and $\sigma(t)$ are known for all values of $t \geq 0$. In the SDE (10.60) for $x(t)$, it only remains to establish the function $\omega(t)$, which shall be done to match observed discount bond prices at time 0.

To make matters more concise, let us set $b(t, T) = B(t, T; 0, 1)$ and $c(t, T) = C(t, T; 0, 1)$ such that, from the definition of $C(t, T)$,

$$P(t, T) = g(t, T; 0, 1) = \frac{P(0, T)}{P(0, t)} e^{c(t, T) - x(t)b(t, T)}. \quad (10.63)$$

The functions b and c obviously satisfy a Riccati system,

$$\frac{dc}{dt} - \omega(t)b + \frac{1}{2}\sigma(t)^2\xi(t)b^2 = 0, \quad (10.64)$$

$$-\frac{db}{dt} + \varkappa(t)b + \frac{1}{2}\sigma(t)^2\beta b^2 = 1, \quad (10.65)$$

where $c(T, T) = b(T, T) = 0$.

Setting $t = 0$ in equation (10.63) establishes the fundamental calibration requirement that $c(0, T) = c(0, T; \omega(\cdot)) = 0$ for all T which, combined with (10.64), defines a so-called *Volterra integral equation* for $\omega(\cdot)$. We can solve it on a time grid $t_0 < t_1 < t_2 < \dots < t_N$ by iterative bootstrapping of the equation $c(0, t_i; \omega(\cdot)) = 0$. Assuming that $\omega(\cdot)$ is piecewise constant at a level ω_i over the time bucket $(t_i, t_{i+1}]$, we can use the following algorithm.

1. As a pre-processing step, find $b(t_i, t_j)$ for all i, j , $j > i$, by solving (10.65). This does not depend on $\omega(\cdot)$.
2. For a given i , assume that ω_j is known for $j < i$.
3. Compute $\Theta(t_i) = \frac{1}{2} \int_0^{t_{i+1}} \sigma(s)^2 \xi(s) b(s, t_{i+1})^2 ds - \int_0^{t_i} \omega(s) b(s, t_{i+1}) ds$.

4. Compute ω_i as the solution to $\Theta(t_i) - \omega_i \int_{t_i}^{t_{i+1}} b(s, t_{i+1}) ds = 0$.
5. Repeat steps 2–4 for all $i = 0, 1, \dots, N - 1$.

Notice that no numerical root search is needed and that the computational complexity of the scheme is $O(N^2)$. By modifying Steps 3 and 4, other interpolation techniques can be supported, although stability issues might come into play. See also Press et al. [1992] for more general schemes to solve Volterra equations.

We should note that there may be cases where the algorithm above will fail, in the sense that the basic regularity condition (10.49) will prevent a valid solution for $\omega(\cdot)$ from existing. This is a fundamental issue with non-Gaussian affine short rate models, but is rarely observed as very strongly downward-sloping yield curves are required to trigger the problem (see the discussion in Hull and White [1994a]).

10.2.4 European Option Pricing

The short rate volatility function $\sigma(t)$ in the affine model (10.60) will normally be determined through calibration against swaptions and caps/floors. For such calibration to be computationally feasible it is, of course, important to establish fast methods for pricing European interest rate options.

For simple options such as caplets or, equivalently, options on zero-coupon bonds, the availability of the moment-generating function for the logarithm of the bond (see Proposition 10.2.2) allows for application of the Fourier methods¹⁰ of Section 8.4. Extensions to swaption pricing through the Jamshidian approach of Section 10.1.3.1 is possible in principle, but the need to perform Fourier integration of a large number of Riccati ODE solutions makes this approach impractical. Several approximation techniques have been proposed in the literature; see, for instance, Collin-Dufresne and Goldstein [2002a] for a survey and details on a method based on Gram-Charlier expansions. Our preferred approach to swaption pricing in the affine model borrows the techniques of Section 10.1.3.2 to work out an approximation for the swap rate martingale dynamics. We shall outline one straightforward and quite accurate approach here; as was the case for the GSR model, we again will stop short of the full-blown projection techniques that will be introduced later in this book for more realistic candidates for actual trading applications.

Let us, as in Section 10.1.3.2, start out by rewriting the swaption payout as

$$V_{\text{swaption}}(T_0) = A(T_0)(S(T_0) - c)^+, \quad (10.66)$$

where

¹⁰For time-homogeneous models, closed-form pricing formulas for options on discount bonds exist for some models, including the CIR model (see Cox et al. [1985]).

$$A(t) = \sum_{i=0}^{N-1} \tau_i P(t, T_{i+1}), \quad S(t) = \frac{P(t, T_0) - P(t, T_N)}{A(t)}.$$

Let Q^A be the measure induced by using $A(t)$ as the numeraire; in this measure $S(t)$ is a martingale. By the reconstitution result (10.63) we have

$$dS(t) = \frac{\partial S(t)}{\partial x} \sigma(t) \sqrt{\xi(t) + \beta x(t)} dW^A(t), \quad (10.67)$$

where $W^A(t)$ is a Q^A -Brownian motion and

$$\begin{aligned} \frac{\partial S(t)}{\partial x} &= -\frac{b(t, T_0)P(t, T_0) - b(t, T_N)P(t, T_N)}{A(t)} \\ &\quad + \frac{S(t)}{A(t)} \sum_{i=0}^{N-1} \tau_i b(t, T_{i+1})P(t, T_{i+1}). \end{aligned}$$

The dynamics (10.67) are generally intractable, but $S(t)$ can — as was the case for the GSR model — be verified to often be well approximated by a linear function of $x(t)$, with slope and intercept being functions of time. Using a Taylor expansion around some point \bar{x} (e.g. $\bar{x} = 0$, but see the discussion in Section 10.1.3.2), we can find $\zeta(t)$, $\chi(t)$ such that

$$S(t) \approx \zeta(t) + \chi(t)x(t),$$

and then (10.67) approximately reduces to an affine SDE for $S(t)$:

$$\begin{aligned} dS(t) &\approx \chi(t)\sigma(t)\sqrt{\xi(t) + \beta x(t)} dW^A(t) \\ &= \chi(t)\sigma(t)\sqrt{\xi(t) + \beta \left(\frac{S(t) - \zeta(t)}{\chi(t)} \right)} dW^A(t) \\ &= \sigma(t)\sqrt{\xi_s(t) + \beta_s(t)S(t)} dW^A(t). \end{aligned} \quad (10.68)$$

While valuation of the payout (10.66) cannot be accomplished in closed form when $S(t)$ follows the time-dependent affine SDE (10.68), we can always rely on transform-based methods. Indeed, it is evident that the characteristic function of $S(T_0)$ can be constructed by applying Proposition 10.2.2 and Remark 10.2.3 to (10.68), whereafter Theorem 8.4.3 gives us a way to calculate the required expected value in

$$V(0) = A(0)E^{Q^A} \left((S(T_0) - c)^+ \right). \quad (10.69)$$

We trust that the reader can see how this would work, so we omit the details. Instead, we proceed to further simplify matters, through time averaging of parameters.

First, we wish to reduce (10.68) to the simplified form

$$dS(t) = \sigma(t) \sqrt{\beta_s(t)} \sqrt{\psi + S(t)} dW^A(t), \quad (10.70)$$

where ψ is some constant. One approach for setting ψ is to simply match quadratic variance of $S(t)$ over $[0, T_0]$, i.e.

$$\int_0^{T_0} \sigma(t)^2 \beta_s(t) \psi dt = \int_0^{T_0} \sigma(t)^2 \beta_s(t) \xi_s(t) dt$$

or

$$\psi = \frac{\int_0^{T_0} \sigma(t)^2 \beta_s(t) \xi_s(t) dt}{\int_0^{T_0} \sigma(t)^2 \beta_s(t) dt}.$$

A more sophisticated alternative would be to rely on a small-noise expansion, as in Chapter 7. In any case, for the SDE (10.70), the expectation in (10.69) can be evaluated in closed form. To see this, simply define $y(t) \equiv \psi + S(t)$ and note that

$$dy(t) = \sigma(t) \sqrt{\beta_s(t)} \sqrt{y(t)} dW^A(t), \quad y(0) = \psi + S(0), \quad (10.71)$$

and

$$V(0) = A(0) E^{Q^A} \left((S(T_0) - c)^+ \right) = A(0) E^{Q^A} \left((y(T_0) - c_y)^+ \right), \quad (10.72)$$

with $c_y \triangleq \psi + c$. Since $y(t)$ in (10.71) is simply a (time-dependent) CEV process with CEV power 1/2, computation of the call option expectation in (10.72) can be carried out by the formulas in Section 7.2. Swaption prices produced this way are, in our experience, accurate and robust, and much more convenient to compute than by competing methods.

10.2.5 Swaption Calibration

As we showed in Section 10.1.4, calibration of the GSR model volatility to swaption prices is a matter of straightforward bootstrapping. Unfortunately, matters are more complicated for general affine models.

10.2.5.1 Basic Problem

To gain insight, let us first consider the simple problem of calibrating the model volatility function $\sigma(t)$ in (10.60) to match the time 0 price of a Δ -tenor zero-coupon bond option maturing at T . Assuming that the initial yield curve is known at time 0, how much volatility information is needed to price this option? The answer to this question depends on the specification of ξ and β .

If $\beta = 0$, we know that the function $b(T, T + \Delta)$ in the bond reconstitution formula (10.63) is independent of σ ; see Proposition 10.1.7 (and adjust notation accordingly). It can also be verified that while $c(T, T + \Delta)$ in

(10.63) depends on the initial discount curve all the way to time $T + \Delta$, it only requires the specification of $\sigma(t)$ to time T . Further, the state of $x(T)$ only depends on $\sigma(t)$, $t < T$. In total, when $\beta = 0$, the discount bond option payout is only affected by $\{\sigma(t)\}_{0 \leq t \leq T}$, irrespective of the magnitude of the bond tenor Δ . This is also obvious from the reconstitution formula (10.18).

If $\beta \neq 0$, however, we see from (10.65) that $b(T, T + \Delta)$ depends¹¹ on the volatility $\{\sigma(t)\}_{0 \leq t \leq T + \Delta}$. This again makes $c(t, T + \Delta)$ depend on volatilities in $[t, T + \Delta]$, requiring the full knowledge of $\{\sigma(t)\}_{0 \leq t \leq T + \Delta}$ to price the option at time 0. This fact has implications for calibration to, say, swaption prices as regular bootstrapping techniques cannot be employed.

10.2.5.2 Calibration Algorithm

Consider now the situation where we wish to calibrate our volatility function $\sigma(t)$ to a swaption strip defined on a maturity grid $0 = T_0 < T_1 < \dots < T_N$. Recall that a swaption strip consists of $N - 1$ swaptions expiring at times T_i , $i = 1, \dots, N - 1$; we here assume that all swaptions are written on swaps that mature at time T_N (coterminal strip). According to the discussion above, pricing any one of these swaptions — even the short-dated ones — in an affine model will require knowledge of $\{\sigma(t)\}_{0 \leq t \leq T_N}$. As it would be too slow to calibrate volatilities by simultaneous, multi-dimensional root search on all levels $\sigma(T_i)$, $i = 0, 1, \dots, N$, we instead notice that while, say, the swaption maturing on date T_i depends on volatilities everywhere on $[0, T_N]$, its dependence on the volatilities in $[0, T_i]$ is much stronger than on the volatilities in the interval $(T_i, T_N]$. Assuming that $\sigma(t)$ is piecewise constant on the maturity grid — with σ_i denoting the flat value on $(T_i, T_{i+1}]$ — we can use this observation to propose the following iterative calibration approach.

1. Start out by setting all σ_i , $i = 0, \dots, N - 1$, equal to a reasonably chosen constant, or equal to values approximated from a calibrated GSR¹² model.
2. Compute $\omega(\cdot)$ to match time 0 prices of the N discount bonds maturing on T_1, T_2, \dots, T_N . One can use the algorithm in Section 10.2.3.2 for this.
3. Set the value σ_0 — but leave all other volatilities σ_i , $i = 1, \dots, N - 1$, unchanged — such that the swaption maturing at time T_1 is priced correctly. We can use the pricing techniques in Section 10.2.4 for this.
4. Repeat Step 3 for $\sigma_1, \sigma_2, \dots, \sigma_{N-2}$, always leaving future (but not past) points on the volatility curve unchanged.
5. Repeat Step 2 and recompute all swaption prices.
6. Repeat Steps 3–5 until all swaptions are priced within given tolerances.

¹¹Recall that we solve $b(t, T + \Delta)$ backward in time from the known boundary condition at $t = T + \Delta$.

¹²For instance, if $\sigma_g(t)$ is the volatility function in the Gaussian model, then we can extract an estimate for $\sigma(t)$ from the relation $\sigma(t)\sqrt{\alpha + f(0, t)\beta} \approx \sigma_g(t)$.

Notice that in Step 4, altering σ_i will slightly distort the prices of swaptions maturing at dates earlier than T_i ; this necessitates the iteration in Step 6.

We (re-)emphasize that the algorithm above, when applied to the Gaussian model, will converge within one iteration in Step 6. Finally, we note that the calibrated model needs to be checked against the regularity conditions discussed in Section 10.2.1.2; if conditions are violated, the problem may potentially be remedied by increasing α .

10.2.6 Quadratic One-Factor Model

In conclusion, let us consider an interesting special case of an affine class. A *quadratic Gaussian one-factor* model is obtained by specifying the short rate to be a quadratic function of a linear Gaussian process,

$$r(t) = \alpha(t) + \beta(t)y(t) + \gamma(t)y(t)^2, \quad (10.73)$$

where

$$dy(t) = -\varkappa(t)y(t)dt + \sigma(t)dW(t), \quad y(0) = 0. \quad (10.74)$$

While this is not immediately obvious, the model is indeed of affine type, albeit in *two* factors. If we denote $u(t) = y(t)^2$, we see that $r(t)$ is a linear function of the state vector $(y(t), u(t))$, which follows the SDEs

$$d \begin{pmatrix} y(t) \\ u(t) \end{pmatrix} = \begin{pmatrix} -\varkappa(t)y(t) \\ \sigma(t)^2 - 2\varkappa(t)u(t) \end{pmatrix} dt + \sigma(t) \begin{pmatrix} 1 \\ 2y(t) \end{pmatrix} dW(t), \quad (10.75)$$

which is affine.

We consider multi-dimensional quadratic Gaussian models in a fair amount of detail in Chapter 12, so we shall be suitably brief here. The affine connection makes it unsurprising that bond reconstruction formulas exist for the quadratic model. In fact, we have that zero-coupon discount bonds are exponentials of a quadratic function of $y(t)$,

$$P(t, T) = P(t, T; y(t)) = e^{a(t, T) - b(t, T)y(t) - c(t, T)y(t)^2}$$

with the coefficients a, b, c satisfying Riccati ODEs.

In some ways the parameterization (10.73)–(10.74) is more convenient than the general affine specification. For example, with the discount bonds known functions of a *Gaussian* factor $y(t)$, a swap rate — or a swap value — is a known function of a Gaussian random variable, which allows us to price a swaption by a simple one-dimensional Gaussian integration. We return to this topic in Chapter 12.

10.2.7 Numerical Methods for the Affine Short Rate Model

Much of the material on numerical methods for the GSR model applies to the affine short rate processes, so we shall be brief. Turning first to finite difference methods, let us again emphasize that the spatial variable should be set to be $x(t) = r(t) - f(0, t)$ rather than $r(t)$ itself. The dynamics for $x(t)$ can be found in (10.60) and lead to the general derivatives pricing PDE

$$\frac{\partial V}{\partial t} + (\omega(t) - \varkappa(t)x) \frac{\partial V}{\partial x} + \frac{1}{2}\sigma(t)^2 (\xi(t) + \beta x) \frac{\partial^2 V}{\partial x^2} = (x + f(0, t)) V,$$

which can be solved by standard methods, given appropriate terminal and boundary conditions. We refer to Section 10.1.5 for general guidelines. Dimensioning of the spatial dimensions of the finite difference grid by probabilistic means will require estimates for the mean and variance of $x(T)$, with T being the terminal horizon. We can compute these from the moment-generating functions established earlier, or, perhaps more easily, by approximating the SDE for $x(t)$ as being approximately Gaussian. If $r(t)$ is close to a CIR process, we can also use the analytical moment results established in Corollary 8.3.3. When establishing the terminal boundary function (i.e. the option payout), we can rely on the reconstitution formulas in (10.63) to turn values of x in the finite difference lattice into the discount bond prices that are required to evaluate the payout.

As for Monte Carlo methods, many of the principles of Section 10.1.6 continue to apply, and we can draw on material in Chapter 8 to design schemes to advance $x(t)$ through time. To elaborate a bit on this, suppose that we are interested in advancing $x(t)$ from time t_i to time t_{i+1} . Assume that all parameters in (10.60) are piecewise constant, such that

$$dx(t) \approx \varkappa_i (q_i - x(t)) dt + \sigma_i \sqrt{\xi_i + \beta x(t)} dW(t), \quad t \in [t_i, t_{i+1}],$$

where¹³ $\varkappa_i = \varkappa(t_i)$, $q_i = \omega(t_i)/\varkappa(t_i)$, $\sigma_i = \sigma(t_i)$, and $\xi_i = \xi(t_i)$. Defining $y(t) = \xi_i + \beta x(t)$, it follows that we can approximate $x(t_{i+1}) \approx (y(t_{i+1}) - \xi_i)/\beta$, where

$$dy(t) = \varkappa_i (\beta q_i + \xi_i - y(t)) dt + \beta \sigma_i \sqrt{y(t)} dW(t), \quad y(t_i) = \xi_i + \beta x(t_i). \quad (10.76)$$

Simulation of this SDE, however, was discussed in detail in Section 9.5 where a number of practical algorithms were introduced. We should notice that typical parameterizations of (10.76) will rarely violate the Feller condition, making this SDE considerably easier to deal with numerically than the stochastic volatility applications in Section 9.5. Additional material on Monte Carlo simulation of generic short rate processes — most of which also applies to affine processes — can be found in Section 11.3.3.

¹³ Alternatively, we can also set $\varkappa_i = (\varkappa(t_i) + \varkappa(t_{i+1}))/2$, and so forth.

One-Factor Short Rate Models II

While the affine specification (including the Gaussian case) of Chapter 10 is, without doubt, the most popular one-factor short rate model in practice, quite a few other models have been proposed in the literature. In this chapter we cover the most important of these models, paying special attention to the case where the short rate is log-normal. We also briefly discuss some issues in the econometric estimation of short rate models, and introduce the important concept of *unspanned stochastic volatility*.

As most of the models introduced in this chapter have no analytical bond reconstitution formulas, their calibration to the initial term structure requires numerical work. Accordingly, the second half of this chapter is dedicated to numerical methods for pricing and, especially, calibration of models based on generic short rate SDE. Particularly useful in this regard is the discussion in Section 11.3.2 on efficient finite difference schemes based on the important concept of *forward induction*.

11.1 Log-Normal Short Rate Models

Given the pervasiveness of the log-normal Black-Scholes model in derivatives pricing theory, it should come as no surprise that many authors have attempted to specify one-factor short rate models where the dynamics of $r(t)$ are of the form $dr(t) = O(dt) + \sigma_r(t)r(t) dW(t)$ with deterministic $\sigma_r(t)$. This section reviews this class of models which, somewhat surprisingly, turns out to have a number of rather severe drawbacks.

11.1.1 The Black-Derman-Toy Model

The *Black-Derman-Toy* (BDT) model was originally specified in a discrete-time binomial setting in Black et al. [1990], but subsequent research has shown the continuous-time limit of the model to be

$$r(t) = U(t)e^{\sigma_r(t)W(t)}, \quad (11.1)$$

where $U(t)$ and $\sigma_r(t)$ are deterministic functions, and $W(t)$ is a scalar Brownian motion in the risk-neutral measure. Notice that $r(t)$ is here an outright function of $W(t)$, a property sometimes known as *path independence* in the short rate dynamics. The following lemma rewrites the BDT model in more familiar terms.

Lemma 11.1.1. *Let $r(t)$ be given by (11.1), and let the prime denote a time derivative. Then*

$$\begin{aligned} d\ln r(t) &= \left(\vartheta(t) + \frac{\sigma'_r(t)}{\sigma_r(t)} \ln r(t) \right) dt + \sigma_r(t) dW(t), \\ dr(t)/r(t) &= \left(\vartheta(t) + \frac{1}{2}\sigma_r(t)^2 + \frac{\sigma'_r(t)}{\sigma_r(t)} \ln r(t) \right) dt + \sigma_r(t) dW(t), \end{aligned} \quad (11.2)$$

where

$$\vartheta(t) = U'(t)/U(t) - \ln(U(t))\sigma'_r(t)/\sigma_r(t).$$

Proof. Set $y(t) = \ln r(t) = \ln U(t) + \sigma(t)W(t)$ and apply Ito's lemma to get

$$\begin{aligned} dy(t) &= \left(\frac{d\ln U(t)}{dt} + \sigma'_r(t) W(t) \right) dt + \sigma_r(t) dW(t) \\ &= (U'(t)/U(t) + (y(t) - \ln U(t))\sigma'_r(t)/\sigma_r(t)) dt + \sigma_r(t) dW(t), \end{aligned}$$

so that (11.2) follows. A second application of Ito's lemma to $r(t) = e^{y(t)}$ then gives the result for $dr(t)$. \square

Of the various specifications of the BDT model, the most convenient is probably (11.2) which describes the logarithm of $r(t)$ as a mean-reverting Gaussian process with a mean reversion speed

$$\kappa(t) = -\frac{\sigma'_r(t)}{\sigma_r(t)}.$$

As $\ln r(t)$ is Gaussian, $r(t)$ is log-normal. In the formulation (11.2), $\vartheta(t)$ can be considered a free parameter, the value of which can — for a fixed volatility function $\sigma_r(t)$ — be determined by calibrating the model to the initial yield curve. As the BDT model has no known bond reconstitution formula, this fit must be done numerically. The original presentation in Black et al. [1990] outlines such a routine, based on brute-force backward induction in a binomial tree; this algorithm is, however, computationally inefficient and should *not* be used. For a much faster approach, see Sections 11.3.2.1 and 11.3.2.2.

Besides the lack of the bond reconstitution formula, the BDT model is plagued by a number of issues that makes it unsuitable for practical applications. One of the issues is explained in Section 11.1.3. In addition, it

is problematic that the mean reversion speed of the model is beyond user control and is fully determined from the short rate volatility and its time derivative. In particular, for those values of t where $\sigma_r(t)$ grows in t , the mean reversion will be negative, i.e. the model will imply “mean-fleeing” behavior. It should be obvious that this feature of the model is undesirable.

11.1.2 Black-Karasinski Model

In order to rectify the problems surrounding the mean reversion in the BDT model, Black and Karasinski [1991] (BK) took the obvious step of introducing a model where mean reversion for $\ln r(t)$ is exogenously specified. In other words, we write

$$d \ln r(t) = \kappa(t) (\vartheta(t) - \ln r(t)) dt + \sigma_r(t) dW(t). \quad (11.3)$$

Equivalently, we may write

$$r(t) = e^{x(t)},$$

where $x(t)$ is a standard mean-reverting Gaussian process; accordingly, the BK short rate dynamics are straightforward to simulate in the Monte Carlo method. The BK dynamics generalize and improve those of the BDT model but still do not allow for a closed-form discount bond reconstitution formula.

11.1.3 Issues in Log-Normal Models

The BK model — and its special case, the BDT model — have short rates that are log-normally distributed. A similar distribution of rates would arise for risk-neutral dynamics of geometric Brownian motion type

$$dr(t) = \mu_r(t)r(t) dt + \sigma_r(t)r(t) dW(t). \quad (11.4)$$

A time-homogeneous version of this model was considered in Rendleman and Bartter [1980]; for the time-homogeneous case (very complicated) formulas¹ for discount bond prices exist, see Dothan [1978] and Hogan and Weintraub [1993]. The model (11.4) has no mean reversion, and cannot be recommended for practical applications, however.

A common problem shared by all the log-normal short rate models (11.2), (11.3), (11.4) is the fact that the expected value of the inverse of future discount bond price is infinite, i.e.

$$\mathbb{E}_t \left(\frac{1}{P(t', T)} \right) = \infty, \quad t < t' < T. \quad (11.5)$$

¹A computationally efficient recursive procedure for bond pricing can be found in Hansen and Jørgensen [1998].

This result is formally shown in Hogan and Weintraub [1993], but is hardly surprising since the expectation of e^{cX} , $c > 0$, is well-known to be infinite when X is a log-normal random variable². An important corollary of this result is listed below, originally due to Sandmann and Sondermann [1997].

Corollary 11.1.2. *Define a forward Libor rate*

$$L(t, T) = \frac{1}{\tau} \left(\frac{P(t, T)}{P(t, T + \tau)} - 1 \right),$$

where $\tau > 0$ is some accrual factor. Assuming that (11.5) holds, then also

$$\mathbb{E}_t (L(T, T)) = \infty, \quad T > t, \quad (11.6)$$

and

$$\mathbb{E}_t \left(e^{\int_{t'}^T r(u) du} \right) = \infty. \quad (11.7)$$

Proof. As $1 + \tau L(T, T) = 1/P(T, T + \tau)$, equation (11.6) follows directly from (11.5). To show (11.7), we use Jensen's inequality³ to show

$$\frac{1}{P(t', T)} = \frac{1}{\mathbb{E}_{t'} \left(e^{- \int_{t'}^T r(u) du} \right)} \leq \mathbb{E}_{t'} \left(e^{\int_{t'}^T r(u) du} \right). \quad (11.8)$$

Taking expectations conditional on the time t filtration yields (11.7). \square

Both (11.6) and (11.7) have unfortunate economic consequences. Formula (11.7) predicts that the expected return of investing in the continuously compounded money market account for a finite period of time is infinite; and (11.6) predicts that all Libor futures rates should be infinite⁴. A related problem was discussed in the context of log-normal HJM models in Section 4.5.3.

11.1.4 Sandmann-Sondermann Transformation

The problems outlined in Corollary 11.1.2 are a significant drawback of log-normal short rate models, and one that should disqualify their use in many applications. On the other hand, market data may dictate that interest rates should, in fact, be log-normal. This is not as big a dilemma as it may appear, as there are ways to build models with a strong log-normal flair, yet avoiding (11.5). One way is to use HJM models of the type

$$df(t, T) = O(dt) + \sigma(t)r(t) dW(t),$$

²Even though the log-normal distribution has finite moments of all orders, the moment-generating function is infinite at any positive argument.

³For details on Jensen's inequality, see the proof of Proposition 11.1.3.

⁴Recall from Section 4.1.2 that the risk-neutral expectation of a random variable must, in the absence of arbitrage, equal its traded futures price.

where f is, as always, the instantaneous forward rate. We return to this type of models in Chapter 13. An interesting alternative was proposed by Sandmann and Sondermann [1997] and in essence involves writing the log-normal dynamics for a *discretely compounded rate* $r_d(t)$, rather than the infinitesimal (continuously compounded) rate $r(t)$. Specifically, we relate r_d and r through the expression

$$e^{r(t)\delta} = 1 + r_d(t)\delta \implies r(t) = \delta^{-1} \ln(1 + r_d(t)\delta), \quad (11.9)$$

where $\delta > 0$ is some finite compounding interval. Sandmann and Sondermann [1997] set $\delta = 1$, but any finite positive value will, in fact, do.

The effect of working with discretely compounded rates is summarized in the following result.

Proposition 11.1.3. *Let $r_d(t)$ be log-normal for all t . Then, for $t < t' < T$,*

$$\mathbb{E}_t(L(T, T)) < \infty, \quad (11.10)$$

$$\mathbb{E}_t\left(e^{\int_{t'}^T r(u)du}\right) < \infty. \quad (11.11)$$

Proof. As the proof of Proposition 11.1.3 is quite instructive, we give full details. From (11.8), to prove both (11.10) and (11.11) it suffices to show that the expectation of $e^{\int_{t'}^T r(u)du}$ is finite. For this, let us recall that Jensen's inequality for integrals states that for a real-valued function g and a concave function φ

$$\varphi\left(\int_{t'}^T g(u)f(u)du\right) \geq \int_{t'}^T \varphi(g(u))f(u)du, \quad (11.12)$$

provided that $f(u)$ is non-negative and $\int_{t'}^T f(u)du = 1$. If φ is convex, the inequality is reversed. Now write

$$\delta^{-1} \int_{t'}^T \ln(1 + r_d(u)\delta)du = \delta^{-1} \int_{t'}^T \frac{1}{T-t'} \ln((1 + r_d(u)\delta)^{T-t'})du$$

and apply (11.12) to the right-hand side with $f(u) = 1/(T-t')$, $\varphi(u) = \ln(u)$, and $g(u) = (1 + r_d(u)\delta)^{T-t'}$ to show that

$$\delta^{-1} \ln\left(\int_t^T \frac{1}{T-t'} (1 + r_d(u)\delta)^{T-t'} du\right) \geq \delta^{-1} \int_{t'}^T \ln(1 + r_d(u)\delta)du. \quad (11.13)$$

From (11.13) it then follows immediately that

$$\begin{aligned}
& \mathbb{E}_t \left(\exp \left(\int_{t'}^T r(u) du \right) \right) \\
&= \mathbb{E}_t \left(\exp \left(\int_{t'}^T \delta^{-1} \ln (1 + r_d(u)\delta) du \right) \right) \\
&\leq \mathbb{E}_t \left(\exp \left(\delta^{-1} \ln \left(\int_{t'}^T \frac{1}{T-t'} (1 + r_d(u)\delta)^{T-t'} du \right) \right) \right) \\
&= \mathbb{E}_t \left(\left(\frac{1}{T-t'} \int_{t'}^T (1 + r_d(u)\delta)^{T-t'} du \right)^{1/\delta} \right).
\end{aligned}$$

Assume that $0 < \delta < 1$, and set $f(u) = 1/(T - t')$, $\varphi(u) = u^{1/\delta}$, and $g(u) = (1 + r_d(u)\delta)^{T-t'}$. By Jensen's inequality (here (11.12) is reversed, since φ is now convex) we must have

$$\left(\frac{1}{T-t'} \int_{t'}^T (1 + r_d(u)\delta)^{T-t'} du \right)^{1/\delta} \leq \frac{1}{T-t'} \int_{t'}^T (1 + r_d(u)\delta)^{(T-t')/\delta} du. \quad (11.14)$$

Therefore

$$\mathbb{E}_t \left(e^{\int_{t'}^T r(u) du} \right) \leq \frac{1}{T-t'} \mathbb{E}_t \left(\int_{t'}^T (1 + r_d(u)\delta)^{(T-t')/\delta} du \right), \quad 0 < \delta < 1. \quad (11.15)$$

Since finite powers of log-normal random variables have finite expected values, (11.11) has been shown for $0 < \delta < 1$. For $\delta \geq 1$ we have $1/\delta \leq 1$, and

$$\begin{aligned}
& \mathbb{E}_t \left(\left(\frac{1}{T-t'} \int_{t'}^T (1 + r_d(u)\delta)^{T-t'} du \right)^{1/\delta} \right) \\
&\leq \left(\mathbb{E}_t \left(\frac{1}{T-t'} \int_{t'}^T (1 + r_d(u)\delta)^{T-t'} du \right) \right)^{1/\delta},
\end{aligned}$$

and (11.11) follows from the same arguments. \square

Comparison of Corollary 11.1.2 and Proposition 11.1.3 shows that models of the BK type

$$d \ln r_d(t) = \kappa(t) (\vartheta(t) - \ln r_d(t)) dt + \sigma(t) dW(t) \quad (11.16)$$

and geometric Brownian motion models

$$dr_d(t)/r_d(t) = \mu(t) dt + \sigma(t) dW(t), \quad (11.17)$$

become significantly more reasonable when the dynamics are written in $r_d(t)$, rather than $r(t)$. We invite the reader to apply Ito's lemma to the

transformation (11.9) to uncover the r -dynamics for the models (11.16)–(11.17).

Remark 11.1.4. When applied to the HJM model class, the “trick” of shifting from continuously compounded to discretely compounded rates lays the foundation for the class of so-called *Libor market* models. We return to these models in Chapters 14 and 15.

11.2 Other Short Rate Models

11.2.1 Power-Type Models and Empirical Model Estimation

A natural extension of the Gaussian and affine short rate SDEs involves retaining the linear mean-reverting drift term of these models, but using a general power function in the diffusion term. That is, we write

$$dr(t) = \kappa(t)(\vartheta(t) - r(t))dt + \sigma(t)r(t)^p dW(t), \quad p > 0. \quad (11.18)$$

The time-inhomogeneous Gaussian and CIR models correspond to the choices $p = 0$ and $p = 1/2$, respectively. The special case of (11.18) where $p = 1$ was suggested in Brennan and Schwartz [1980] and Courtadon [1982], and is quite similar to the BK model — to the extent that the case $p = 1$ shares⁵ the unfortunate properties of the BK model listed in Corollary 11.1.2.

The general model (11.18) is similar to the CEV model described in Chapter 7, and manipulation of the parameter p may allow for a better fit of the model to observed volatility smiles in interest rate options. Due to its intractability — no bond reconstitution formula exists for $p \notin \{0, 1/2\}$ — the model is, however, rarely used in derivatives pricing applications. (For related, but significantly more tractable, models with power-type diffusion terms, see Chapter 13.) Starting with the influential article by Chan et al. [1992] (CKLS), the specification (11.18) has, however, been quite popular in econometric work. As the CKLS paper is one of the most frequently cited references on estimation of one-factor short rate models, let us make a (very) brief foray into econometrics, to review the CKLS conclusions and some of the criticism their work has subsequently drawn.

The CKLS estimation procedure is based on US Treasury bond data from the period 1964–1989. Assuming that $r(t)$ can be approximated by the one-month yield on US Treasury bonds, they estimate eight models with parameter restrictions, and one model with no restrictions. Generally speaking, the empirical results indicate that the value of the parameter p is the most important in determining whether a model is accepted or rejected.

⁵The proof of this statement is straightforward, and follows from a standard comparison theorem for SDEs (see p. 293 of Karatzas and Shreve [1997]).

The unrestricted estimate of p is close to 1.5, and values less than around $p = 1$ are rejected in their tests. The Vasicek and CIR models are, for instance, rejected, whereas the Brennan-Schwarz/Courtadon model (with $p = 1$) is accepted.

The fact that the CKLS estimates suggest that $p \geq 1$ is surprising in light of the generally downward-sloping volatility skews in fixed income derivatives, and also raises considerable questions about model regularity, as one would expect from Corollary 11.1.2. Indeed, as shown in Honore [1998b], $r(t)$ will a.s. explode (i.e. reach ∞) in finite time if $p > 1$. Beyond this, we notice that the CKLS analysis has received criticism on a number of procedural and data-related fronts. For instance, Honore [1998a] (and quite a few others) point out that the 1 month Treasury yield may be an unreliable proxy for the short rate. Repeating the analysis with a carefully computed value of $r(t)$ (obtained by exploiting the fact that a one-factor model implies a one-to-one correspondence between any discount yield and the short rate), Honore revises the CKLS estimate for p significantly downward, to around $p = 0.8$. Bliss and Smith [1997] also point out that the data used by CKLS covers the period October 1979 — September 1982, when the US Federal Reserve followed unusual monetary policies (“The Fed Experiment”). Properly accounting for this, Bliss and Smith revise the CKLS estimate for p down to around 1.0. Applying different estimation methods and different data sets, Andersen and Lund [1997] and Christensen et al. [2001] estimate p to around 0.7 and 0.8, respectively.

Moving away from observations of only a short rate proxy, Gibbons and Ramaswamy [1993] test the ability of the CIR model to simultaneously describe the evolution in four zero-coupon rates; with data covering the same period as the CKLS study, they accept the hypothesis $p = 1/2$. Finally, to muddy waters even further, Ait-Shalia [1996] points out (in an analysis that has subsequently been criticized as lacking robustness in Chapman and Pearson [2000]) that the simple linear drift term in (11.18) is fundamentally misspecified and should be adjusted to include non-linear terms such as $1/r(t)$ and $r(t)^2$.

By now, it should be clear to the reader that the problem of estimating short rate models is not close to being conclusively solved, despite an impressively long list of papers associated to it. In much contemporary empirical research, the importance of the choice of p is generally downplayed, with the affine class ($p = 1/2$) enjoying considerable current popularity due to its analytical tractability.

11.2.2 The Black Shadow Rate Model

As described in Chapter 10, one drawback of the Gaussian short rate model class is the implication that interest rates can become negative with positive probability. As long as investors can stuff their mattresses with currency, (nominal) interest rates must, however, always remain non-negative

to preclude arbitrage. In practice, the probability of negative rates may be small enough to ignore, but as argued in Rogers [1996] prices of certain contingent claims may be highly sensitive to even a remote probability of negative rates, in which case the Gaussian model should obviously be avoided. Possible model alternatives with non-negative rates include the log-normal models in Section 11.1, as well as the affine models in Section 10.2.

Rather than altering the underlying model, an alternative “fix” of the Gaussian model involves simply taking the positive part of the Gaussian short rate process, i.e. writing

$$r(t) = (r^*(t))^+, \quad (11.19)$$

where $r^*(t)$ is a Gaussian process

$$dr^*(t) = \kappa(t) (\vartheta(t) - r^*(t)) dt + \sigma_r(t) dW(t). \quad (11.20)$$

This approach was first proposed in Black [1995] and Rogers [1995]. The form of (11.19) suggests an analogy where the interest rate $r(t)$ is an *option*, granting a choice between an underlying *shadow short rate* $r^*(t)$ and zero. In other words, whenever an interest rate product has a negative rate, we invest our money in currency instead.

The truncation in (11.19) may at first glance appear rather crude. For instance, the process for $r(t)$ retains full volatility $\sigma_r(t)$ as long as $r(t) > 0$, and can then suddenly get extinguished completely for potentially long stretches of time. In contrast, alternative models with non-negative rates such as BDT, BK, and CIR generally all imply that the interest rate volatility will gradually vanish (linearly for BDT/BK, as a square-root for CIR) as the short rate tends to zero. Interestingly, there is some empirical evidence (from the US in the 1930’s and Japan in the 1990’s) that suggests that very low interest rates are often accompanied by higher absolute rate volatility than standard models would predict, see Goldstein and Keirstead [1997]. Such evidence may lend some credibility to models such as (11.19).

The process (11.19) is not analytically tractable, and numerical methods (such as those of Section 11.3) must be applied to price discount bonds and other fixed income securities. For the case where all parameters in (11.20) are constants — i.e. $r^*(t)$ follows a Vasicek model — Gorovoi and Linetsky [2004] list⁶ a complicated eigenfunction expansion series for discount bond prices. Of course, as the constant-parameter model will not be able to match the current yield curve, the result in Gorovoi and Linetsky [2004] has limited uses in practical applications.

Finally, the reader may very well ask whether perhaps (11.19) could be replaced by a reflecting or absorbing barrier at zero, or by the application of

⁶The authors also develop eigenfunction expansions for cases where the shadow rate $r^*(t)$ follow time-homogeneous diffusions more complicated than the Vasicek model.

a suitable transformation such as $r(t) = r^*(t)^2$. The latter idea was discussed in Section 10.2.6 and the former is investigated in Goldstein and Keirstead [1997] where eigenfunction expansions are derived for the time-homogeneous case. In Black [1995], the author objects to reflecting barriers on economic grounds.

11.2.3 Spanned and Unspanned Stochastic Volatility: the Fong and Vasicek Model

In previous chapters, we demonstrated how to incorporate stochastic volatility as a mechanism to model the volatility smile in vanilla models. One wonders how to proceed with such a construction for term structure models. We postpone much of this discussion to later chapters, but shall here take the opportunity to discuss what constitutes a *true* stochastic volatility model for interest rate evolution, and why the short rate framework is *not* particularly amendable to stochastic volatility extensions. More specifically, we shall introduce the notion of *spanned* and *unspanned* stochastic volatility. For concreteness, our discussion focuses on a model proposed by Fong and Vasicek [1991] which, despite initial appearances, is in fact not a true stochastic volatility model.

The Fong-Vasicek (FV) model is characterized by risk-neutral SDEs

$$\begin{aligned} dr(t) &= \kappa_r (\vartheta_r - r(t)) dt + \sqrt{z(t)} dW_1(t), \\ dz(t) &= \kappa_z (\vartheta_z - z(t)) dt + \eta \sqrt{z(t)} dW_2(t), \end{aligned}$$

where W_1 and W_2 are correlated Brownian motions, $\langle dW_1(t), dW_2(t) \rangle = \rho dt$, and $\kappa_r, \vartheta_r, \kappa_z, \vartheta_z, \eta$ are positive constants. We recognize the FV model as being essentially a time-homogeneous GSR model augmented by a stochastic variance process $z(t)$; the process for $z(t)$ is identical to that of the Heston model (see Chapter 8). Bond prices in the FV model can be shown to satisfy

$$P(t, t + \delta) = e^{A(\delta) + r(t)B(\delta) + z(t)C(\delta)}, \quad (11.21)$$

for deterministic functions A, B, C satisfying a coupled system of ODEs. Details about these ODEs⁷ and their rather complicated analytical solution can be found in Selby and Strickland [1995]. For our purposes here, what matters is not the precise form of A, B , and C , but rather the fact that $P(t, T)$ in the FV model is a deterministic function of the two state variables $r(t)$ and $z(t)$. As a consequence, one can theoretically hedge out exposure to both $r(t)$ and $z(t)$ by simply taking positions in two discount bonds with different maturities. Equivalently, given observations at time t of the prices of two discount bonds with different maturities, we can invert (11.21) to uncover the current values of the two variables $r(t)$ and $z(t)$.

⁷The ODEs are easily derived by substituting the right-hand side of (11.21) into the no-arbitrage PDE for $P(t, T)$ in the FV model.

When a “stochastic volatility” variable $z(t)$ can be hedged by positions in discount bonds, we say that $z(t)$ is *spanned* by the discount curve. If all stochastic volatility variables are spanned by the discount curve, moves in, say, implied volatilities of caps and swaptions would always be accompanied by moves in the yield curve, making vega hedging theoretically unnecessary. In reality, however, there is much evidence that interest rate option volatilities cannot be perfectly hedged by trading only discount bonds; see Casassus et al. [2005] for a review of the literature. Formally, this implies that the volatilities of discount bond prices depend on a vector of random state variables $(z_1(t), \dots, z_n(t))$ that are not included in the state variables used in reconstitution formulas for the discount curve. That is,

$$\frac{\partial P(t, T)}{\partial z_i(t)} = 0, \quad i = 1, \dots, n, \quad (11.22)$$

yet, for $Q(t, T) \equiv E_t(d(P(t, T))^2)/dt$,

$$\frac{\partial Q(t, T)}{\partial z_i(t)} \neq 0, \quad i = 1, \dots, n. \quad (11.23)$$

Random variables satisfying (11.22)–(11.23) are said to represent *unspanned* stochastic volatility (USV). Whenever we talk about true stochastic volatility in this book, we always refer to models with USV, i.e. to models that prescribe moves in rate volatility that cannot be inferred from moves in the level and the shape of the discount curve.

A detailed, and highly recommended, account of USV can be found in Collin-Dufresne and Goldstein [2002b]. Among many results, the authors prove that in a time-homogeneous setting, bivariate affine models for the term structure of interest rates cannot exhibit USV. This explains our results for the FV model, and also demonstrates that several other classical stochastic-volatility models (e.g. Longstaff and Schwartz [1992]) are not in the USV class. Further analysis of a number of stochastic volatility models proposed in the literature (a surprising number of which do not, in fact, allow for USV) can be found in Collin-Dufresne and Goldstein [2002b] and Casassus et al. [2005].

11.3 Numerical Methods for General One-Factor Short Rate Models

The models described in Sections 11.1 and 11.2 all (with the exception of the Fong-Vasicek model) involve risk-neutral SDEs for the short rate of the type

$$dr(t) = \mu_r(t, r(t)) dt + \sigma_r(t, r(t)) dW(t), \quad (11.24)$$

for certain user-specified deterministic functions $\mu_r(t, r)$ and $\sigma_r(t, r)$. As discussed, many of these models do not allow for closed-form expressions

relating discount bonds to the state of $r(t)$, necessitating the application of numerical methods for even simple tasks such as calibrating the model to the initial yield curve. While this issue ultimately should make one pause when it comes to deciding whether a model is suited for practical applications, we still want to cover some techniques that are useful in handling generic models such as (11.24). Some of the techniques we shall discuss, e.g. calibration through forward induction, have broad applicability.

11.3.1 Finite Difference Methods

As always, we set $x(t) = r(t) - f(0, t)$, such that the generic pricing PDE for a derivative $V = V(t, x)$ becomes

$$\frac{\partial V}{\partial t} + \mu_x(t, x) \frac{\partial V}{\partial x} + \frac{1}{2} \sigma_x(t, x)^2 \frac{\partial^2 V}{\partial x^2} = (x + f(0, t)) V, \quad (11.25)$$

where we have defined

$$\mu_x(t, x) = \mu_r(t, x + f(0, t)) - \frac{\partial f(0, t)}{\partial t}, \quad \sigma_x(t, x) = \sigma_r(t, x + f(0, t)). \quad (11.26)$$

Given terminal conditions for $V(T, x)$, as well as suitable boundary conditions in the x -domain, we can solve this equation by the generic finite difference methods of Chapter 2. For later use, let us quickly recall that a standard θ -method finite difference scheme on an equidistant x -grid $\{x_j\}_{j=0}^{m+1}$ would result in a matrix scheme of the type

$$[\mathbf{I} - \theta \Delta_t \mathbf{A} ((1 - \theta)t_{i+1} + \theta t_i)] \widehat{\mathbf{V}}(t_i) = \\ [\mathbf{I} + (1 - \theta) \Delta_t \mathbf{A} ((1 - \theta)t_{i+1} + \theta t_i)] \widehat{\mathbf{V}}(t_{i+1}) + \mathbf{B}(t_i, t_{i+1}), \quad (11.27)$$

where $\widehat{\mathbf{V}}(t) = (\widehat{V}(t, x_1), \dots, \widehat{V}(t, x_m))^T$ with $\widehat{V}(t, x)$ denoting the approximation to the true solution $V(t, x)$, $\mathbf{A}(t)$ is a tri-diagonal matrix, and $\mathbf{B}(t_i, t_{i+1})$ is a vector representing any boundary conditions that cannot be folded into the matrix \mathbf{A} . We solve the matrix system (11.27) backward in time, starting from a given value of $\widehat{\mathbf{V}}(T)$. Determination of the spatial boundaries in the x -domain — that is, the values of x_0 and x_{m+1} — can, as always, be set by probabilistic arguments, through estimation of the first and second moment of $x(T)$. While these estimates should be targeted to the specifics of the model at hand, if all else fails we can always rely on Gaussian estimate, e.g. something like

$$x_0 = \mu_x(0, 0) T - \alpha \sigma_x(0, 0) \sqrt{T}, \quad x_{m+1} = \mu_x(0, 0) T + \alpha \sigma_x(0, 0) \sqrt{T},$$

where α is some confidence interval multiplier (e.g. 4 or 5).

The discussion in Section 10.1.5.2 about x -domain boundary conditions apply to (11.27) as well, but determining the terminal condition $\widehat{\mathbf{V}}(T)$ can

be problematic⁸, as option payouts will often involve several discount factors at time T (e.g. to price a swap or to compute a Libor rate). As there is generally no way of computing such discount bonds analytically for the model (11.24), we are forced to compute the discount bond prices themselves by finite difference methods. To compute the value of the discount bond $P(T, T^*, x)$, $T^* > T$, at time T we:

1. Extend the finite difference grid to time T^* .
2. Set the boundary condition $P(T^*, T^*, x_j) = 1$, $j = 0, \dots, m + 1$.
3. On some suitable time grid, use (11.27) on P to step backward to time T .
4. Use the finite difference estimates $\hat{P}(T, T^*, x_j)$, $j = 0, \dots, m + 1$, to fill in $\hat{\mathbf{V}}(T)$.

To the extent that $\mathbf{V}(T)$ involves multiple discount bonds, we perform the algorithm above on all of the required discount bonds; the grid must then be extended to the maturity of the longest-dated discount bond needed in the payout computation. In some cases this can dramatically increase computation times relative to models where closed-form discount bond reconstitution formulas exist. For instance, for a 3 month option on a 30 year swap, a model with an analytical formula for discount bonds (e.g. the affine short rate model) would require us only to build the finite difference grid out to 3 months; when such a formula does not exist, we are forced to use a 30 year finite difference grid.

11.3.2 Calibration to Initial Yield Curve

Assume that the volatility function $\sigma_r(t, r)$ has been fixed, but $\mu_r(t, r) = \mu_r(t, r; \vartheta(t))$ has a free time-dependent parameter $\vartheta(t)$, the value of which we wish to set in such a way that the initial discount curve $P(0, T)$, $T > 0$, is correctly recovered by the model. Suppose that we assume, as is common, that $\vartheta(t)$ is piecewise constant on some time grid $0 = t_0 < t_1 < \dots < t_N$, with ϑ_{i-1} denoting the value of ϑ that applies on the interval $[t_{i-1}, t_i]$. A brute force approach to the calibration of $\vartheta(t)$ could proceed as follows.

1. Assuming that $\vartheta_0, \dots, \vartheta_{i-2}$ are known, make a guess⁹ for ϑ_{i-1} .
2. Setting the terminal boundary value to $V(t_i, x) = 1$ for all x , use the backward finite difference grid algorithm (11.27) to compute the value of $P(0, t_i)$.
3. If the computed value of $P(0, t_i)$ equals that quoted in the market stop; otherwise return to Step 1.

⁸The same holds for exercise values of callable securities.

⁹An initial guess for ϑ_{i-1} could be $\vartheta_{i-1} = \vartheta_{i-2}$. Subsequent guesses would be performed by a root-search algorithm.

This approach is similar to an algorithm proposed in Black et al. [1990] (albeit the authors worked only with binomial trees) and involves very high computational costs as the numerical search for each of the parameters $\vartheta_0, \vartheta_1, \dots$ involves solving a full finite difference grid in each loop. Specifically, assume that on average each search for ϑ_i involves M iterations over steps 1–3 above. With N ϑ -values to find, the computational effort of a finite difference grid with m spatial steps (say) is $O(Nm)$; it follows that the total computational cost for the calibration of $\vartheta(t)$ is

$$O(MN^2m),$$

which is often prohibitively expensive in practice.

11.3.2.1 Forward Induction

In the setting of binomial and trinomial trees, the brute-force BDT algorithm was markedly improved upon in Jamshidian [1991a], using a technique known as *forward induction*. The basic idea is to work with a forward equation, rather than with the backward equation (11.27). Two varieties of this approach are feasible in a finite difference setting, depending on whether the forward equation is developed by direct discretization of the continuous-time Fokker-Planck equation, or by rearrangement of a discretized backward equation. We cover the former approach here, and the latter in Section 11.3.2.2 below.

Let $G(t, x; s, y)$, $s \geq t$, denote the time t value of a security that pays out a Dirac delta amount iff $x(s) = y$, given that $x(t) = x$. Clearly then

$$P(0, T) = \int_{-\infty}^{\infty} G(0, 0; T, y) dy. \quad (11.28)$$

We already saw this financial contract, the so-called Arrow-Debreu security, at the end of Section 1.8. Its price G — being the value of a perfectly valid derivative contract — certainly satisfies a backward Kolmogorov equation

$$\frac{\partial G}{\partial t} + \mu_x(t, x; \vartheta(t)) \frac{\partial G}{\partial x} + \frac{1}{2} \sigma_x(t, x)^2 \frac{\partial^2 G}{\partial x^2} = (x + f(0, t)) G, \quad (s, y) \text{ fixed}, \quad (11.29)$$

with the terminal value condition $G(s, y; s, y) = \delta(y)$. Clearly $G(t, z; s, y)$ is closely related to the transition density $p(t, z; s, y)$ defined in Section 1.8, and can be expected to satisfy a forward Kolmogorov equation, too. The correct equation is

$$\begin{aligned} -\frac{\partial G}{\partial s} - \frac{\partial}{\partial y} (\mu_x(s, y; \vartheta(s))G) + \frac{1}{2} \frac{\partial^2}{\partial y^2} (\sigma_x(t, y)^2 G) \\ = (y + f(0, t)) G, \quad (t, x) \text{ fixed}, \end{aligned} \quad (11.30)$$

subject to the initial condition $G(t, x; t, x) = \delta(x)$. This PDE is identical to the Fokker-Planck equation listed in Section 1.8, except for the fact that the right-hand side is not zero, but contains a discounting term.

Fixing $(t, x) = (0, 0)$, the PDE (11.30) can be discretized by finite difference methods in standard fashion, although we keep in mind that the PDE is to be solved *forward* in time from its initial (Dirac) condition, rather than backward. In most cases¹⁰, the appropriate spatial boundary conditions for the finite difference solution \widehat{G} are

$$\widehat{G}(0, 0; s, y_0) = \widehat{G}(0, 0; s, y_{m+1}) = 0,$$

which assumes that y_0 and y_{m+1} have been set such that the probability density at these levels is negligible. In a discrete setting, the initial condition $G(t, x; t, x) = \delta(x)$ is translated to $\widehat{G}(0, 0; 0, y) = (\Delta y)^{-1} 1_{\{y=0\}}$ for Δy suitably defined, in agreement with the averaging principles of Section 2.5.2. Due to the strongly discontinuous initial condition, Rannacher stepping (see Section 2.5) should always be used.

We are now ready to state our revised algorithm for calibration of $\vartheta(t)$, working again with the assumption that $\vartheta(t)$ is piecewise constant on a time grid $0 = t_0 < t_1 < \dots < t_N$. We assume that $\vartheta_0, \dots, \vartheta_{i-2}$ have been found, as has been $\widehat{G}(0, 0; t_{i-1}, y)$ for all $y \in \{y_j\}_{j=1}^m$ in the finite difference grid.

1. Make a guess for ϑ_{i-1} .
2. Solve (11.30) one time step forward, to time t_i , and save $\widehat{G}(0, 0; t_i, y_j)$, $j = 1, \dots, m$.
3. Compute the discount bond price $P(0, t_i) = \sum_j \widehat{G}(0, 0; t_i, y_j)(y_j - y_{j-1})$.
4. If the computed value of $P(0, t_i)$ equals that quoted in the market stop; otherwise return to Step 1.

The cost of Steps 2 and 3 in this algorithm is $O(m)$, where m is the number of points in the y -direction of the finite difference grid. The total computational effort of this algorithm is therefore

$$O(MmN),$$

where M is the average number of root search iterations over Steps 1–4 above. We recall that the effort of the brute-force backward equation approach was $O(MmN^2)$, so the use of forward induction saves us a factor of N . As N is often of the order of $N = 100$, these savings can be considerable. In typical applications, M is often in the order $M = 2$ to 4 , so calibration of the model (11.24) to the initial discount curve should only be a few times slower than pricing by finite difference methods a single option maturing at time t_N (the cost of which we recall to be $O(mN)$).

¹⁰For some models, the density can grow to infinity at the boundary, notably in the CIR model for $x \rightarrow 0+$ when the Feller condition is violated. Should that be the case, a more careful analysis of boundary conditions is required, see e.g. Lucic [2008].

11.3.2.2 Forward-from-Backward Induction

The backward and forward Kolmogorov equations (11.29)–(11.30) are consistent in the continuous-time limit, but not necessarily so when discretized, finite difference style. As a result, the function $\vartheta(t)$ uncovered from the algorithm in Section 11.3.2.1 will generally not allow a finite difference grid based on the backward equation (11.25) to recover the initial term structure of discount bonds without errors, even if discretized on t - and x -grids that are identical to those used for the calibration of $\vartheta(t)$. As long as the time line is sufficiently finely spaced, the errors tend to be very small, however, and rarely a cause for concern. Nevertheless, it should be noted that it is, in fact, possible to restate the forward induction algorithm in such a way that the algorithm becomes precisely compatible with the brute-force backward equation approach we discussed earlier.

To develop this approach, we start out with the discretized backward equation (11.27) and rearrange it to yield

$$\widehat{\mathbf{V}}(t_i) = \mathbf{T}_i^{i+1} \widehat{\mathbf{V}}(t_{i+1}) + \mathbf{G}_i^{i+1}, \quad (11.31)$$

where, with $\mathbf{A}_i^{i+1} \triangleq \mathbf{A}((1-\theta)t_{i+1} + \theta t_i)$,

$$\begin{aligned} \mathbf{T}_i^{i+1} &\triangleq [\mathbf{I} - \theta \Delta_t \mathbf{A}_i^{i+1}]^{-1} [\mathbf{I} + (1-\theta) \Delta_t \mathbf{A}_i^{i+1}] \\ &= \mathbf{I} + [\mathbf{I} - \theta \Delta_t \mathbf{A}_i^{i+1}]^{-1} \mathbf{A}_i^{i+1} \Delta_t, \end{aligned}$$

and

$$\mathbf{G}_i^{i+1} = [\mathbf{I} - \theta \Delta_t \mathbf{A}_i^{i+1}]^{-1} \mathbf{B}(t_i, t_{i+1}).$$

Repeated application of (11.31) yields, for some l ,

$$\widehat{\mathbf{V}}(0) = \mathbf{T}_0^l \widehat{\mathbf{V}}(t_l) + \mathbf{G}_0^l, \quad (11.32)$$

where \mathbf{T}_0^l and \mathbf{G}_0^l can be found iteratively from the equations

$$\mathbf{T}_0^{i+1} = \mathbf{T}_0^i \mathbf{T}_i^{i+1}, \quad \mathbf{T}_0^0 = \mathbf{I}, \quad (11.33)$$

$$\mathbf{G}_0^{i+1} = \mathbf{T}_0^i \mathbf{G}_i^{i+1} + \mathbf{G}_0^i, \quad \mathbf{G}_0^0 = \mathbf{0}. \quad (11.34)$$

Before we proceed, let us introduce some notation. First, we let $\mathbf{1}_j$ denote an m -dimensional column vector with j -th element equal to 1 and all other elements equal to zero. Also we set $\mathbf{1}$ to mean a column vector with all elements equal to 1. Consider now the zero-coupon bond maturing at time t_i , and assume that for the grid $\{x_j\}_{j=0}^{m+1}$ the initial value of $x = 0$ sits in position β , i.e. $x_\beta = 0$. By the definition of a discount bond, (11.32) allows us to write

$$P(0, t_i) = \mathbf{1}_\beta^\top (\mathbf{T}_0^i \mathbf{1} + \mathbf{G}_0^i) = \mathbf{1}^\top \mathbf{D}_0^i + g_0^i, \quad (11.35)$$

where \mathbf{D}_0^i is an m -dimensional vector and g_0^i is a scalar:

$$\mathbf{D}_0^i = (\mathbf{T}_0^i)^\top \mathbf{1}_\beta, \quad g_0^i = (\mathbf{G}_0^i)^\top \mathbf{1}_\beta.$$

Assuming that the influence from the finite difference grid boundary is small, we evidently have

$$P(0, t_i) \approx \sum_{j=1}^m (\mathbf{D}_0^i)_j,$$

where $(\mathbf{D}_0^i)_j$ denotes the j -th element of \mathbf{D}_0^i . Comparison with (11.28) shows that \mathbf{D}_0^i can be interpreted as the discrete-time Arrow-Debreu security vector for maturity t_i (up to the scaling Δx). From (11.33)–(11.34) we have

$$\mathbf{D}_0^{i+1} = (\mathbf{T}_0^i \mathbf{T}_i^{i+1})^\top \mathbf{1}_\beta = (\mathbf{T}_i^{i+1})^\top \mathbf{D}_0^i, \quad \mathbf{D}_0^0 = \mathbf{1}_\beta, \quad (11.36)$$

$$g_0^{i+1} = (\mathbf{T}_0^i \mathbf{G}_i^{i+1} + \mathbf{G}_0^i)^\top \mathbf{1}_\beta = (\mathbf{G}_i^{i+1})^\top \mathbf{D}_0^i + g_0^i, \quad g_0^0 = 0. \quad (11.37)$$

Recalling the definition of \mathbf{T}_i^{i+1} , it follows that

$$\begin{aligned} \mathbf{D}_0^{i+1} &= (\mathbf{T}_i^{i+1})^\top \mathbf{D}_0^i \\ &= [\mathbf{I} + (1 - \theta) \Delta_t \mathbf{A}_i^{i+1}]^\top [\mathbf{I} - \theta \Delta_t (\mathbf{A}_i^{i+1})^\top]^{-1} \mathbf{D}_0^i \\ &\triangleq [\mathbf{I} + (1 - \theta) \Delta_t (\mathbf{A}_i^{i+1})^\top] \mathbf{Y}_0^i, \end{aligned} \quad (11.38)$$

where \mathbf{Y}_0^i is a vector satisfying a tri-diagonal matrix equation

$$[\mathbf{I} - \theta \Delta_t (\mathbf{A}_i^{i+1})^\top] \mathbf{Y}_0^i = \mathbf{D}_0^i. \quad (11.39)$$

With this, we are ready to state our revised fitting algorithm for $\vartheta(t)$. We assume that $\vartheta_0, \dots, \vartheta_{i-2}$ have been found, as has \mathbf{D}_0^{i-1} and g_0^{i-1} .

1. Make a guess for ϑ_{i-1} .
2. Solve (11.39) for \mathbf{Y}_0^{i-1} .
3. Compute \mathbf{D}_0^i from (11.38); and g_0^i from (11.37).
4. Compute the discount bond price $P(0, t_i)$ from (11.35).
5. If the computed value of $P(0, t_i)$ equals that quoted in the market, then stop; otherwise return to Step 1.

The computational efforts for Steps 2, 3, and 4 are all $O(m)$, so the algorithm above is of the same computational complexity as the algorithm in Section 11.3.2.1.

11.3.2.3 Yield Curve and Volatility Calibration

Volatility calibration of a general short rate model of the type (11.24) is a rather involved affair. The typical scheme — moving model volatilities around until the prices of calibration targets match the market — is beset with complications such as

- Bond reconstitution formulas are unavailable so the model needs to be numerically recalibrated to the initial yield curve after each update of the model volatilities.
- All bonds used in the payoffs of calibration targets need to be computed numerically for each volatility update.
- Bond and calibration target values depend on the entire volatility function, making decoupling of individual target calibrations difficult.

These difficulties, however, have not prevented some major investment banks from risk-managing large derivatives portfolios with such a setup. As we cannot possibly do justice to all the tricks that would be required to make this operational, we content ourselves with presenting a mere outline of a possible algorithm.

We consider the model of the type (11.24) but, for notational convenience, write the volatility term in a separable form

$$dr(t) = \mu_r(t, r(t); \vartheta(t)) dt + \sigma_r(t)\psi(r(t)) dW(t). \quad (11.40)$$

Here, the purely-time-dependent function $\sigma_r(t)$ is used to calibrate to swaptions, and $\vartheta(t)$ is used to match the initial yield curve. Implicitly, $\vartheta(t)$ depends on $\sigma_r(t)$.

As in Section 10.1.4, we assume that we are given a collection of $N - 1$ swaptions defined on a maturity grid $0 = T_0 < T_1 < \dots < T_N$ such that the i -th swaption expires at times T_i , $i = 1, \dots, N - 1$. For concreteness assume that all underlying swaps mature on T_N . We further assume that $\sigma_r(t)$ is discretized in a piecewise constant manner on the maturity grid, with σ_i denoting the flat value on $[T_i, T_{i+1})$.

Before discussing calibration, let us outline an efficient algorithm for pricing all swaptions in the calibration set given a collection of volatilities $\sigma_0, \dots, \sigma_{N-1}$. We implicitly assume that the model is rewritten using a state variable x as in (11.25), and the x -domain is discretized with $\{x_j\}_{j=0}^{m+1}$. As the algorithm involves both forward and backward induction, it is important to follow the approach of Section 11.3.2.2 and use a forward algorithm that is fully compatible with the backward one.

1. Update the volatility function of the model with the new values $\sigma_0, \dots, \sigma_{N-1}$.
2. Using the forward induction algorithm from Section 11.3.2.2, calibrate $\vartheta(t)$ for all $t \in [0, T_N]$.
3. On Step 2, Arrow-Debreu prices $\widehat{G}(0, 0; T_i, \cdot)$ are calculated; *save them* for all $i = 1, \dots, N - 1$.
4. For each $n = N, \dots, 2$:
 - a) Create a new copy of the finite difference grid.
 - b) Populate a payoff $P(T_n, T_n) = 1$ at time T_n .
 - c) Calculate $P(T_i, T_n)$ from $P(T_{i+1}, T_n)$ by backward induction for $i = n - 1, \dots, 1$.

5. For each $i = 1, \dots, N - 1$:

- a) Create the T_i -expiry swaption payoff from $P(T_i, T_n)$, $n = i+1, \dots, N$, calculated on Step 4.
- b) Integrate the payoff against $\hat{G}(0, 0; T_i, \cdot)$ stored on Step 3.
- c) This gives the value of the T_i -expiry swaption.

This algorithm describes how to map a set of model volatilities $\sigma_0, \dots, \sigma_{N-1}$ into model prices of calibration targets. In principle, one can now perform a multi-dimensional optimization to match swaption prices from the model to the market. Given that the number of swaptions could be large — 30 or 40 would not be uncommon — and each valuation is rather expensive, the resulting algorithm, while not necessarily completely impractical, would require significant computational resources.

A further improvement entails adopting the iterative bootstrap algorithm outlined in Section 10.2.5. We recall that the main premise of this algorithm approach was that the value of the T_i -expiry swaption depended on $\sigma_r(s)$ for $s \in [0, T_i]$ in a much stronger way than on $\sigma_r(s)$ for $s \in [T_i, T_N]$; this tends to also be true for the model of the type (11.40) for a wide range of volatility specifications. To incorporate this observation into the algorithm above, we would work our way forward from $i = 0$ and determine σ_i by one-dimensional root-search to match the market value of the T_{i+1} -expiry swaption; all values σ_j , $j \neq i$ would be kept constant in the root search. The bootstrap loop would work its way from $i = 0$ to $i = N - 1$ and would be repeated a few times until convergence, in the manner described in Section 10.2.5. We trust that the reader gets the idea and can fill in remaining details.

11.3.2.4 The Dybvig Parameterization

In some models, it may be the case that the SDE (11.24) can be reformulated as

$$r(t) = s(t) + \vartheta(t), \quad (11.41)$$

where $\vartheta(t)$ is a free time-dependent function to be fitted against discount bond prices, and $s(t)$ satisfies an SDE

$$ds(t) = \mu_s(t, s(t)) dt + \sigma_s(t, s(t)) dW(t). \quad (11.42)$$

We notice that

$$P(0, T) = E \left(e^{- \int_0^T r(u) du} \right) = e^{- \int_0^T \vartheta(u) du} E \left(e^{- \int_0^T s(u) du} \right),$$

such that

$$\int_0^T \vartheta(u) du = \ln \left(\frac{E \left(e^{- \int_0^T s(u) du} \right)}{P(0, T)} \right). \quad (11.43)$$

To the extent that the numerator in the right-hand side of (11.43) is easy to compute — e.g. if the SDE for $s(t)$ permits a closed-form solution — calibration of $\vartheta(t)$ can conveniently be found by direct differentiation, see (11.44).

The specification in (11.41) was proposed in Dybvig [1997] and is, as we have already seen in Section 10.1.1.2, quite natural in the context of Gaussian models. For non-Gaussian models, a “fudge” approach in (11.41) may be less desirable, as the domain of $r(t)$ is hard to control. For instance, suppose that $s(t)$ is a time-homogeneous CIR process

$$ds(t) = \kappa(s_0 - s(t)) dt + \sigma \sqrt{s(t)} dW(t),$$

which is guaranteed to produce only non-negative values of $s(t)$. The combined process¹¹ $r(t) = s(t) + \vartheta(t)$, however, will have domain $r(t) \in [\vartheta(t), \infty)$ which is rather awkward as $\vartheta(t)$ is largely out of the user’s control. This is reflected in the SDE for $r(t)$, where now $\vartheta(t)$ enters into the volatility term:

$$dr(t) = \kappa(s_0 + \vartheta(t) + \vartheta'(t)/\kappa - r(t)) dt + \sigma \sqrt{r(t) - \vartheta(t)} dW(t),$$

where $\vartheta'(t) = d\vartheta(t)/dt$. By affecting the short rate volatility, the interpretation of $\vartheta(t)$ as only serving to fit the yield curve can no longer be maintained. This conclusion holds not only for CIR models, but for all models where σ_s depends on $s(t)$, since

$$dr(t) = \mu_s(t, r(t) - \vartheta(t)) dt + \vartheta'(t) dt + \sigma_s(t, r(t) - \vartheta(t)) dW(t).$$

While sometimes very convenient, the Dybvig parameterization should consequently be approached with considerable care.

11.3.2.5 Link to HJM Models

By construction the Dybvig parameterization in Section 11.3.2.4 ensures that the model is calibrated to the initial forward curve $f(0, t)$, $t > 0$. As the resulting model is driven by Brownian motions, we know that it must be in the HJM class. An interesting question arises: what is the type of an HJM model that is defined through the Dybvig procedure. To answer this, let $s(t)$ satisfy (11.42) and define

$$Q(t, T, s) = \mathbb{E} \left(e^{-\int_t^T s(u) du} \middle| s(t) = s \right).$$

From (11.43), we then get

¹¹This model has been advocated by Brigo and Mercurio [2001] as an easy-to-implement alternative to a true time-dependent CIR process. For reasons explained above, this model has certain drawbacks that require careful evaluation.

$$\vartheta(t) = \frac{\partial}{\partial t} \ln \left(\frac{Q(0, t, s(0))}{P(0, t)} \right) \quad (11.44)$$

and

$$P(t, T) = e^{-\int_t^T \vartheta(u) du} Q(t, T, s(t)) = \frac{Q(0, t, s(0))}{Q(0, T, s(0))} \frac{P(0, T)}{P(0, t)} Q(t, T, s(t)).$$

As

$$f(t, T) = -\frac{\partial}{\partial T} \ln P(t, T) = -\frac{\partial}{\partial T} Q(t, T, s(t)) + \vartheta(T),$$

it follows that

$$df(t, T) = O(dt) + \sigma_f(t, T) dW(t),$$

where

$$\begin{aligned} \sigma_f(t, T) &= -\frac{\partial^2}{\partial T \partial s} Q(t, T, s(t)) \sigma_s(t, s(t)) \\ &= -\frac{\partial^2}{\partial T \partial s} Q(t, T, f(t, t) - \vartheta(t)) \sigma_s(t, f(t, t) - \vartheta(t)). \end{aligned}$$

At time t , the forward rate volatility structure generated by the short rate model evidently depends on the forward curve at time t (through $f(t, t)$) as well as the function $\vartheta(t)$. As $\vartheta(t)$ depends (through (11.44)) on the forward curve at time 0, it follows that the HJM dynamics here generally have “memory” of the initial condition at time 0. If one were to alter these initial conditions, the form of HJM dynamics would fundamentally change.

Looking back to Section 10.1.2.2 where the Gaussian short rate model was developed from a separable HJM model, no dependence in the HJM dynamics on initial conditions arose. Hence, it is clear that not *all* short rate models generate “memory” in the HJM dynamics. This raises the obvious question: under which circumstances will a finite-dimensional Markov HJM model have no dependence in its dynamics of the initial forward curve at time 0? The answer to this is listed in Filipovic and Teichmann [2004] which shows that “essentially” all such models must be time-inhomogeneous *affine* models. There are considerable technical details involved in the exact statement of the result, all of which can be found in Filipovic and Teichmann [2004].

11.3.2.6 The Hagan and Woodward Parameterization

Many of the short rate models considered so far have a free time-dependent parameter in the *drift* of the state variable dynamics; the forward induction algorithm can then be used to set this parameter to match the initial yield curve. Hagan and Woodward [1999a] propose an interesting twist on the idea, where the free parameter is introduced into the *numeraire*.

Hagan and Woodward [1999a] start with an observation that only two ingredients are required to create an interest rate model:

- A set of stochastic processes that drive the evolution of interest rates.
- A positive-valued process that is used as a numeraire.

Once the numeraire is specified, the values of all instruments are recovered by the standard pricing formula, see Chapter 1. Critically, the numeraire does not need to be the money market account, or any other “identifiable” security such as a discount bond or an annuity — a positive process is all that is required (so in fact we need not a numeraire but a *deflator* as defined in Section 1.3, but in this section we use the two terms interchangeably).

We define the stochastic process that drives interest rates by a general one-dimensional¹² process

$$dx(t) = \mu_x(t, x(t)) dt + \sigma_x(t, x(t)) dW(t), \quad x(0) = 0. \quad (11.45)$$

Furthermore, we choose the deflator to be a function of the state variable $x(t)$. Without loss of generality, we specify

$$N(t) = \frac{1}{P(0, t)} e^{h(t, x(t)) + a(t)}. \quad (11.46)$$

Here $h(t, x)$ is user-specified, and $a(t)$ is used to fit the model to the initial yield curve. It is often natural to normalize the parameters such that

$$h(t, 0) = 0, \quad a(0) = 0.$$

Once the deflator is specified, we assume that (11.45) is, in fact, given under the measure Q^N associated with this deflator.

Let E^N denote expectations in measure Q^N . The time t price of a T -maturity discount bond in this model, as a function of the state variable x , is given by

$$\begin{aligned} P(t, T, x) &= N(t) E^N(N(T)^{-1} | x(t) = x) \\ &= \frac{P(0, T)}{P(0, t)} e^{h(t, x) + a(t) - a(T)} E_t^N \left(e^{-h(T, x(T))} \right). \end{aligned} \quad (11.47)$$

Consistency with the initial yield curve requires

$$P(0, T, x(0)) = P(0, T), \quad T \geq 0,$$

and we obtain the following condition on $a(T)$,

$$a(T) = \ln E^N \left(e^{-h(T, x(T))} \right), \quad T \geq 0. \quad (11.48)$$

Hagan and Woodward [1999a] show that if the model (11.45)–(11.46) is consistent with the initial yield curve, i.e. the condition (11.48) is satisfied, the model is in fact arbitrage free.

¹²Multi-dimensional extensions are possible.

To obtain $a(t)$ from (11.48), one should not solve for $\mathbb{E}^N(e^{-h(T,x(T))})$ with a backward PDE. Instead, similarly to Section 11.3.2, one should use forward PDE to obtain $p(t, x)$, the density $\mathbb{Q}^N(x(t) \in dx)/dx$, for $(t, x) \in \mathbb{R}^+ \times \mathbb{R}$. The forward Kolmogorov equation states that

$$-\frac{\partial p}{\partial t}(t, x) - \frac{\partial}{\partial x}(\mu_x(t, x)p(t, x)) + \frac{1}{2}\frac{\partial^2}{\partial x^2}(\sigma_x(t, x)^2 p(t, x)) = 0, \quad p(0, x) = \delta(x).$$

Once $p(t, x)$ is determined, we obtain

$$a(T) = \ln \int e^{-h(T,x)} p(T, x) \, dx, \quad T \geq 0,$$

where the integral is taken over the range of the random variable $x(T)$. Interestingly, the calibration of $a(t)$ is independent of the initial yield curve $P(0, T)$, $T \geq 0$. The calibration is somewhat faster than the forward induction algorithm of Section 11.3.2.1 for general short rate models as it requires only a *single* forward pass of the finite difference scheme.

Zero-coupon discount bonds are obtained via (11.47), i.e., generally, numerically. For special choices of $\mu_x(t, x)$, $\sigma_x(t, x)$ and $h(t, x)$, closed-form formulas could be available.

Let us define $\gamma(t, x; T)$ by

$$\gamma(t, x(t); T) = \mathbb{E}_t^N \left(e^{-h(T,x(T))} \right),$$

so that

$$P(t, T, x) = \frac{P(0, T)}{P(0, t)} e^{h(t,x)+a(t)-a(T)} \gamma(t, x; T).$$

Then instantaneous forward rates are given by

$$f(t, T, x) = f(0, T) + a'(T) - \frac{\partial \ln \gamma(t, x; T)}{\partial T}.$$

Applying Ito's lemma to $e^{-h(T,x(T))}$ and setting $T = t$, we observe that the short rate $r(t) = f(t, t, x(t)) = r(t, x(t))$ is given by

$$r(t, x) = f(0, t) + a'(t) - \left. \frac{\partial \ln \gamma(t, x; T)}{\partial T} \right|_{T=t} \quad (11.49)$$

$$\begin{aligned} &= f(0, t) + a'(t) + \frac{\partial h(t, x)}{\partial t} + \mu_x(t, x) \frac{\partial h(t, x)}{\partial x} \\ &\quad + \frac{1}{2} \sigma_x(t, x)^2 \left(\frac{\partial^2 h(t, x)}{\partial x^2} - \left(\frac{\partial h(t, x)}{\partial x} \right)^2 \right). \end{aligned} \quad (11.50)$$

To make matters more concrete, let us now specialize to the case where $h(t, x) = h(t)x$ and $\mu_x(t, x) = 0$ in the general framework (11.45)–(11.46). One can show that this class of models includes the one-factor Gaussian

short rate model and the affine models. If one relaxes the requirement that $\mu_x(t, x) \equiv 0$, then the BK model is also in the class. With this restricted parameterization (11.50) yields, after ignoring small convexity terms,

$$r(t) \approx f(0, t) + a'(t) + h'(t)x(t), \quad (11.51)$$

and, approximately,

$$dr(t) \approx (\vartheta_r(t) - \varkappa_r(t)r(t)) dt + \sigma_r(t, r(t)) dW(t),$$

where

$$\varkappa_r(t) = -\frac{h''(t)}{h'(t)}, \quad (11.52)$$

with $\vartheta_r(t)$, $\sigma_r(t, r)$ appropriately defined. Hence, the numeraire scaling $h(t)$ can be conceptually linked to the mean reversion parameter for the short rate.

As a practical example of the general approach, Hagan and Woodward [1999a] propose the following class of “ $\beta - \eta$ ” models:

$$dx(t) = \lambda(t)(1 + \beta x(t))^\eta dW(t),$$

$$N(t) = \frac{1}{P(0, t)} e^{h(t)x(t) + a(t)}.$$

For this specification, the transition density of $x(t)$ is known in closed form, allowing for (more or less) analytical calibration to the initial yield curve. The parameters β , η are used to match the skew of the volatility smile. Note the resemblance between the volatility term in this model and that of a vanilla displaced-CEV model in Section 7.2.4. As pointed out in that section, adding a displacement to the CEV function does not significantly alter the range of available volatility smiles. Hence, should this approach be pursued, we recommend the specialization of the $\beta - \eta$ model with $\eta = 1$, and perhaps an extension of the skew parameter β to be time dependent. Specifying a constant mean reversion $\varkappa_r(t) \equiv \varkappa_r$ and solving (11.52) for $h(t)$ (which we normalize, conveniently and without loss of generality, by $h(0) = 0$ and $h'(0) = 1$), we obtain

$$dx(t) = \lambda(t)(1 + \beta(t)x(t)) dW(t),$$

$$N(t) = \frac{1}{P(0, t)} \exp \left(\frac{1 - e^{-\varkappa_r t}}{\varkappa_r} x(t) + a(t) \right).$$

In conclusion, we note a fairly strong resemblance between the Dybvig parameterization and the approach of Hagan and Woodward. Indeed, comparing (11.51) and (11.41) we see that $a'(t)$ in the former plays pretty much the same role as $\vartheta(t)$ in the latter. The initial yield curve fit conditions, (11.48) and (11.43), are also rather similar. Hence, our words of caution with regards to the Dybvig parameterization apply here as well.

11.3.3 Monte Carlo Simulation

11.3.3.1 SDE Discretization

For the purposes of securities pricing by Monte Carlo methods, we are generally interested in advancing not only $r(t)$ through time, but also the inverse of the money market numeraire

$$\exp\left(-\int_0^t r(u)du\right) = P(0,t) \exp\left(-\int_0^t x(u)du\right) \triangleq P(0,t)Y(t), \quad (11.53)$$

where we recall that $x(t) \equiv r(t) - f(0,t)$. Our starting point can be¹³ the vector SDE

$$d \begin{pmatrix} x(t) \\ Y(t) \end{pmatrix} = \begin{pmatrix} \mu_x(t, x(t)) \\ -x(t)Y(t) \end{pmatrix} dt + \begin{pmatrix} \sigma_x(t, x(t)) \\ 0 \end{pmatrix} dW(t), \quad (11.54)$$

where the functions μ_x and σ_x were defined in (11.26) above.

In general, simulation of (11.54) requires usage of discretization methods, several of which were introduced in Chapter 3. For instance, the *Euler scheme* for (11.54) would advance the SDE for time t_i to t_{i+1} according to

$$\begin{pmatrix} \widehat{x}_{i+1} \\ \widehat{Y}_{i+1} \end{pmatrix} = \begin{pmatrix} \widehat{x}_i \\ \widehat{Y}_i \end{pmatrix} + \begin{pmatrix} \mu_x(t_i, \widehat{x}_i) \\ -\widehat{Y}_i \widehat{x}_i \end{pmatrix} \Delta_i + \begin{pmatrix} \sigma_x(t_i, \widehat{x}_i) \\ 0 \end{pmatrix} Z_i \sqrt{\Delta_i}, \quad \Delta_i = t_{i+1} - t_i,$$

where $\widehat{x}_i = \widehat{x}(t_i)$, $\widehat{Y}_i = \widehat{Y}(t_i)$, and $Z_i \sim \mathcal{N}(0, 1)$ is a sample from a standard Gaussian distribution. The Euler scheme is of (weak) convergence order one in the time step. To improve this, a second-order Milstein scheme can be constructed by Ito-Taylor expanding (11.54) to second order, using the technique in Section 3.2.6.3. The construction is tedious but straightforward (see Chapter IV in Andersen [1996] for the details), so we skip it and only show the final result

$$\begin{aligned} \widehat{x}_{i+1} &= \widehat{x}_i + \left(\mu_x(t_i, \widehat{x}_i) - \frac{1}{2} \mathcal{L}_1 \sigma_x(t_i, \widehat{x}_i) \right) \Delta_i + \sigma_x(t_i, \widehat{x}_i) Z_i \sqrt{\Delta_i} \\ &\quad + \frac{1}{2} (\mathcal{L}_1 \mu_x(t_i, \widehat{x}_i) + \mathcal{L}_0 \sigma_x(t_i, \widehat{x}_i)) Z_i \Delta_i \sqrt{\Delta_i} \\ &\quad + \frac{1}{2} \mathcal{L}_0 \mu_x(t_i, \widehat{x}_i) \Delta_i^2 + \frac{1}{2} \mathcal{L}_1 \sigma_x(t_i, \widehat{x}_i) Z_i^2 \Delta_i, \end{aligned} \quad (11.55)$$

$$\widehat{Y}_{i+1} = \widehat{Y}_i \left(1 - \widehat{x}_i \Delta_i + \frac{1}{2} (\widehat{x}_i^2 - \mu_x(t_i, \widehat{x}_i)) \Delta_i^2 - \frac{1}{2} \sigma_x(t_i, \widehat{x}_i) Z_i \Delta_i \sqrt{\Delta_i} \right), \quad (11.56)$$

where we have introduced differential operators

¹³Instead of discretizing $Y(t)$, we could also discretize $I(t) = \ln Y(t)$, as in Section 10.1.6.1. For variation we use $Y(t)$ in this section.

$$\mathcal{L}_0 = \frac{\partial}{\partial t} + \mu_x(t, x) \frac{\partial}{\partial x} + \frac{1}{2} \sigma_x(t, x)^2 \frac{\partial^2}{\partial x^2}, \quad \mathcal{L}_1 = \sigma_x(t, x) \frac{\partial}{\partial x}.$$

The Milstein scheme (11.55)–(11.56) is rather formidable-looking, and its practical efficiency tends to be quite model-dependent. Still, using an affine model as a test case Andersen [1996] shows that the Milstein scheme outperforms the Euler scheme handily, even after taking the additional computational burden of (11.55)–(11.56) into consideration. Schemes with order higher than two can be constructed along similar principles, but will, in our experience, rarely be worth the hassle. We also remind the reader that higher-order schemes can be constructed by Richardson extrapolation, as discussed in Section 3.2.7. Andersen [1996] reports modest gains for a third-order scheme constructed by Richardson extrapolation of the Milstein scheme above.

At this point, let us note that for European-style securities paying a single cash flow at time T , the ideas of Section 10.1.6.3 can be applied here, and the burden of simulating $Y(t)$ could be avoided by a change to the T -forward measure Q^T . For securities that pay intermediate cash flows, however, matters are more complicated as these flows must effectively be future-valued to time T . For instance, a random coupon c paid at time $T' < T$ will require us to compute the numeraire-deflated value $c/P(T', T)$. But here, unfortunately, the quantity $P(T', T)$ is generally not known analytically at time T' as a function of the model state variables. Of course, without affecting the economics of the trade one could invest the proceeds c into a money market account β at time T' , yielding the payout $c\beta(T)/\beta(T')$ at time T . Evaluating this payout, however, would again require us to keep track of $Y(t)$, at least on the interval $[T', T]$. This problem, however, can be avoided by using the spot measure instead of the forward measure, as outlined in Section 10.1.6.3. Much more material about numeraire simulation strategies can be found in Chapter 14.

Finally, a note on variance reduction for short rate model simulation. A systematic discussion of variance reduction techniques for short rate models can be found in Chapter IV of Andersen [1996] and Andersen and Boyle [2000]. Most of the methods discussed in these sources can be found in the survey of Section 3.4 and shall not be repeated. We do, however, highlight here the particularly useful idea of applying importance sampling based on information extracted from a tractable (e.g. Gaussian or affine) approximation to the short rate SDE. We postpone the discussion of this technique, which relies on the material in Section 3.4.4.3, to Chapter 25.

11.3.3.2 Practical Issues with Monte Carlo Methods

As was the case for finite difference methods (see discussion in Section 11.3.1), whenever an explicit bond reconstitution formula is lacking, the effort and complexity required to price derivatives by Monte Carlo methods

increase significantly. For instance, consider applying Monte Carlo methods to the pricing of short-dated expiry T call option on a long-dated maturity T^* discount bond. This price (V) is computed as a risk-neutral expectation

$$\begin{aligned} V(0) &= P(0, T) \mathbb{E} \left(Y(T)^{-1} (P(T, T^*) - K)^+ \right) \\ &= P(0, T) \mathbb{E} \left(Y(T)^{-1} \left(\mathbb{E}_T \left(e^{-\int_T^{T^*} x(u) du} - \int_T^{T^*} f(0, u) du \right) - K \right)^+ \right) \\ &= P(0, T^*) \mathbb{E} \left(Y(T)^{-1} \left(\mathbb{E}_T \left(e^{-\int_T^{T^*} x(u) du} \right) - K \frac{P(0, T)}{P(0, T^*)} \right)^+ \right). \end{aligned}$$

An immediate problem is here the fact that the inner time T expectation

$$\mathbb{E}_T \left(e^{-\int_T^{T^*} x(u) du} \right)$$

is not explicitly known as a function of $x(T)$, but must itself be computed by numerical methods. A brute-force approach involves estimating the expectation by Monte Carlo methods, launching a “simulation-within-a-simulation” at time T . The computational expense involved in such a scheme would most likely be prohibitive. Alternatives involve using a regression on a space of basis functions to estimate the function

$$Q(T, T^*, x) = \mathbb{E} \left(e^{-\int_T^{T^*} x(u) du} \mid x(T) = x \right);$$

we discuss this approach in some detail in Chapter 18.

Alternatively, we can always estimate $Q(T, T^*, x)$ by finite difference methods, as in Section 11.3.1. Combining finite difference methods and Monte Carlo methods for the purposes of pricing a European option on a discount bond makes little practical sense, of course, as we would always prefer finite difference methods for this payout. For path-dependent options, however, this idea may in fact be the best way of computing option prices. Loosely, such a scheme would use a finite difference grid to pre-compute zero-coupon bond prices $P(t_i, \cdot, x)$ on a grid $x \in \{x_j\}_{j=0}^{m+1}$, at all dates t_i , $i = 1, 2, \dots, N$, required by the path-dependent payout function considered. When paths for $x(t)$ in a subsequent step are generated by Monte Carlo simulation, interpolation of the N discount bond price vectors of dimension $m+2$ available in the finite difference grid would allow us to compute rapidly discount bond prices $Q(t_i, \cdot, x(t_i))$, $i = 1, 2, \dots$, at all relevant dates. We can use the schemes in Section 11.3.3.1 for the purpose of drawing paths of $x(t)$. If, however, we wish to make the dynamics for $x(t)$ perfectly consistent with the finite difference grid, we can use the forward induction techniques of Sections 11.3.2.1 and 11.3.2.2 to work out the (discrete) transition probabilities for $x(t)$ implied by the finite difference grid. Paths for $x(t)$ can then be generated directly from these probabilities. With this approach, we only draw values of $x(t)$ on the spatial grid $\{x_j\}_{j=0}^{m+1}$ of the finite difference grid, and therefore never have to apply interpolation methods when looking up discount bond prices.

11.A Appendix: Markov-Functional Models

The purpose of this appendix is to give a brief account of the class of *Markov-functional* (MF) models. We only consider the one-factor case. Extensions to higher dimensions are possible, but practical implementation challenges tend to increase substantially for dimensions higher than one. MF models were introduced in Kennedy et al. [2000] and while their popularity is generally waning, they are still used in some banks.

11.A.1 State Process and Numeraire Mapping

We have already observed in Section 11.3.2.6 that to define an arbitrage-free interest rate model we really only need two ingredients: a stochastic process that drives the evolution of interest rates, and a functional form that maps that process into a numeraire. The development of Markov-functional models normally starts with specializing this setup to a numeraire taken to be the discount bond $P(\cdot, T^*)$ to the final maturity of interest T^* , and a Markov stochastic process that is Gaussian in the corresponding terminal measure Q^* (see Section 4.2.4). Assuming arbitrarily that $x(t)$ is a Q^* -martingale, we write

$$dx(t) = \sigma(t)e^{\kappa t} dW^*(t), \quad x(0) = 0, \quad (11.57)$$

where W^* is a one-dimensional Brownian motion in the terminal measure, and where we for simplicity have assumed that the mean reversion κ is constant. The role of κ is to control inter-temporal correlations in the model; see Section 13.1.8.1 for the importance of this. The transition density of $x(t)$ in Q^* is, trivially,

$$\begin{aligned} p(y, t; z, s) &\triangleq Q^*(x(s) \in [z, z + dz] | x(t) = y) / dz \\ &= \frac{1}{\sqrt{2\pi}v(t, s)} \exp\left(-\frac{(z - y)^2}{2v(t, s)^2}\right), \quad s > t, \end{aligned}$$

where

$$v(t, s)^2 = \int_t^s \sigma(u)^2 e^{2\kappa u} du.$$

In the spirit of Section 11.3.2.6, we define $P(t, T^*)$ to be a deterministic function of the state variable process $x(t)$,

$$P(t, T^*) = P(t, T^*, x(t)), \quad P(t, T^*, x) = H(t, x), \quad (11.58)$$

for some exogenously given function $H : \mathbb{R}^2 \rightarrow [0, 1]$. As we recall, this is sufficient to define all discount bonds in the model since, for any $0 \leq t < T \leq T^*$,

$$P(t, T) = P(t, T^*) E_t^* \left(\frac{1}{P(T, T^*)} \right). \quad (11.59)$$

where E^* denotes expectation in measure Q^* . This allows us to express all discount bonds as functions of $x(t)$:

$$\begin{aligned} P(t, T) &= P(t, T, x(t)), \\ P(t, T, x) &= P(t, T^*, x) E^* \left(\frac{1}{P(T, T^*, x(T))} \middle| x(t) = x \right) \\ &= H(t, x) \int_{-\infty}^{\infty} \frac{p(x, t; z, T)}{H(T, z)} dz. \end{aligned} \quad (11.60)$$

The formula (11.59) can be specialized to $t = 0$, yielding

$$P(0, T) = P(0, T^*) E^* \left(\frac{1}{P(T, T^*)} \right) = P(0, T^*) E^* \left(\frac{1}{H(T, x(T))} \right), \quad (11.61)$$

which constitutes a *no-arbitrage condition* on the mapping function $H(\cdot, \cdot)$. This condition is often used to choose a particular function $H(\cdot, \cdot)$ from a given parametric family; compare this to condition (11.48) in Section 11.3.2.6.

In practice, the numeraire mapping function $H(t, x)$ in (11.58) is often specified only indirectly, through definition of functional forms for market rates, such as Libor or swap rates. The following two sections explore variations on this idea.

11.A.2 Libor MF Parameterization

Let us assume that a tenor structure

$$0 = T_0 < T_1 < \dots < T_N = T^*, \quad \tau_n = T_{n+1} - T_n,$$

is given, and define spanning forward Libor rates by

$$L_n(t) \triangleq L(t, T_n, T_{n+1}) = \left(\frac{P(t, T_n)}{P(t, T_{n+1})} - 1 \right) \tau_n^{-1}, \quad n = 0, \dots, N-1, \quad (11.62)$$

(see (4.2)). It turns out that, if we can specify the mapping of the state process $x(\cdot)$ into Libor rates on their fixing dates, $L_n(T_n)$, for all $n = 1, \dots, N-1$:

$$L_n(T_n) = l_n(x(T_n)), \quad n = 1, \dots, N-1,$$

then this is sufficient to recover the numeraire-mapping function $H(T_n, \cdot)$, $n = 1, \dots, N$, on tenor dates¹⁴ and consequently define the MF model by

¹⁴With this approach, the numeraire-mapping function is undefined for times that are not in the tenor structure. The “discrete” nature of the resulting model is one of the common criticisms of the MF approach. Pragmatically, it means that *all* dates of interest for a particular derivative security should be added to the tenor structure, or interpolation schemes not unlike those considered in Section 15.1 need to be designed.

(11.59). We show this by induction on the fixing time T_n , for $n = N-1, \dots, 1$. The starting point of the induction follows directly from (11.62) as we have

$$H(T_{N-1}, x) = P(T_{N-1}, T^*, x) = P(T_{N-1}, T_N, x) = (1 + \tau_{N-1} l_{N-1}(x))^{-1}. \quad (11.63)$$

For the induction step, let us assume that $H(T_i, x)$ are known for $i = n+1, \dots, N-1$. By (11.59) we have

$$\frac{P(T_n, T_{n+1})}{P(T_n, T^*)} = E_{T_n}^* \left(\frac{1}{P(T_{n+1}, T^*)} \right)$$

which implies that

$$\frac{1}{H(T_n, x)} = (1 + \tau_n l_n(x)) E_{T_n}^* \left(\frac{1}{H(T_{n+1}, x(T_{n+1}))} \middle| x(T_n) = x \right), \quad (11.64)$$

and the statement follows.

The consistency condition (11.61) is often used to select a particular function $l_n(\cdot)$ from a parametric family, for each n . To explain this, let us first consider what functional forms for $l_n(\cdot)$ are typically used. Suppose we desire to build a model where Libor rates on the tenor structure are close to log-normal¹⁵. Then, with $v_n = v(0, T_n)$, we fundamentally would want something like

$$L_n(T_n) \approx L_n(0) \exp \left(k_n x(T_n) - \frac{1}{2} k_n^2 v_n^2 \right), \quad (11.65)$$

where

$$k_n = \frac{e^{-\varkappa T_n} - e^{-\varkappa T_{n+1}}}{\varkappa \tau_n},$$

to hold for all n . The particular form of k_n 's, as well as the volatility parameterization (11.57), are strongly inspired by the Gaussian short rate model, see Proposition 10.1.7. (To see the connection more clearly the reader should note that the state variable $x(t)$ here is related to $x(t)$ in Proposition 10.1.7 by a multiplicative scaling of $e^{\varkappa t}$; compare (11.57) to (10.16) and disregard $y(t)$ in the latter.) Note that the quantity

$$T_n^{-1/2} k_n v_n$$

after calibration should be close to the implied Black volatility of a caplet maturing at time T_n . To preclude arbitrage, (11.65) cannot be used as is for all n (since only $L_{N-1}(t)$ is a martingale in Q^*), so we could, for instance, add a “convexity multiplier” c_n and write

$$l_n(x; c_n) = c_n L_n(0) \exp \left(k_n x - \frac{1}{2} k_n^2 v_n^2 \right).$$

¹⁵It is straightforward to extend the arguments to the case of displaced log-normal Libor rates, say. We leave this to the reader.

While v_n 's (or, equivalently, the mean reversion \varkappa and the model volatility $\sigma(t)$) can be treated as free constants to be calibrated to option prices, we would use (11.61) to set c_n 's such that the initial yield curve is replicated by the model. It is trivial to see that $c_{N-1} = 1$ and

$$H_{N-1}(T_{N-1}, x) = (1 + \tau_{N-1} l_{N-1}(x; 1))^{-1},$$

wherefore other c_n 's may be obtained as solutions to

$$\mathbb{E}_{T_n}^* \left((1 + \tau_n l_n(x(T_n); c_n)) \mathbb{E}_{T_n}^* \left(\frac{1}{H(T_{n+1}, x(T_{n+1}))} \middle| x(T_n) \right) \right) = P(0, T_n)$$

for $n = N - 2, \dots, 1$.

11.A.3 Swap MF Parameterization

Defining an MF model in terms of Libor rates is especially convenient if the model is meant to price a security that depends primarily on Libor rates (on their fixing dates), e.g. a TARN (see Section 5.15 and Chapter 20). In particular, the Libor MF parameterization allows one direct control over volatilities and other distributional characteristics of Libor rates, which makes it fairly straightforward to set up a calibration scheme that is suitable for the security (see e.g. Section 20.1.3). Libor rates, however, are not always the primary driving factors; for example, prices of Bermudan swaptions are arguably more directly linked to distributions of swap rates (see Section 19.2). Fortunately, MF models can be formulated in terms of swap rates as well, as we shall now demonstrate. It is worth noting that the relationship between Libor and swap MF models is similar to that between Libor and swap *market* models which we explore in Section 15.4.

For concreteness let us consider a set of so-called “core” swap rates,

$$S_n(t) \triangleq S_{n, N-n}(t) = \frac{P(t, T_n) - P(t, T^*)}{A_n(t)}, \quad (11.66)$$

$$A_n(t) \triangleq A_{n, N-n}(t) = \sum_{i=n}^{N-1} \tau_i P(t, T_{i+1}),$$

$n = 1, \dots, N - 1$, where we used the notations (4.8), (4.10) for $A_{k,m}$ and $S_{k,m}$. We assume that the core swap rates on their fixing dates are specified as deterministic functions $s_n(x)$ of the state process,

$$S_n(T_n) = s_n(x(T_n)), \quad n = 1, \dots, N - 1.$$

As for the Libor-based specification above, we claim that the knowledge of the functions $\{s_n(\cdot)\}_{n=1}^{N-1}$ (together with the dynamics of the state process $x(t)$) is sufficient to define the numeraire mapping $H(\cdot, \cdot)$ and, therefore, an arbitrage-free model of interest rates.

The proof also proceeds by induction. As we have that $S_{N-1}(T_{N-1}) = L_{N-1}(T_{N-1})$, the starting point of induction is given by (11.63), i.e.

$$H(T_{N-1}, x) = (1 + \tau_{N-1} s_{N-1}(x))^{-1}.$$

For the induction step $n+1 \rightarrow n$, we note from (11.66) that

$$\begin{aligned} \frac{1}{P(T_n, T^*)} &= 1 + S_n(T_n) \frac{A_n(T_n)}{P(T_n, T^*)} \\ &= 1 + S_n(T_n) \sum_{i=n}^{N-1} \tau_i \frac{P(T_n, T_{i+1})}{P(T_n, T^*)} \end{aligned}$$

and so we have (compare to (11.64))

$$\frac{1}{H(T_n, x)} = 1 + s_n(x) \sum_{i=n}^{N-1} \tau_i E_{T_n}^* \left(\frac{1}{H(T_{i+1}, x(T_{i+1}))} \middle| x(T_n) = x \right).$$

As in the Libor specification, we can choose functions $\{s_n(\cdot)\}$ to approximate log-normal (or displaced log-normal) distribution of swap rates. Also in direct analogy to the Libor case, we typically have some no-arbitrage conditions to satisfy, usually by means of setting some parameters in the specific functional form of $\{s_n(\cdot)\}$. We leave these details for the reader to explore.

11.A.4 Non-Parametric Calibration

So far we defined Markov-functional models by specific parametric mappings of the state process into Libor and swap rates. Originally, however, the class of models was introduced in a non-parametric way (see Hunt and Kennedy [2000] for a typical treatment), where mapping functions are deduced from market prices of caplets or swaptions across all strikes. While we typically prefer the parametric approach (for reasons we touch upon below), let us nevertheless quickly review the non-parametric method for completeness.

Through equation (11.60), we can turn the payout of any T -maturity security that depends on the state of the yield curve into a function of $x(T)$, $g(x(T); K, T)$ say, where K is some payout parameter (virtually always a strike). The time 0 price of this security is

$$\begin{aligned} V(0; K, T) &= P(0, T^*) E^* \left(\frac{g(x(T))}{P(T, T^*)} \right) \\ &= P(0, T^*) \int_{-\infty}^{\infty} p(0, 0; z, T) \frac{g(z; K)}{H(T, z)} dz. \end{aligned} \quad (11.67)$$

Assuming that $H(T, x)$ is invertible in x , we may write (11.67) as

$$V(0; K, T) = P(0, T^*) \int_{-\infty}^{\infty} p(0, 0; z, T) q(H(T, z); K, T) dz, \quad (11.68)$$

where

$$q(H(T, z); K, T) \triangleq \frac{g(z; K, T)}{H(T, z)}.$$

If $V(0; K, T_n)$ is known¹⁶ for a continuum of parameters (strikes) K , (11.68) defines an integral equation that may allow one to uncover the function $H(T_n, \cdot)$ (or, often more conveniently, $l_n(\cdot)$ or $s_n(\cdot)$).

Solution of (11.68) is typically done for a fixed number of M strikes, with $H(T, x)$ solved for on a grid $\{x_j\}_{j=1}^M$. In practice, this procedure is difficult to make fully robust, and the numerical solution is often prone to instabilities at long maturities, even if sophisticated special-purpose numerical techniques are employed (see Hunt and Kennedy [2000] for such techniques, many of which rely on the fact that polynomials can be integrated exactly against the Gaussian density). Even when numerically stable, a non-parametric solution for $H(\cdot, \cdot)$ may imply unrealistic evolution of the volatility smile through time, a general feature of local volatility models as explained in Section 7.1.3. To avoid these issues we may either pre-smooth the option prices used for calibration purposes (e.g. by best-fitting a CEV or a displaced log-normal model to the market smile), or, preferably in our opinion, we may use a low-dimensional parametric form for $H(t, x)$ as in Sections 11.A.2 and 11.A.3.

11.A.5 Numerical Implementation

Numerical securities valuation in an MF model is typically quite simple as the state process is Gaussian. Let us assume that the function $H(t, x)$ has been established, and consider, say, implementation of the model in a finite difference grid. Let the derivative value function be $V(t, x(t))$, and set $V^*(t, x) = V(t, x)/H(t, x)$ (such that $V(0, 0) = P(0, T^*)V^*(0, 0)$). As $V^*(t, x)$ must be a Q^* -martingale, we can write

$$\frac{\partial V^*}{\partial t} + \frac{1}{2}\sigma(t)^2 e^{2xt} \frac{\partial^2 V^*}{\partial x^2} = 0, \quad (11.69)$$

subject to appropriate terminal and intermediate jump payout conditions. In evaluating terminal and intermediate payout conditions, we would typically need to apply the numerical expression (11.60) to establish the state of the yield curve. We should note that the MF literature generally prefers to use Gaussian integration methods (rather than standard PDE solvers) to evaluate the PDE (11.69), see Hunt and Kennedy [2000] for details.

¹⁶We could use market prices for this, or we could use option prices computed from a vanilla model that we wish for the MF model to emulate.

11.A.6 Comments and Comparisons

The one-factor MF model competes with a number of models in this book, especially the quasi-Gaussian class in Chapter 13. The quasi-Gaussian model allows for arbitrary local volatility (as does the MF model), but has closed-form formulas that allow for reconstituting the term structure of discount bonds *analytically* from the underlying state variable, rather than through numerical integration (see (11.60)). In addition, the quasi-Gaussian model is substantially more “direct” in its modeling of the forward curve and has an easy-to-state term structure of instantaneous forward rate volatilities. This, in turn, makes the model more transparent in its causality structure — especially when it comes to the evolution of the volatility term structure and smile — and often makes it easy to devise good closed-form approximations for swaption and cap prices. As a consequence, calibration of quasi-Gaussian models to option and bond prices is virtually always much faster than for MF models. In addition, quasi-Gaussian models are quite straightforward to extend to high dimensions and to stochastic volatility dynamics; these extensions are far more difficult¹⁷ for MF models. On the flip side, a quasi-Gaussian model involving a single Brownian motion involves a *two-dimensional* state vector process, which makes derivatives pricing by finite difference methods slower than for MF models¹⁸. For most applications, total computation time of calibration and valuation is, however, less for the quasi-Gaussian model.

Due to its flexibility, extensibility, transparency, and ease of numerical implementation (no integration tricks are required), we generally prefer the quasi-Gaussian model over the MF model, and consequently dedicate an entire chapter to the former — and only this appendix to the latter. For those interested in learning more about MF models than what we offered here, Kennedy et al. [2000] and Hunt and Kennedy [2000] are good starting points.

¹⁷Indeed, we are unaware of the existence of any published MF models with stochastic volatility.

¹⁸As we shall see in Chapter 13, one component in the state vector is locally deterministic (i.e. it involves no Brownian motion term), so in a sense the quasi-Gaussian model has a state process dimension of “one-and-a-half”, allowing for significant speed-ups in the numerical implementation.

Multi-Factor Short Rate Models

Short rate models with only a single driving Brownian motion imply that the instantaneous correlation between forward rates at different maturities is one, a prediction that is demonstrably contrary to reality, as we show in Chapter 14. While many standard securities are, as it turns out, only weakly affected by correlations across the term structure of forward rates, this may not be the case for exotic securities, especially the ones that depend in a non-linear way on the spread between rates of different maturities. Indeed, as a general rule all derivatives that have payouts¹ exhibiting significant convexity to non-parallel moves of the forward curve must *not* be priced in a one-factor model.

In this chapter, we proceed to extend the material from Chapters 10 and 11 to cover the case of multiple driving Brownian motions. This will allow us to properly deal with securities that depend on non-parallel forward curve moves, and will also entail more subtle benefits, including the ability to model non-monotonic volatility term structures in fully time-stationary fashion. The role of the traditional multi-factor short rate models in modern derivatives pricing is, even more so than for the one-factor models, increasingly limited, as more sophisticated multi-factor frameworks have emerged over the last decade. We shall have ample opportunity to address these developments in future chapters, but a brief treatment here of the multi-factor short rate model class is still worthwhile.

As multi-factor short rate models are typically substantially more demanding to handle numerically than are one-factor models, analytical tractability is key to making multi-factor models operational. For instance, a completely generic SDE specification of a multi-factor model (along the

¹The judgment of whether a security is convex in forward rate twists and tilts can often be quite difficult. Some securities that one might guess should be sensitive to forward rate correlation in fact only display material sensitivity to forward rate *auto*-correlation. Bermudan swaptions are a good example; see Chapter 19 for more details.

lines of Section 11.3) will require significant computational effort to calibrate to market yields and volatilities, rarely leading to a usable result. As a consequence, we here elect to stay entirely in the realm of models that will allow discount bonds to be priced in closed form from the state variables of the model.

This chapter is broken into three parts. The first part develops the multi-factor Gaussian model in considerable generality, in the process demonstrating a number of features and techniques that apply to all short rate models. The second, much shorter, part provides a brief description of the multi-factor affine class, and the third part considers a particular class of quadratic-affine models that are well-suited for practical applications. For a fuller treatment of the multi-factor affine and affine-quadratic models, we refer the reader to Duffie et al. [2000], Duffie and Kan [1996], Duffie et al. [2003], Leippold and Wu [2002] and Ahn et al. [2002].

12.1 The Gaussian Model

As was the case for the one-factor Gaussian model, the multi-factor Gaussian model can be developed in two different ways: the “classical” way (from the bottom up) and the “modern” way (from a separability condition). As either technique leads to useful insights, we here show both.

12.1.1 Development from Separability Condition

A general d -factor Gaussian model can be written as

$$dP(t, T)/P(t, T) = r(t) dt - \sigma_P(t, T)^\top dW(t),$$

where $\sigma_P(t, T)$ is a bounded d -dimensional function of time, and $W(t)$ a d -dimensional Brownian motion in the risk-neutral measure Q . Written in terms of instantaneous forward rates, we get, from the HJM results in Chapter 4,

$$\begin{aligned} df(t, T) &= \sigma_f(t, T)^\top \sigma_P(t, T) dt + \sigma_f(t, T)^\top dW(t) \\ &= \sigma_f(t, T)^\top \int_t^T \sigma_f(t, u) du dt + \sigma_f(t, T)^\top dW(t). \end{aligned} \quad (12.1)$$

This model is generally not Markovian, unless we impose additional restrictions. A relevant result is the following.

Proposition 12.1.1. *Assume that $\sigma_f(t, T)$ is “separable”, in the sense that it can be written as*

$$\sigma_f(t, T) = g(t)h(T), \quad (12.2)$$

where g is a $d \times d$ deterministic matrix-valued function, and h is a d -dimensional deterministic vector. Then

$$f(t, T) = f(0, T) + \Omega(t, T) + h(T)^\top z(t),$$

where $\Omega(t, T)$ is a deterministic scalar given in (12.6) and $z(t)$ is a d -dimensional random vector satisfying

$$dz(t) = g(t)^\top dW(t), \quad z(0) = 0. \quad (12.3)$$

In particular, we have

$$r(t) = f(0, t) + \Omega(t, t) + h(t)^\top z(t). \quad (12.4)$$

Proof. Inserting (12.2) into (12.1) and integrating over time, we get

$$f(t, T) = f(0, T) + \Omega(t, T) + h(T)^\top z(t), \quad (12.5)$$

where $z(t) = \int_0^t g(u)^\top dW(u)$ and

$$\Omega(t, T) = h(T)^\top \int_0^t g(s)^\top g(s) \int_s^T h(u) du ds. \quad (12.6)$$

□

Notice that the discount bond price volatility for the model in Proposition 12.1.1 becomes simply

$$\sigma_P(t, T) = g(t) \int_t^T h(u) du.$$

12.1.1.1 Mean-Reverting State Variables

Proposition 12.1.1 demonstrates that if (12.2) is satisfied, then the forward curve can be reconstructed from d Gaussian martingale variables $z_i(t)$, $i = 1, \dots, d$, with joint SDE (12.3). The choice of d state variables is, however, not unique, and may in fact have disadvantages in a numerical implementation since often the components of $g(t)$ grow exponentially with time. As a result, it is common to shift variables to explicitly have a mean-reverting drift. To demonstrate one particular construction, set

$$H(t) = \text{diag}(h(t)) = \begin{pmatrix} h_1(t) & 0 & \cdots & 0 \\ 0 & h_2(t) & \cdots & \cdots \\ \ddots & \ddots & \ddots & 0 \\ 0 & \ddots & 0 & h_d(t) \end{pmatrix}. \quad (12.7)$$

Assuming that for all t we have $h_i(t) \neq 0$, $i = 1, \dots, d$, then $H(t)$ is invertible, and we can define a diagonal $d \times d$ matrix $\kappa(t)$ by

$$\varkappa(t) = -\frac{dH(t)}{dt} H(t)^{-1}. \quad (12.8)$$

Let us also set

$$x(t) = H(t) \int_0^t g(s)^\top g(s) \int_s^t h(u) du ds + H(t)z(t), \quad (12.9)$$

$$y(t) = H(t) \left(\int_0^t g(s)^\top g(s) ds \right) H(t). \quad (12.10)$$

Notice that $x(t)$ is a d -dimensional random vector, and $y(t)$ is a deterministic $d \times d$ symmetric matrix. It is easily verified that $y(t)$ solves the ODE

$$dy(t)/dt = H(t)g(t)^\top g(t)H(t) - \varkappa(t)y(t) - y(t)\varkappa(t).$$

Proposition 12.1.2. *Let the forward rate volatility be separable, as in Proposition 12.1.1. Let $\varkappa(t)$, $x(t)$ and $y(t)$ be defined as in (12.8)–(12.10), and assume that $H(t) = \text{diag}(h(t))$ is invertible. Also define $\mathbf{1} = (1, 1, \dots, 1)^\top \in \mathbb{R}^d$. Then*

$$dx(t) = (y(t)\mathbf{1} - \varkappa(t)x(t)) dt + \sigma_x(t)^\top dW(t), \quad \sigma_x(t) = g(t)H(t),$$

and, with $M(t, T) \triangleq H(T)H(t)^{-1}\mathbf{1}$,

$$f(t, T) = f(0, T) + M(t, T)^\top \left(x(t) + y(t) \int_t^T M(t, u) du \right). \quad (12.11)$$

In particular, we have

$$r(t) = f(t, t) = f(0, t) + \mathbf{1}^\top x(t) = f(0, t) + \sum_{i=1}^d x_i(t).$$

Proof. Applying the Leibniz integration rule to the definition of $x(t)$ yields

$$\begin{aligned} dx(t) &= \left[\frac{dH(t)}{dt} \int_0^t g(s)^\top g(s) \int_s^t h(u) du ds \right] dt \\ &\quad + \left[H(t) \int_0^t g(s)^\top g(s) h(t) ds \right] dt + \frac{dH(t)}{dt} z(t) dt + H(t) dz(t) \\ &= \frac{dH(t)}{dt} H(t)^{-1} x(t) dt + \left[H(t) \left(\int_0^t g(s)^\top g(s) ds \right) H(t)\mathbf{1} \right] dt \\ &\quad + H(t) g(t)^\top dW(t) \\ &= (y(t)\mathbf{1} - \varkappa(t)x(t)) dt + \sigma_x(t)^\top dW(t). \end{aligned}$$

Using the forward curve reconstitution formula in Proposition 12.1.1, we get

$$\begin{aligned}
f(t, T) &= f(0, T) + h(T)^\top \int_0^t g(s)^\top g(s) \int_s^T h(u) du ds + h(T)^\top z(t) \\
&= f(0, T) + \mathbf{1}^\top H(T) \int_0^t g(s)^\top g(s) \int_s^T h(u) du ds + \mathbf{1}^\top H(T) z(t) \\
&= f(0, T) + \mathbf{1}^\top \left(H(T) \int_0^t g(s)^\top g(s) \int_s^T h(u) du ds \right) \\
&\quad + \mathbf{1}^\top \left(H(T) H(t)^{-1} x(t) - H(T) \int_0^t g(s)^\top g(s) \int_s^t h(u) du ds \right) \\
&= f(0, T) + \mathbf{1}^\top H(T) H(t)^{-1} x(t) \\
&\quad + \mathbf{1}^\top H(T) H(t)^{-1} y(t) H(t)^{-1} \int_t^T h(u) du.
\end{aligned}$$

The result (12.11) follows from the definition of $M(t, T)$, the symmetry of $H(t)$, and the fact that $H(t)^{-1} h(u) = M(t, u)$. \square

If $H(t)$ fails to be invertible, it must be because one or more of the elements in $h(t)$ are (locally) equal to zero. From Proposition 12.1.1, if this is the case it follows that some of the z_i 's must be locally redundant, in turn demonstrating that the model is not truly d -dimensional for all t . As this strongly hints at a mis-specification, the invertibility condition in the proposition above is not a strong one.

In Propositions 12.1.1 and 12.1.2, reconstitution of the discount curve from the Markov state variables is done through the instantaneous forward curve. Obviously, we can also proceed to write explicit expressions for discount bond prices. For instance, using Proposition 12.1.2 we get:

Corollary 12.1.3. *In the setting of Proposition 12.1.2, define*

$$G(t, T) = \int_t^T M(t, u) du.$$

Then

$$P(t, T) = \frac{P(0, T)}{P(0, t)} \exp \left(-G(t, T)^\top x(t) - \frac{1}{2} G(t, T)^\top y(t) G(t, T) \right).$$

Proof. From the expression (12.11) for $f(t, T)$, we get

$$\begin{aligned}
P(t, T) &= \exp \left(- \int_t^T f(t, u) du \right) \\
&= \exp \left(- \int_t^T f(0, u) du - \left(\int_t^T M(t, u)^\top du \right) x(t) \right) \\
&\quad \times \exp \left(- \int_t^T M(t, u)^\top y(t) \int_t^u M(t, s) ds du \right),
\end{aligned}$$

so that

$$\begin{aligned} P(t, T) &= \frac{P(0, T)}{P(0, t)} \\ &\times \exp \left(-G(t, T)^T x(t) - \int_t^T M(t, u)^T y(t) \int_t^u M(t, s) ds du \right). \end{aligned}$$

But here

$$\int_t^T M(t, u)^T y(t) \int_t^u M(t, s) ds du = \int_t^T \frac{\partial G(t, u)}{\partial u}^T y(t) G(t, u) du.$$

As $y(t)$ is symmetric, standard matrix calculus shows that

$$\begin{aligned} \frac{\partial}{\partial u} (G(t, u)^T y(t) G(t, u)) &= \frac{\partial G(t, u)}{\partial u}^T y(t) G(t, u) + G(t, u)^T y(t) \frac{\partial G(t, u)}{\partial u} \\ &= 2 \frac{\partial G(t, u)}{\partial u}^T y(t) G(t, u), \end{aligned}$$

such that, finally,

$$\begin{aligned} \int_t^T \frac{\partial G(t, u)}{\partial u}^T y(t) G(t, u) du &= \frac{1}{2} \int_t^T \frac{\partial}{\partial u} (G(t, u)^T y(t) G(t, u)) du \\ &= \frac{1}{2} G(t, T)^T y(t) G(t, T). \end{aligned}$$

□

Let us examine some of the matrices involved in the multi-dimensional Gaussian model. As $\varkappa(t)$ is diagonal, we must have

$$\varkappa(t) = \text{diag} \left((\varkappa_1(t), \varkappa_2(t), \dots, \varkappa_d(t))^T \right),$$

in which case (12.8) implies that

$$h(t) = \left(e^{- \int_0^t \varkappa_1(s) ds}, e^{- \int_0^t \varkappa_2(s) ds}, \dots, e^{- \int_0^t \varkappa_d(s) ds} \right)^T. \quad (12.12)$$

Each element in the forward volatility vector $\sigma_f(t, T)$ is a time-weighted average of these d exponentiated integrals. Also, we note that

$$\begin{aligned} M(t, T) &= H(T) H(t)^{-1} \mathbf{1} \\ &= \left(e^{- \int_t^T \varkappa_1(s) ds}, e^{- \int_t^T \varkappa_2(s) ds}, \dots, e^{- \int_t^T \varkappa_d(s) ds} \right)^T, \end{aligned}$$

and

$$\begin{aligned} y(t) &= \int_0^t H(t)H(s)^{-1}\sigma_x(s)^\top\sigma_x(s)H(s)^{-1}H(t) ds \\ &= \int_0^t \text{diag}(M(s,t))\sigma_x(s)^\top\sigma_x(s)\text{diag}(M(s,t)) ds. \end{aligned}$$

As all quantities in the dynamics for $x(t)$ and in the reconstitution formula for $f(t, T)$ and $P(t, T)$ evidently can be computed from knowledge of the d deterministic mean reverersions $\kappa_1(t), \kappa_2(t), \dots, \kappa_d(t)$ and the $d \times d$ volatility matrix $\sigma_x(t)$, it follows that specification of $\kappa(t)$ and $\sigma_x(t)$ fully determines our d -dimensional Gaussian model.

A brief comment about computing the bond reconstitution formula in Corollary 12.1.3 is in order. The vector $G(t, T)$ takes the form

$$G(t, T) = \left(\int_t^T e^{-\int_t^u \kappa_1(s) ds} du, \dots, \int_t^T e^{-\int_t^u \kappa_d(s) ds} du \right)^\top,$$

where the individual components can, importantly, be rewritten as

$$\begin{aligned} \int_t^T e^{-\int_t^u \kappa_i(s) ds} du &= \left(\int_0^T e^{-\int_0^u \kappa_i(s) ds} du - \int_0^t e^{-\int_0^u \kappa_i(s) ds} du \right) e^{\int_0^t \kappa_i(s) ds} \\ &\triangleq (\Lambda_i(T) - \Lambda_i(t)) e^{\int_0^t \kappa_i(s) ds}. \end{aligned} \quad (12.13)$$

In implementations, we would typically pre-cache the $2d$ scalar functions

$$\Lambda_1(t), \Lambda_2(t), \dots, \Lambda_d(t), \exp\left(\int_0^t \kappa_1(s) ds\right), \dots, \exp\left(\int_0^t \kappa_d(s) ds\right)$$

on a suitable time grid, allowing subsequent discount bond pricing to be done quickly and conveniently for arbitrary t and T .

Remark 12.1.4. The risk-neutral process for the discount bond $P(t, T)$ is log-normal

$$dP(t, T)/P(t, T) = r(t) dt - G(t, T)^\top \sigma_x(t)^\top dW(t).$$

12.1.1.2 Further Changes of Variables

Going back to Proposition 12.1.2, we note that its form is rather convenient as it writes the short rate $r(t)$ as its forward value $f(0, t)$ plus a straight sum of d Gaussian mean-reverting variables, with each variable having a drift depending only on itself (since $\kappa(t)$ is diagonal). This representation is, however, just one of many. If we allow the expression for $r(t)$ to be somewhat more complicated, then we are, for instance, free to use any mean reversion matrix — diagonal or not — that we would like. Before stating

this result, we need a little extra notation. Specifically, let us consider the generic homogeneous ODE system

$$\frac{dp(t)}{dt} = -Q(t)p(t),$$

where $Q(t)$ is a deterministic $d \times d$ matrix and p a d -dimensional (column) vector. It is well-known that the solution to this equation can always be represented as

$$p(T) = J_Q(T)p(0), \quad (12.14)$$

where $J_Q(T)$ is a $d \times d$ deterministic matrix satisfying

$$\frac{dJ_Q(t)}{dt} = -Q(t)J_Q(t). \quad (12.15)$$

The matrix $J_Q(t)$ is computable by classical ODE methods² and satisfies the boundary condition $J_Q(0) = I$, where I is the identity matrix. For the special case where Q is independent of time, we have

$$J_Q(t) = \exp(-Qt),$$

as one would expect³. For later use, let us notice that, in general,

$$\frac{d(J_Q(t)^{-1})}{dt} = (J_Q(t)^{-1})Q(t). \quad (12.16)$$

Lemma 12.1.5. *In the setup of Proposition 12.1.1, let us introduce some $d \times d$ mean reversion matrix $k(t)$ and assume that $J_k(t)$ (see (12.15)) exists and is invertible for all t . Then*

$$r(t) = f(0, t) + \Omega(t, t) + h(t)^\top J_k(t)^{-1}x(t),$$

where

$$dx(t) = -k(t)x(t)dt + \sigma_x(t)^\top dW(t), \quad \sigma_x(t) = g(t)J_k(t)^\top.$$

Proof. Set

$$x(t) = J_k(t)z(t),$$

such that $z(t) = J_k(t)^{-1}x(t)$ and, from (12.15),

$$dx(t) = -k(t)x(t)dt + J_k(t)g(t)^\top dW(t).$$

²Some readers may recognize $J_Q(T)$ as the *product integral* of $-Q(t)$ on $[0, T]$; see Dollard and Friedman [1979]. In the probability literature, the product integral is often referred to as the *fundamental matrix*, see Arnold [1974] or Karatzas and Shreve [1997].

³Recall that the exponential of a square matrix A is defined as $e^A = \sum_{k=0}^{\infty} A^k/k!$.

The result for $r(t)$ follows directly from Proposition 12.1.1. \square

The lemma shows that we can incorporate essentially *any* mean reversion matrix k into the basic martingale setup in Proposition 12.1.1 by proper scaling of i) the weighting of the state variables in the expression of $r(t)$; and ii) the volatility matrix g^\top . For numerical applications, the best choice of mean reversion is typically one that leaves both $k(t)$ and $\sigma_x(t)$ close to constant.

12.1.2 Classical Development

The traditional approach to specification of a multi-dimensional short rate model does not go through a separability condition, but instead involves postulating that $r(t)$ is an affine function of a set of state variables satisfying a linear system of SDEs. That is, one would write

$$r(t) = b_q(t) + c_q(t)^\top q(t), \quad (12.17)$$

where $b_q(t) \in \mathbb{R}$ and $c_q(t) \in \mathbb{R}^d$ are deterministic, and the d -dimensional vector-valued process $q(t)$ satisfies the risk-neutral SDE

$$dq(t) = k(t) (m(t) - q(t)) dt + \sigma(t) dW(t), \quad (12.18)$$

with $m(t) \in \mathbb{R}^d$ and $k(t), \sigma(t) \in \mathbb{R}^{d \times d}$ all being deterministic. Using the definition of $J_k(t)$ given above, we can solve (12.18) explicitly.

Lemma 12.1.6. *Let $q(t)$ be as given in (12.18). Then*

$$q(t) = J_k(t) \left(q(0) + \int_0^t J_k(s)^{-1} k(s) m(s) ds + \int_0^t J_k(s)^{-1} \sigma(s) dW(s) \right), \quad (12.19)$$

i.e. $q(t)$ has a d -dimensional Gaussian distribution, with mean

$$\mu_q(t) = J_k(t) \left(q(0) + \int_0^t J_k(s)^{-1} k(s) m(s) ds \right),$$

and covariance matrix

$$\Sigma_q(t) = J_k(t) \left(\int_0^t J_k(s)^{-1} \sigma(s) \sigma(s)^\top (J_k(s)^{-1})^\top ds \right) J_k(t)^\top.$$

Proof. Set

$$u(t) = J_k(t)^{-1} q(t),$$

and observe, from (12.16), that

$$du(t) = J_k(t)^{-1} k(t) m(t) dt + J_k(t)^{-1} \sigma(t) dW(t).$$

Setting $q(t) = J_k(t)u(t)$ and observing that $u(0) = q(0)$ leads to the result in the lemma. \square

Given Lemma 12.1.5 above, we would expect the class of models spanned by specification (12.17)–(12.18) to be identical to that of the separability condition in Proposition 12.1.1. For completeness, let us make this connection explicit.

Lemma 12.1.7. *Let $r(t)$ and $q(t)$ be as in (12.17) and (12.18), and define the martingale process*

$$dz(t) = \sigma_z(t)^\top dW(t), \quad z(0) = q(0), \quad \sigma_z(t) = \sigma(t)^\top (J_k(t)^{-1})^\top.$$

Then

$$r(t) = b_z(t) + c_z(t)^\top z(t),$$

where

$$\begin{aligned} b_z(t) &= b_q(t) + c_q(t)^\top J_k(t) \int_0^t J_k(s)^{-1} k(s) m(s) ds, \\ c_z(t) &= J_k(t)^\top c_q(t). \end{aligned}$$

Proof. If we set

$$z(t) = J_k(t)^{-1} q(t) - w(t), \quad w(t) \triangleq \int_0^t J_k(s)^{-1} k(s) m(s) ds,$$

then Ito's lemma shows that

$$dz(t) = J_k(t)^{-1} \sigma(t) dW(t).$$

Notice that

$$q(t) = J_k(t) (z(t) + w(t));$$

insertion of this expression into (12.17) proves the lemma. \square

We emphasize that the form of the expression for $r(t)$ in Lemma 12.1.7 is identical to Proposition 12.1.1 once we align notation:

$$b_z(t) = f(0, t) + \Omega(t, t), \quad c_z(t) = h(t), \quad g(t) = \sigma_z(t) = \sigma(t)^\top (J_k(t)^{-1})^\top.$$

Besides confirming that the classical approach is, indeed, equivalent to the approach in Section 12.1.1, we also note from Proposition 12.1.1 and Lemma 12.1.5 that we are free to change state variables to something other than $q(t)$ or $z(t)$.

12.1.2.1 Diagonalization of Mean Reversion Matrix

While we are looking at the traditional approach to multi-factor Gaussian models, let us for later use consider a standard question about this model class: if the model for $r(t)$ is time-homogeneous and the mean reversion matrix $k(t) = k$ is non-diagonal, can we transform the state variables in

such a way that the model remains time-homogeneous but has a *diagonal* mean reversion matrix? We know from Lemma 12.1.5 that such a change of variables is always possible if we can accept that the resulting model is not time-homogeneous. To retain time-homogeneity, however, we need to impose some regularity on k , as we show below.

Proposition 12.1.8. *Consider the model (12.17) and (12.18) with parameters c_q, k, σ being independent of time. Assume that k is diagonalizable,*

$$k = LKL^{-1},$$

where K is a $d \times d$ diagonal matrix. Set $c_Q = L^\top c_q$ and $\sigma_Q = L^{-1}\sigma$. Then

$$r(t) = b_q(t) + c_Q^\top Q(t),$$

where

$$dQ(t) = K(L^{-1}m(t) - Q(t)) dt + \sigma_Q dW(t).$$

Proof. Follows immediately from the variable transformation $Q(t) = L^{-1}q(t)$. \square

In Proposition 12.1.8, we emphasize that the new mean reversion matrix K , as well as the volatility σ_Q and the scaling vector c_Q all are independent of time. We also remind the reader that a sufficient condition for k to be diagonalizable is that k has d distinct real eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_d$; in this case $K = \text{diag}((\lambda_1, \lambda_2, \dots, \lambda_d)^\top)$. See also Section 3.1.3.

A closely related question is as follows: if the model (12.17) for $r(t)$ has a constant, non-diagonal mean reversion matrix k (but is otherwise not necessarily time-homogeneous), under which circumstances can we write $r(t) = f(0, t) + \mathbf{1}^\top x(t)$ where $x(t)$ has a constant *diagonal* mean reversion matrix \varkappa ? From Proposition 12.1.2, we know that this re-write is generally possible if we allow \varkappa to depend on t . For \varkappa to additionally be constant, the following result suffices.

Proposition 12.1.9. *Consider the model (12.17) and (12.18) with k independent of time and diagonalizable, i.e.*

$$k = LKL^{-1},$$

where K is a constant $d \times d$ diagonal matrix. Assume that $B(t) = \text{diag}(e^{-Kt}L^\top c_q)$ is invertible. Then

$$r(t) = f(0, t) + \mathbf{1}^\top x(t),$$

where

$$dx(t) = (y(t)\mathbf{1} - Kx(t)) dt + \sigma_x(t)^\top dW(t),$$

with

$$\sigma_x(t) = \sigma(t)^\top (L^{-1})^\top \text{diag}(L^\top c_q),$$

$$y(t) = \int_0^t e^{-K(t-s)} \sigma_x(s)^\top \sigma_x(s) e^{-K(t-s)} ds.$$

Proof. Using the same steps as in Proposition 12.1.8 above, we know that

$$\begin{aligned} r(t) &= b_q(t) + c_Q^\top Q(t), \\ dQ(t) &= K(L^{-1}m(t) - Q(t)) dt + L^{-1}\sigma(t) dW(t). \end{aligned}$$

An application of Ito's lemma to $e^{Kt}Q(t)$ reveals that, in the notation of Section 12.1.1, this model is characterized by

$$h(t) = e^{-K^\top t} c_Q = e^{-Kt} L^\top c_q, \quad g(t) = \sigma(t)^\top (L^{-1})^\top e^{Kt}.$$

As K is diagonal, $\varkappa(t)$ in (12.8) becomes

$$\varkappa(t) = \text{diag}(Ke^{-Kt}L^\top c_q) \text{diag}(e^{-Kt}L^\top c_q)^{-1} = K,$$

and

$$g(t)H(t) = \sigma(t)^\top (L^{-1})^\top e^{Kt} \text{diag}(e^{-Kt}L^\top c_q) = \sigma(t)^\top (L^{-1})^\top \text{diag}(L^\top c_q).$$

The result follows from Proposition 12.1.2. \square

12.1.3 Correlation Structure

As discussed earlier, one important motivation for the introduction of a multi-factor interest rate model is the ability to control correlations among various points on the forward curve. Let $\rho(t, T_1, T_2)$ denote the time t instantaneous correlation between the forward rates $f(t, T_1)$ and $f(t, T_2)$. From the representation in Proposition 12.1.1, we get the following result, the proof of which is straightforward.

Lemma 12.1.10. *Let the model for $r(t)$ be as in Proposition 12.1.1. Then*

$$\rho(t, T_1, T_2) = \frac{h(T_1)^\top g(t)^\top g(t) h(T_2)}{\sqrt{h(T_1)^\top g(t)^\top g(t) h(T_1)} \sqrt{h(T_2)^\top g(t)^\top g(t) h(T_2)}}.$$

In a practical model, we generally would strongly prefer for this correlation structure to be *time-stationary*, in the sense that ρ does not depend outright on t , but only on time to maturity $T_1 - t$ and $T_2 - t$, i.e.

$$\rho(t, T_1, T_2) = \rho(T_1 - t, T_2 - t).$$

This restriction, if enforced, imposes a number of constraints on the model parameters; we return to this topic in Section 12.1.4.2 below.

While on the topic of correlation, let us remark that multi-factor Gaussian models are sometimes specified with the use of correlated Brownian motions. This setup, of course, can be translated to our setting quite easily. Specifically suppose that we start with a setup similar to that of (12.17) and (12.18), but now write

$$dq(t) = k(t)(m(t) - q(t)) dt + \sigma(t) dW^*(t),$$

where W^* is a d -dimensional vector of correlated Brownian motions. Let $R(t)$ be the relevant correlation matrix of increments to $W^*(t)$ and let

$$R(t) = C(t)C(t)^\top$$

for a square root matrix $C(t)$. Then, from results in Section 3.1.2.1, we may write $dW^*(t) = C(t)dW(t)$ for a standard (uncorrelated) vector-valued Brownian motion $W(t)$, and thereby

$$dq(t) = k(t)(m(t) - q(t)) dt + \sigma(t) C(t) dW(t).$$

It follows, not surprisingly, that we can incorporate correlation in Brownian increments by a simple rotation of the volatility matrix (by $C(t)$).

12.1.4 The Two-Factor Gaussian Model

Having now outlined the general theory, let us make matters more concrete (and more practical) by focusing on the important case of $d = 2$.

12.1.4.1 Some Basics

In practical applications, a reasonable way to specify a two-dimensional model would be to set, in the notation of Proposition 12.1.1,

$$h(t) = \begin{pmatrix} e^{-\int_0^t \kappa_1(u)du} \\ e^{-\int_0^t \kappa_2(u)du} \end{pmatrix}, \quad g(t) = \begin{pmatrix} \sigma_{11}(t)e^{\int_0^t \kappa_1(u)du} & \sigma_{12}(t)e^{\int_0^t \kappa_2(u)du} \\ \sigma_{21}(t)e^{\int_0^t \kappa_1(u)du} & \sigma_{22}(t)e^{\int_0^t \kappa_2(u)du} \end{pmatrix}.$$

We may, without loss of generality, assume that $g(t)$ is lower diagonal⁴, i.e. we can set $\sigma_{12}(t) = 0$ always. In this case, from Proposition 12.1.2 we have

$$r(t) = f(0, t) + x_1(t) + x_2(t),$$

where $x(t) = (x_1(t), x_2(t))^\top$ satisfies (with $x(0) = 0$)

$$dx(t) = (y(t)\mathbf{1} - \kappa(t)x(t)) dt + \sigma_x(t)^\top dW(t), \quad \sigma_x(t) = \begin{pmatrix} \sigma_{11}(t) & 0 \\ \sigma_{21}(t) & \sigma_{22}(t) \end{pmatrix}, \quad (12.20)$$

and where $\kappa(t) = \text{diag}((\kappa_1(t), \kappa_2(t))^\top)$ and $y(t)$ is a deterministic 2×2 matrix.

We notice that the instantaneous correlation between $x_1(t)$ and $x_2(t)$ is

$$\rho_x(t) = \frac{\sigma_{22}(t)\sigma_{21}(t)}{\sqrt{\sigma_{11}(t)^2 + \sigma_{21}(t)^2}\sqrt{\sigma_{22}(t)^2}},$$

⁴If $g(t)$ is not lower diagonal, we can always change variables (via a Cholesky decomposition, say) to make it so.

so for convenience we may, as in Section 12.1.3, rewrite our model as

$$dx(t) = (y(t)\mathbf{1} - \kappa(t)x(t))dt + \sigma_x^*(t)dW^*(t), \quad (12.21)$$

where $\langle dW_1^*(t), dW_2^*(t) \rangle = \rho_x(t)dt$ and $\sigma_x^*(t)$ is diagonal with non-negative elements,

$$\sigma_x^*(t) = \begin{pmatrix} \sqrt{\sigma_{11}(t)^2 + \sigma_{21}(t)^2} & 0 \\ 0 & |\sigma_{22}(t)| \end{pmatrix} \triangleq \text{diag}\left((\sigma_1(t), \sigma_2(t))^T\right).$$

The term $y(t)$ in (12.21) can be computed by numerical integration from the results in Proposition 12.1.2.

From Corollary 12.1.3, the reconstitution formula for the yield curve is

$$f(t, T) = f(0, T) + M(t, T)^T(x(t) + y(t)G(t, T)), \quad (12.22)$$

$$P(t, T) = \frac{P(0, T)}{P(0, t)} \exp\left(-G(t, T)^T x(t) - \frac{1}{2} G(t, T)^T y(t) G(t, T)\right),$$

where

$$G(t, T) = \int_t^T M(t, u) du, \quad M(t, T) = \left(e^{-\int_t^T \kappa_1(u) du}, e^{-\int_t^T \kappa_2(u) du}\right)^T.$$

The specification (12.21)–(12.22) is, we feel, the most intuitive representation of the two-factor Gaussian short rate model. For a complete model specification, we evidently must specify 5 functions of time: $\rho_x(t)$, $\kappa_1(t)$, $\kappa_2(t)$, $\sigma_1(t)$, and $\sigma_2(t)$. As discussed earlier, however, we may want to ensure that the model has a time-stationary correlation structure, in the sense defined in Section 12.1.3. We turn to this in Section 12.1.4.2 below.

12.1.4.2 Variance and Correlation Structure

According to Proposition 12.1.1, it follows that the model (12.21) has the forward rate process

$$df(t, T) = O(dt) + \begin{pmatrix} \sigma_1(t)e^{-\int_t^T \kappa_1(u) du} \\ \sigma_2(t)e^{-\int_t^T \kappa_2(u) du} \end{pmatrix}^T dW^*(t), \quad (12.23)$$

where we recall that $\langle dW_1^*(t), dW_2^*(t) \rangle = \rho_x(t)dt$. From this representation, or from the results in Section 12.1.3, we get the following lemma.

Lemma 12.1.11. *For the model (12.21), let*

$$\begin{aligned} b(t, T_1, T_2) = 1 + \rho_x(t) \frac{\sigma_2(t)}{\sigma_1(t)} &\left(e^{-\int_t^{T_1} (\kappa_2(u) - \kappa_1(u))du} + e^{-\int_t^{T_2} (\kappa_2(u) - \kappa_1(u))du} \right) \\ &+ \left(\frac{\sigma_2(t)}{\sigma_1(t)}\right)^2 e^{-\int_t^{T_1} (\kappa_2(u) - \kappa_1(u))du - \int_t^{T_2} (\kappa_2(u) - \kappa_1(u))du}. \end{aligned}$$

Then

$$\text{Var}_t(df(t, T)) = \sigma_1(t)^2 e^{-2 \int_t^T \kappa_1(u) du} b(t, T, T),$$

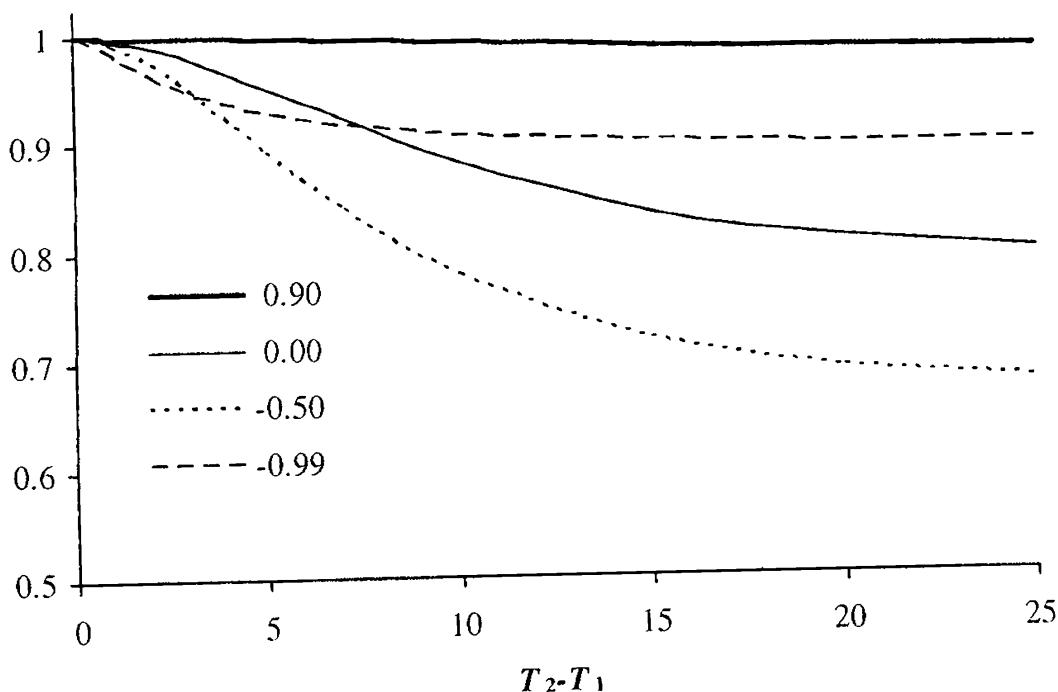
$$\rho(t, T_1, T_2) = \text{Corr}_t(df(t, T_1), df(t, T_2)) = \frac{b(t, T_1, T_2)}{\sqrt{b(t, T_1, T_1)b(t, T_2, T_2)}}.$$
(12.24)

In particular, $\rho(t, T_1, T_2)$ is time-stationary if $\rho_x(t)$, $\kappa_2(t) - \kappa_1(t)$, and $\sigma_2(t)/\sigma_1(t)$ are all constant.

We emphasize that if we chose to make our correlation structure time-stationary, then Lemma 12.1.11 shows that only two functions of time ($\sigma_1(t)$ and $\kappa_1(t)$, say) and three constants (ρ_x , $\kappa_2(t) - \kappa_1(t)$, and $\sigma_2(t)/\sigma_1(t)$) may be specified freely. Notice that if either i) $\rho_x = 1$; ii) $\kappa_2(t) - \kappa_1(t) = 0$; or iii) $\sigma_2(t)/\sigma_1(t) = 0$; then $\rho(t, T_1, T_2) = 1$, i.e. the model is reduced to having only one factor.

Figure 12.1 below shows some examples of the types of correlation term structures that can be generated in our two-factor Gaussian model.

Fig. 12.1. Forward Rate Correlation Term Structure



Notes: For the model (12.21), the figure graphs $\rho(0, T_1, T_2)$ from Lemma 12.1.11 against $T_2 - T_1$, using four different values of the parameter ρ_x . Other parameters were: $T_1 = 0.1$, $\kappa_1 = 0.1$, $\kappa_2 = 0.25$, $\sigma_1 = 0.025$, and $\sigma_2 = 0.02$.

Note the fact that the forward rate correlation is *not* necessarily a monotonic function of ρ_x . For parameterization purposes, it is often helpful

to consider $\rho(t, t, \infty)$, the correlation between the short rate and a (very) long-dated forward rate. From (12.24) we get, assuming time-homogeneity in the correlation structure,

$$\rho(t, t, \infty) = \frac{1 + \rho_x c_\sigma}{\sqrt{1 + 2\rho_x c_\sigma + c_\sigma^2}}, \quad c_\sigma = \sigma_2(t)/\sigma_1(t),$$

an expression that does not depend on mean reversion speeds. Given either ρ_x or c_σ , an exogenous specification of $\rho(t, t, \infty)$ allows us to back out c_σ or ρ_x , respectively.

12.1.4.3 Volatility Hump

Besides allowing us to properly model the correlation between various points of the forward curve, the use of a two-factor Gaussian model has another benefit relative to a one-factor model: the ability to produce a time-stationary, non-monotonic term structure of forward rate volatilities. We recall from Section 10.1.2.3 that this was not possible in a one-factor model, where non-constant mean reversion was required to produce a caplet volatility “hump”. To provide some details, assume that κ_1 and κ_2 are fixed at non-negative constant values, with at least one being positive. Also assume that σ_1 and σ_2 are fixed at constant positive values. From Lemma 12.1.11 we have

$$\begin{aligned} \text{Var}_t(df(t, T)) &= \sigma_1^2 e^{-2\kappa_1(T-t)} + \sigma_2^2 e^{-2\kappa_2(T-t)} + 2\rho_x \sigma_1 \sigma_2 e^{-(\kappa_1+\kappa_2)(T-t)} \\ &\triangleq g(T-t), \end{aligned}$$

where

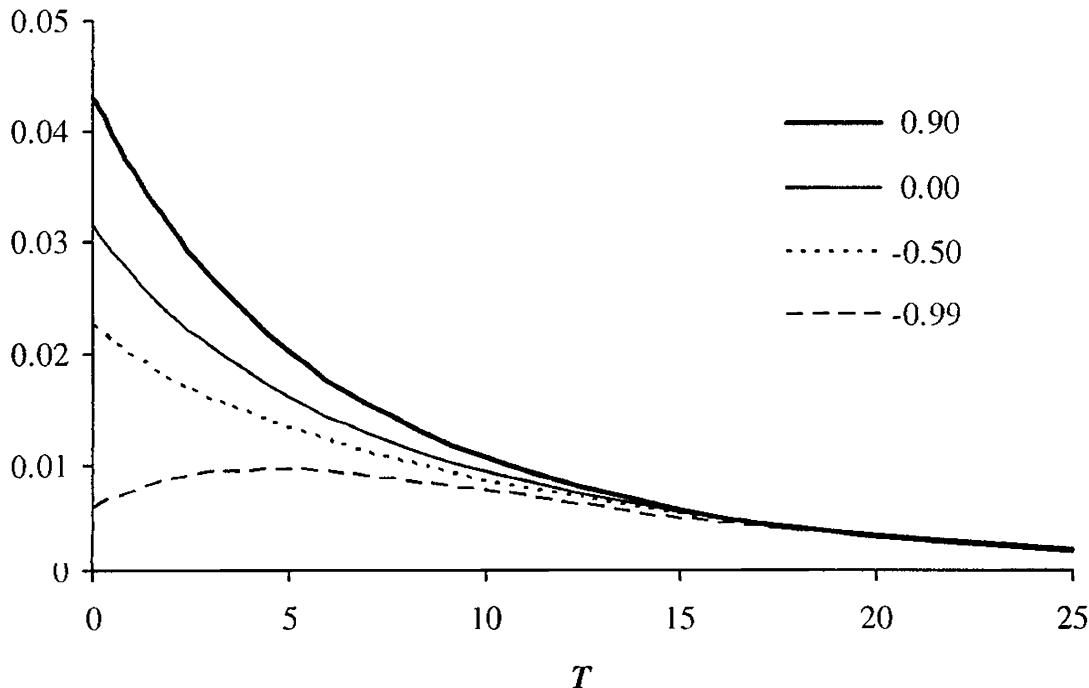
$$\frac{1}{2} \frac{dg(\tau)}{d\tau} = -\kappa_1 \sigma_1^2 e^{-2\kappa_1 \tau} - \kappa_2 \sigma_2^2 e^{-2\kappa_2 \tau} - \rho_x \sigma_1 \sigma_2 (\kappa_1 + \kappa_2) e^{-(\kappa_1+\kappa_2)\tau}.$$

For positive values of ρ_x , the forward rate variance term structure will thus always be downward-sloping ($dg(\tau)/d\tau \leq 0$). However, if we set ρ_x sufficiently negative, there may be intermediate values for $\tau = T - t$ for which the variance will increase in τ ; Figure 12.2 shows an example.

12.1.4.4 Another Formulation of the Two-Factor Model

To round off our treatment of the two-factor Gaussian model, we note that the model traditionally has been developed and presented in a manner quite different from ours. Indeed, a common starting point (e.g. Hull and White [1994b]) for the model is the doubly mean-reverting form:

$$\begin{aligned} dr(t) &= \kappa_r (\vartheta(t) + \varepsilon(t) - r(t)) dt + \sigma_r dW_r(t), \\ d\varepsilon(t) &= -\kappa_\varepsilon \varepsilon(t) dt + \sigma_\varepsilon dW_\varepsilon(t), \end{aligned} \tag{12.25}$$

Fig. 12.2. Forward Rate Volatility Term Structure

Notes: For the model (12.21), the figure graphs $\sqrt{\text{Var}(df(0, T))}$ against T , using four different values of the parameter ρ_x . Other parameters were: $\kappa_1 = 0.1$, $\kappa_2 = 0.25$, $\sigma_1 = 0.025$, and $\sigma_2 = 0.02$.

where $\vartheta(t)$ is deterministic and $\langle dW_r(t), dW_\varepsilon(t) \rangle = \rho dt$, for some constant ρ . The model can be extended to time-dependent σ_r and σ_ε , but we omit this for the sake of simplicity.

To write (12.25) in terms that are more compatible with our notation, let $q(t) = (q_1(t), q_2(t))^\top$ be defined by

$$\begin{aligned} dq(t) &= \begin{pmatrix} \kappa_r & -\kappa_r \\ 0 & \kappa_\varepsilon \end{pmatrix} \left(\begin{pmatrix} \vartheta(t) \\ 0 \end{pmatrix} - q(t) \right) dt + \begin{pmatrix} \sigma_r & 0 \\ \sigma_\varepsilon \rho & \sigma_\varepsilon \sqrt{1 - \rho^2} \end{pmatrix} dW(t) \\ &\triangleq k(m(t) - q(t)) dt + \sigma(t) dW(t), \end{aligned} \quad (12.26)$$

where the elements of $W(t) = (W_1(t), W_2(t))^\top$ are independent. Clearly, if we set

$$r(t) = q_1(t), \quad \varepsilon(t) = q_2(t), \quad (12.27)$$

we replicate the model (12.25) above. But (12.26)–(12.27) is of the form in Section 12.1.2, with $b_q(t) = 0$ and $c_q(t) = (1, 0)^\top$, so we can use Propositions 12.1.8 and 12.1.9 to re-write the model in alternative formats. Some relevant characterizations are listed below.

Lemma 12.1.12. Assume that $\kappa_r \neq \kappa_\varepsilon$. The model (12.25) can be written as

$$r(t) = Q_1(t) + Q_2(t) \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon},$$

where

$$\begin{aligned} dQ(t) &= \begin{pmatrix} \kappa_r & 0 \\ 0 & \kappa_\varepsilon \end{pmatrix} \left(\begin{pmatrix} \vartheta(t) \\ 0 \end{pmatrix} - Q(t) \right) dt \\ &\quad + \begin{pmatrix} \sigma_r - \rho \kappa_r \frac{\sigma_\varepsilon}{\kappa_r - \kappa_\varepsilon} & -\kappa_r \sigma_\varepsilon \frac{\sqrt{1-\rho^2}}{\kappa_r - \kappa_\varepsilon} \\ \rho \sigma_\varepsilon & \sigma_\varepsilon \sqrt{1-\rho^2} \end{pmatrix} dW(t). \end{aligned}$$

Proof. If $\kappa_r \neq \kappa_\varepsilon$, the matrix k can be diagonalized as $k = LKL^{-1}$, where

$$L = \begin{pmatrix} 1 & \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} \\ 0 & 1 \end{pmatrix}, \quad K = \begin{pmatrix} \kappa_r & 0 \\ 0 & \kappa_\varepsilon \end{pmatrix}, \quad L^{-1} = \begin{pmatrix} 1 & -\frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} \\ 0 & 1 \end{pmatrix}.$$

The lemma then follows from Proposition 12.1.8. \square

Lemma 12.1.13. *Assume that $\kappa_r \neq \kappa_\varepsilon$ and let K be given in Lemma 12.1.12. The model (12.25) can be written as*

$$r(t) = f(0, t) + x_1(t) + x_2(t),$$

where

$$\begin{aligned} dx(t) &= (y(t)\mathbf{1} - Kx(t))dt + \sigma_x^\top dW(t), \tag{12.28} \\ \sigma_x &= \begin{pmatrix} \sigma_r - \rho \sigma_\varepsilon \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} & \rho \sigma_\varepsilon \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} \\ -\sigma_\varepsilon \sqrt{1-\rho^2} \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} & \sigma_\varepsilon \sqrt{1-\rho^2} \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} \end{pmatrix}, \end{aligned}$$

and

$$y(t) = \int_0^t e^{-K(t-s)} \sigma_x^\top \sigma_x e^{-K(t-s)} ds.$$

The forward rate volatility is

$$\sigma_f(t, T) = \begin{pmatrix} e^{-\kappa_r(T-t)} \sigma_r + (e^{-\kappa_\varepsilon(T-t)} - e^{-\kappa_r(T-t)}) \rho \sigma_\varepsilon \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} \\ (e^{-\kappa_\varepsilon(T-t)} - e^{-\kappa_r(T-t)}) \sqrt{1-\rho^2} \sigma_\varepsilon \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} \end{pmatrix}.$$

Proof. The representation for $dx(t)$ follows directly from Proposition 12.1.9, after applying Lemma 12.1.12. The forward rate volatility can be recovered from Proposition 12.1.1. \square

In Lemma 12.1.13, we note that $y(t)$ can be written in closed form; we leave this as an exercise to the reader. We also point out that in the dynamics for $x(t)$, we may translate back to the original correlated Brownian motions W_r and W_ε by an inverse Cholesky decomposition, writing

$$d \begin{pmatrix} W_1(t) \\ W_2(t) \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ -\frac{\rho}{\sqrt{1-\rho^2}} & \frac{1}{\sqrt{1-\rho^2}} \end{pmatrix} d \begin{pmatrix} W_r(t) \\ W_\varepsilon(t) \end{pmatrix}.$$

Inserting this expression into (12.28) gives the simpler expression

$$dx(t) = (y(t)\mathbf{1} - Kx(t)) dt + \begin{pmatrix} \sigma_r & -\frac{\kappa_r - \kappa_\varepsilon}{\kappa_r - \kappa_\varepsilon} \sigma_\varepsilon \\ 0 & \frac{\kappa_r - \kappa_\varepsilon}{\kappa_r - \kappa_\varepsilon} \sigma_\varepsilon \end{pmatrix} d \begin{pmatrix} W_r(t) \\ W_\varepsilon(t) \end{pmatrix}. \quad (12.29)$$

By the same token, we can write

$$df(t, T) = O(dt) + \left(\frac{\sigma_r e^{-\kappa_r(T-t)}}{\frac{\kappa_r - \kappa_\varepsilon}{\kappa_r - \kappa_\varepsilon} \sigma_\varepsilon (e^{-\kappa_\varepsilon(T-t)} - e^{-\kappa_r(T-t)})} \right)^\top d \begin{pmatrix} W_r(t) \\ W_\varepsilon(t) \end{pmatrix}. \quad (12.30)$$

The discount bond reconstitution formula for the model (12.25) can be recovered from Corollary 12.1.3, using the representation in Lemma 12.1.13. The reconstitution formula can, of course, be re-stated in terms of the original q_i variables, should we desire to do so.

Finally, let us note that the special case of $\sigma_r = 0$ may be useful as a way to model the fact that central bank activities (which govern the dynamics of the short end of the forward curve) are often largely predictable. For the case $\sigma_r = 0$, notice that we get, from (12.29),

$$\begin{aligned} dx(t) &= O(dt) + \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} \sigma_\varepsilon \begin{pmatrix} -1 \\ 1 \end{pmatrix} dW_\varepsilon(t), \\ df(t, T) &= O(dt) + \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} \sigma_\varepsilon \left(e^{-\kappa_\varepsilon(T-t)} - e^{-\kappa_r(T-t)} \right) dW_\varepsilon(t). \end{aligned} \quad (12.31)$$

In other words, the two state variables x_1 and x_2 here become perfectly anti-correlated, and the forward curve dynamics are reduced to depending on only one Brownian motion. Despite this, notice that the model still requires two state variables (x_1 and x_2); this is a consequence of having two mean reverersions in the forward rate volatility.

12.1.5 Multi-Factor Statistical Gaussian Model

The single-Brownian-motion, two-state model discussed at the end of Section 12.1.4 highlights an interesting and important interpretation of the model parameters. Let us re-write the model slightly to make our point. We have, from (12.31), that

$$df(t, t + \tau) = O(dt) + l(\tau) dz(t),$$

where we have denoted

$$\begin{aligned} l(\tau) &= \frac{\kappa_r}{\kappa_r - \kappa_\varepsilon} (e^{-\kappa_\varepsilon \tau} - e^{-\kappa_r \tau}), \\ dz(t) &= \sigma_\varepsilon dW_\varepsilon(t). \end{aligned} \quad (12.32)$$

We can interpret $z(t)$ as the (single) factor that affects the movements of the forward rate curve $\{f(t, t + \tau)\}_{\tau \geq 0}$, and the function $l(\tau)$ as the response function, or a *loading*, whose value at time τ determines the impact of the

factor shock on a rate of tenor τ . Note that in this parameterization, the loading is a function of the tenor τ only, i.e. is time-homogeneous. This opens up the possibility of linking it to the statistically-estimated properties of the movements of the yield curve, the connection that we shall explore momentarily⁵. First, however, let us develop some technical tools.

The exponential functions $\{e^{-\kappa\tau}\}_{\kappa \in \mathbb{R}}$ are dense in the space of all continuous functions. Hence, any continuous function $l(\tau)$ can be approximated by a linear combination of exponential functions to an arbitrary degree of precision. Assume such a function (recycling the notations) $l(\tau)$ is given. Then, we can find a set of coefficients $\{v_i\}_{i=1}^n$ and exponents $\{\kappa_i\}_{i=1}^n$ such that, approximately,

$$l(\tau) \approx \sum_{i=1}^n v_i e^{-\kappa_i \tau}. \quad (12.33)$$

A moment of reflection on (12.32) and (12.33) shows that a model with the loading (12.33) could be represented as an n -state, single-Brownian-motion Gaussian model; we formalize this result as a proposition.

Proposition 12.1.14. *Let*

$$df(t, t + \tau) = O(dt) + l(\tau) dz(t), \quad (12.34)$$

$$l(\tau) = \sum_{i=1}^n v_i e^{-\kappa_i \tau}, \quad dz(t) = \sigma_1(t) dW_1(t),$$

where $W_1(t)$ is a one-dimensional Brownian motion and $\sigma_1(t)$ is a one-dimensional function of time. Then this model admits a Markovian representation in n state variables

$$r(t) = f(0, t) + \sum_{i=1}^n x_i(t),$$

where, with $x(t) = (x_1(t), \dots, x_n(t))^\top$ and $\kappa = \text{diag}((\kappa_1, \dots, \kappa_n)^\top)$, we have

$$dx(t) = (y(t)\mathbf{1} - \kappa x(t)) dt + \sigma_x(t)^\top dW(t),$$

$$\sigma_x(t) = \sigma_1(t) \begin{pmatrix} v_1 & v_2 & \cdots & v_n \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 \end{pmatrix}, \quad W(t) = (W_1(t), 0, \dots, 0)^\top,$$

and

$$y(t) = H(t) \left(\int_0^t \sigma_1(s)^2 H(s)^{-1} U H(s)^{-1} ds \right) H(t),$$

$$H(t) = \text{diag}((e^{-\kappa_1 t}, \dots, e^{-\kappa_n t})^\top),$$

$$U = \{v_k v_j\}_{k,j=1}^n,$$

⁵The material in this section is largely inspired by Balasanov [1996].

Proof. From (12.34),

$$df(t, T) = O(dt) + \sum_{i=1}^n e^{-\kappa_i T} (e^{\kappa_i t} \sigma_1(t) v_i dW_1(t)).$$

Hence the model can be written in a separable form with

$$h(t) = ((e^{-\kappa_1 t}, \dots, e^{-\kappa_n t})^\top), \quad g(t) = \sigma_x(t) H(t)^{-1},$$

where $H(t)$ and $\sigma_x(t)$ are given in the statement of the proposition. The result follows from Proposition 12.1.2, definition (12.10) for $y(t)$, and the fact that

$$\sigma_x(t)^\top \sigma_x(t) = \sigma_1(t)^2 U,$$

where U is the $n \times n$ matrix defined above. \square

The model (12.34) allows for an essentially arbitrary loading $l(\tau)$, but employs only one factor to describe the dynamics of the yield curve, a restriction that we can easily relax. Suppose we believe that m factors are needed to describe the dynamics of an interest rate curve. Also assume that we are given m loadings, each describing the (linear) response of the forward rate curve to a given factor. Approximating each loading by a linear combination of exponentials, we arrive at a model of the form

$$df(t, t + \tau) = O(dt) + \sum_{j=1}^m l_j(\tau) dz_j(t), \quad (12.35)$$

$$l_j(\tau) = \sum_{i=1}^{n_j} v_{j,i} e^{-\kappa_{j,i} \tau}, \quad dz_j(t) = \sigma_j(t) dW_j(t),$$

where W_j 's are independent Brownian motions. By a simple (but laborious, and left as an exercise to the reader) extension of Proposition 12.1.14, the model (12.35) can be shown to be Markovian in a total of

$$n = \sum_{j=1}^m n_j$$

state variables.

Proposition 12.1.15. *The model (12.35) admits a Markovian representation*

$$r(t) = f(0, t) + \mathbf{1}^\top x(t),$$

where $x(t) = (x_1(t), \dots, x_n(t))^\top$ satisfies

$$dx(t) = (y(t)\mathbf{1} - \kappa x(t)) dt + \sigma_x(t)^\top dW(t),$$

with

$$\begin{aligned}\boldsymbol{\kappa} &= \text{diag}((\kappa_{1,1}, \dots, \kappa_{1,n_1}, \kappa_{2,1}, \dots, \kappa_{2,n_2}, \dots, \kappa_{m,n_m})^\top), \\ \sigma(t) &= \text{diag}((\sigma_1(t), \dots, \sigma_1(t), \sigma_2(t), \dots, \sigma_2(t), \dots, \sigma_m(t))^\top), \\ h(t) &= (e^{-\kappa_{1,1}t}, \dots, e^{-\kappa_{1,n_1}t}, e^{-\kappa_{2,1}t}, \dots, e^{-\kappa_{2,n_2}t}, \dots, e^{-\kappa_{m,n_m}t})^\top,\end{aligned}$$

and

$$\sigma_x(t) = v\sigma(t),$$

where

$$v = \begin{pmatrix} v_{1,1} & \cdots & v_{1,n_1} & 0 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ 0 & \cdots & 0 & v_{2,1} & \cdots & v_{2,n_2} & \cdots & \vdots & \vdots & \vdots \\ \vdots & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & 0 & 0 & \cdots & v_{m,1} & \cdots & v_{m,n_m} \end{pmatrix}.$$

Here $W(t)$ is an m -dimensional vector of independent Brownian motions

$$W(t) = (W_1(t), W_2(t), \dots, W_m(t))^\top,$$

and

$$y(t) = H(t) \left(\int_0^t H(s)^{-1} \sigma(s) v^\top v \sigma(s) H(s)^{-1} ds \right) H(t), \quad H(t) = \text{diag}(h(t)).$$

The representation (12.35) allows us to link the interest rate model parameterization to statistical properties of the movements of the yield curve. To demonstrate, let us fix N_τ , the number of tenors of interest, and specify a set of tenors $\{\tau_1, \dots, \tau_{N_\tau}\}$. Then we can observe from history how the vector of rates $f(t) = (f(t, t + \tau_1), \dots, f(t, t + \tau_{N_\tau}))^\top$ changed over time. With the application of principal components (PC) analysis⁶ we can identify a set of m , $m \leq N_\tau$, factors $\zeta(t) = (\zeta_1(t), \dots, \zeta_m(t))^\top$, and m loadings $\lambda_j = (\lambda_j(\tau_1), \dots, \lambda_j(\tau_{N_\tau}))^\top$, $j = 1, \dots, m$, that we can use to represent

$$\Delta f(t) \approx \sum_{j=1}^m \lambda_j \Delta \zeta_j(t); \tag{12.36}$$

here Δ is a time-differencing operator, i.e. the day-to-day change of a given quantity. The PC analysis guarantees that the m factors will be optimal in the sense that the m factors will explain the largest possible variability of the vector of rates. As we shall see in Chapter 14, we can typically use a value m that is much smaller than N_τ , allowing for a significant reduction in dimension of the model. The loading vectors λ_j here define the shapes of the forward curve shocks associated with each factor.

⁶PC analysis was introduced in Section 3.1.3; its application to statistical analysis of interest rate curve movements will be described in details later, in Chapter 14.

Once we have identified loading vectors λ_j , $j = 1, \dots, m$, the transition from (12.36) to (12.35) is merely an interpolation exercise where the functions $l_j(\tau)$ are extracted from the vectors λ_j by tenor-interpolation and a best-fit approximation with a linear combination of exponentials. After this step, the model can be represented, and efficiently implemented, in Markovian state variables as outlined in Proposition 12.1.15. The remaining factor volatility parameters $\sigma_j(t)$, $j = 1, \dots, m$, may, for instance, be found by calibrating the model to market prices of caps/swaptions; see Section 12.1.6 for swaption pricing formulas that would be needed for such a calibration.

As we demonstrate later in Chapter 14, in the model (12.35) it is often sufficient to choose m to be 3 or 4, i.e. the yield curve movements through time are usually well explained by 3 to 4 factors. For each loading, the number of exponential terms required to match its shape is, typically, between 2 and 4 — the higher the loading number, the more complicated its shape typically is, and the more exponential terms are required. So, the overall number of state variables in the Markovian representation is typically around 10 or so.

The combination of statistical and risk neutral calibration, where some parameters (loadings) are obtained from historical data and others (factor volatilities) are market-implied, is an appealing characteristic of the model (12.35) and the parameterization strategy outlined above. Ultimately, however, (12.35), being nearly time-homogeneous, does not constitute a setup flexible enough for a precise calibration to all, or the majority of, market-quoted swaptions. In particular, historical loading shapes are often at odds with those consistent with the implied swaption volatilities. While, perhaps, not fully suitable as a model for interest rate exotics, (12.35) is still useful in settings where incorporation of historical information into pricing is of primary importance, such as risk management applications (such as VaR calculations, see Section 22.3), proprietary trading, or mortgage bonds valuation.

Remark 12.1.16. As an implementation note, we observe that working with instantaneous forward rates in the historical setting is inconvenient. Fortunately, (12.35) can be integrated in τ to obtain a similar linear representation for continuously compounded forward yields (see (4.1)), for which historical analysis can be performed more easily.

12.1.6 Swaption Pricing

While the parameterization (12.35) is just one of many for multi-factor interest rate models, it demonstrates a common strategy of specifying forward rate correlations exogenously and then calibrating the overall levels of model volatilities to European swaptions. We elaborate more on such calibration ideas in Section 12.1.7 and in Chapter 14. To make these ideas operational, we need to establish efficient pricing formulas for European swaptions. For

concreteness, we consider a payer swaption maturing at time $T_0 > 0$, with the underlying swap paying an annualized coupon c at times $T_1 < T_2 < \dots < T_N$. The swaption payout at time T_0 is thereby

$$V_{\text{swaption}}(T_0) = \left(1 - P(T_0, T_N) - c \sum_{i=0}^{N-1} \tau_i P(T_0, T_{i+1}) \right)^+, \quad \tau_i = T_{i+1} - T_i. \quad (12.37)$$

12.1.6.1 Two-Dimensional Jamshidian Decomposition

Let us consider to what extent we can use the Jamshidian decomposition in Section 10.1.3.1 in the multi-dimension case. For simplicity, we focus on the two-factor case, and later indicate how to extend our arguments to higher dimensions. Throughout we work with the parameterization⁷ in Section 12.1.1.1, i.e. $r(t) = f(0, t) + x_1(t) + x_2(t)$, where x_1 and x_2 are state variables satisfying⁸

$$dx_1(t) = (\vartheta_1(t) - \varkappa_1(t)x_1(t)) dt + \sigma_{11}(t) dW_1(t) + \sigma_{12}(t) dW_2(t), \quad (12.38)$$

$$dx_2(t) = (\vartheta_2(t) - \varkappa_2(t)x_2(t)) dt + \sigma_{21}(t) dW_1(t) + \sigma_{22}(t) dW_2(t). \quad (12.39)$$

Expressions for $\vartheta_1(t) = y_{11}(t) + y_{12}(t)$ and $\vartheta_2(t) = y_{21}(t) + y_{22}(t)$ can be found in Section 12.1.1.1. We recall, in particular, the reconstitution formula

$$P(T, T + \Delta) = \frac{P(0, T + \Delta)}{P(0, T)} e^{A(T, T + \Delta) - x_1(T)G_1(T, T + \Delta) - x_2(T)G_2(T, T + \Delta)}, \quad (12.40)$$

where A, G_1, G_2 are known deterministic functions given in Corollary 12.1.3.

We shall first need to establish the following result.

Lemma 12.1.17. *Consider a put option on a discount bond, i.e. a derivative with T_0 payout*

$$p_i(T_0; K) = (K - P(T_0, T_i))^+, \quad T_i > T_0.$$

Let E^{T_0} denote expectation in the T_0 -forward measure Q^{T_0} , and define the x_2 -conditional option price

$$p_i(0; K, x_2) = P(0, T_0) E^{T_0} (p_i(T_0; K) | x_2(T_0) = x_2).$$

Then

⁷This choice is made largely for reasons of familiarity. We indicate later (at the very end of this section) how the choice of different state variables may streamline the method.

⁸We here use $\sigma(t) = \sigma_x(t)^\top$.

$$\begin{aligned} p_i(0; K, x_2) &= P(0, T_i) e^{A(T_0, T_i) - x_2 G_2(T_0, T_i)} \\ &\quad \times \left(K^* \Phi(-d_+) - e^{\Omega(T_0, T_i, x_2)} \Phi(-d_-) \right), \end{aligned}$$

where

$$\begin{aligned} d_{\pm} &= \frac{\Omega(T_0, T_i, x_2) - \ln K^* \pm \frac{1}{2} G_1^2(T_0, T_i) s_1(T_0, x_2)^2}{G_1(T_0, T_i) s_1(T_0, x_2)}, \\ K^* &= \frac{P(0, T_0)}{P(0, T_i)} e^{-A(T_0, T_i) + x_2 G_2(T_0, T_i)} K, \\ \Omega(T_0, T_i, x_2) &= -\mu_1(T_0, x_2) G_1(T_0, T_i) + \frac{1}{2} G_1^2(T_0, T_i) s_1(T_0, x_2)^2, \end{aligned}$$

and $\mu_1(T_0, x_2)$ and $s_1(T_0, x_2)$ are given in (12.41)–(12.42) below.

Proof. A discount bond $P(t, T)$ has risk-neutral dynamics of the form (see Remark 12.1.4)

$$dP(t, T)/P(t, T) = r(t) dt - \sigma_{P,1}(t, T) dW_1(t) - \sigma_{P,2}(t, T) dW_2(t),$$

where

$$\begin{aligned} \sigma_{P,1}(t, T) &= G_1(t, T) \sigma_{11}(t) + G_2(t, T) \sigma_{21}(t), \\ \sigma_{P,2}(t, T) &= G_1(t, T) \sigma_{12}(t) + G_2(t, T) \sigma_{22}(t). \end{aligned}$$

From standard results (see Chapter 4), we know that $dW^{T_0}(t) = dW(t) + \sigma_P(t, T_0) dt$, where $\sigma_P(t, T_0) = (\sigma_{P,1}(t, T_0), \sigma_{P,2}(t, T_0))^T$, is a Brownian motion in Q^{T_0} . The Q^{T_0} dynamics for $x_1(t)$ and $x_2(t)$ therefore are

$$\begin{aligned} dx_1(t) &= \left(\vartheta_1^{T_0}(t) - \kappa_1(t) x_1(t) \right) dt + \sigma_{11}(t) dW_1^{T_0}(t) + \sigma_{12}(t) dW_2^{T_0}(t), \\ dx_2(t) &= \left(\vartheta_2^{T_0}(t) - \kappa_2(t) x_2(t) \right) dt + \sigma_{21}(t) dW_1^{T_0}(t) + \sigma_{22}(t) dW_2^{T_0}(t), \end{aligned}$$

where

$$\begin{aligned} \vartheta_1^{T_0}(t) &= \vartheta_1(t) - \sigma_{11}(t) \sigma_{P,1}(t, T_0) - \sigma_{12}(t) \sigma_{P,2}(t, T_0), \\ \vartheta_2^{T_0}(t) &= \vartheta_2(t) - \sigma_{21}(t) \sigma_{P,1}(t, T_0) - \sigma_{22}(t) \sigma_{P,2}(t, T_0). \end{aligned}$$

In measure Q^{T_0} , $x_1(T_0)$ and $x_2(T_0)$ are jointly Gaussian, with moments

$$\begin{aligned} E^{T_0}(x_i(T_0)) &= \int_0^{T_0} e^{-\int_s^{T_0} \kappa_i(u) du} \vartheta_i^{T_0}(s) ds, \quad i = 1, 2, \\ \text{Var}^{T_0}(x_1(T_0)) &= \int_0^{T_0} e^{-\int_s^{T_0} 2\kappa_1(u) du} (\sigma_{11}(s)^2 + \sigma_{12}(s)^2) ds, \\ \text{Var}^{T_0}(x_2(T_0)) &= \int_0^{T_0} e^{-\int_s^{T_0} 2\kappa_2(u) du} (\sigma_{21}(s)^2 + \sigma_{22}(s)^2) ds, \\ \text{Cov}^{T_0}(x_1(T_0), x_2(T_0)) &= \int_0^{T_0} e^{-\int_s^{T_0} (\kappa_1(u) + \kappa_2(u)) du} \\ &\quad \times (\sigma_{11}(s) \sigma_{21}(s) + \sigma_{12}(s) \sigma_{22}(s)) ds. \end{aligned}$$

Conditional upon $x_2(T_0)$, $x_1(T_0)$ must therefore be Gaussian with moments

$$\begin{aligned} \mathbb{E}^{T_0}(x_1(T_0)|x_2(T_0) = x_2) &= \mathbb{E}^{T_0}(x_1(T_0)) + \frac{\text{Cov}^{T_0}(x_1(T_0), x_2(T_0))}{\text{Var}^{T_0}(x_2(T_0))}x_2 \\ &\triangleq \mu_1(T_0, x_2), \end{aligned} \quad (12.41)$$

and

$$\begin{aligned} \text{Var}^{T_0}(x_1(T_0)|x_2(T_0) = x_2) &= \text{Var}^{T_0}(x_1(T_0)) - \frac{\text{Cov}^{T_0}(x_1(T_0), x_2(T_0))^2}{\text{Var}^{T_0}(x_2(T_0))} \\ &\triangleq s_1(T_0, x_2)^2. \end{aligned} \quad (12.42)$$

Using (12.40) we get, after a little rearrangement,

$$\begin{aligned} p_i(0; K, x_2) &= P(0, T_i)e^{A(T_0, T_i) - x_2 G_2(T_0, T_i)} \\ &\times \mathbb{E}^{T_0}\left(\left(K^* - e^{-x_1(T_0)G_1(T_0, T_i)}\right)^+ \middle| x_2(T_0) = x_2\right), \end{aligned}$$

where K^* was defined above. The result of the lemma now follows from the standard results for log-normal random variables. \square

Conditional on a value for $x_2(T_0)$, the payer swaption payout is monotonically increasing in $x_1(T_0)$, allowing for application of the Jamshidian decomposition to break the (conditional) swaption price into a sum of (conditional) discount bond options. A subsequent numerical integration against the density of $x_2(T_0)$ will then uncover the unconditional swaption price.

To formally state our result for the swaption price, define a function $x_1^*(x_2)$ as the solution to the equation $V_{\text{swaption}}(T_0, x_1^*(x_2), x_2) = 0$, or

$$1 - P(T_0, T_N; x_1^*(x_2), x_2) - c \sum_{i=0}^{N-1} \tau_i P(T_0, T_{i+1}; x_1^*(x_2), x_2) = 0,$$

where the functions $P(T_0, T_i; x_1^*(x_2), x_2)$ are given in Corollary 12.1.3. Given $x_1^*(x_2)$, we also define x_2 -dependent strikes

$$K_i(x_2) = P(T_0, T_i; x_1^*(x_2), x_2), \quad i = 1, \dots, N. \quad (12.43)$$

Proposition 12.1.18. *In the two-factor Gaussian model (12.38)–(12.39), let K_i be given by (12.43), and let x_2 -conditional discount bond put prices be given as in Lemma 12.1.17. Then, the swaption in (12.37) has price*

$$\begin{aligned} V_{\text{swaption}}(0) &= \int_{-\infty}^{\infty} \frac{p_N(0; K_N(x_2), x_2) + c \sum_{i=0}^{N-1} \tau_i p_{i+1}(0; K_{i+1}(x_2), x_2)}{\sqrt{\text{Var}^{T_0}(x_2(T_0))}} \\ &\quad \times \phi\left(\frac{x_2 - \mathbb{E}^{T_0}(x_2(T_0))}{\sqrt{\text{Var}^{T_0}(x_2(T_0))}}\right) dx_2, \end{aligned}$$

where $\phi(x)$ is the standard Gaussian density. The moments $E^{T_0}(x_2(T_0))$ and $\text{Var}^{T_0}(x_2(T_0))$ are given in Lemma 12.1.17.

Proof. Let $V(T_0, x_2)$ denote the swaption price at time T_0 , conditional on $x_2(T_0) = x_2$. If $x_2(T_0) = x_2$, we note that the swaption only pays out a positive amount if $x_1(T_0) > x_1^*(x_2)$. Following the argument in Section 10.1.3.1, we can then easily decompose the swaption payout as follows,

$$\begin{aligned} V(T_0, x_2) &= \left(1 - P(T_0, T_N; x_1(T_0), x_2) - c \sum_{i=0}^{N-1} \tau_i P(T_0, T_{i+1}; x_1(T_0), x_2) \right) \\ &\quad \times 1_{\{x_1(T_0) > x_1^*(x_2)\}} \\ &= (K_N(x_2) - P(T_0, T_N; x_1(T_0), x_2))^+ \\ &\quad + c \sum_{i=0}^{N-1} \tau_i (K_{i+1}(x_2) - P(T_0, T_{i+1}; x_1(T_0), x_2))^+. \end{aligned}$$

Clearly, then

$$\begin{aligned} V_{\text{swaption}}(0) &= P(0, T_0) E^{T_0}(V_{\text{swaption}}(T_0)) \\ &= P(0, T_0) \int_{\mathbb{R}} E^{T_0}(V(T_0, x_2(T_0))) Q^{T_0}(x_2(T_0) \in dx_2), \end{aligned}$$

and the result follows from the observation that $x_2(T_0)$ is Gaussian in measure Q^{T_0} , with moments given in Lemma 12.1.17. \square

The technique behind Proposition 12.1.18 extends in straightforward fashion to dimension $d > 2$, with the “unconditioning” step involving numerical integration against a $(d-1)$ -dimensional Gaussian density. This is rarely practical — especially since the integrand involves root-search to establish trigger levels for exercise — so in a real application we would typically never use Jamshidian decomposition, but instead introduce fast approximations. We list one such approximation in the next section.

As indicated earlier (in footnote 7), it is possible to make the derivation of the two-dimensional Jamshidian decomposition a little smoother by choosing another set of Markov state variables. To sketch how one might proceed, notice that, in the T_0 -forward measure,

$$\begin{aligned} dP(t, T_0, T)/P(t, T_0, T) &= - (G(t, T) - G(t, T_0))^{\top} \sigma_x(t)^{\top} dW^{T_0}(t) \\ &= - (\Lambda(T) - \Lambda(T_0))^{\top} H(t) \sigma_x(t)^{\top} dW^{T_0}(t), \end{aligned} \tag{12.44}$$

where $\Lambda(t) = (\Lambda_1(t), \Lambda_2(t))^{\top}$ is the two-dimensional vector given in (12.13), and $H(t)$ is the 2×2 diagonal matrix

$$H(t) = \text{diag} \left(\exp \left(\int_0^t \varkappa_1(s) ds \right), \exp \left(\int_0^t \varkappa_2(s) ds \right) \right).$$

Defining the two-dimensional Gaussian process

$$dz(t) = H(t)\sigma_x(t)^\top dW^{T_0}(t),$$

it follows from (12.44) that we can express forward bonds as closed form expressions of the Q^{T_0} -martingale process $z(t)$. The Q^{T_0} -dynamics for $z(t)$ are simpler than those of $x(t)$ (listed in the proof of Lemma 12.1.17), making subsequent manipulations easier.

12.1.6.2 Gaussian Swap Rate Approximation

We now return to the d -dimensional setting of Section 12.1.1. As in Section 5.10, we start by rewriting the swaption payout to

$$V_{\text{swaption}}(T_0) = A(T_0)(S(T_0) - c)^+,$$

where $A(t) = A_{0,N}(t)$ and $S(t) = S_{0,N}(t)$ are the swap annuity and par rate, respectively:

$$A(t) = \sum_{i=0}^{N-1} \tau_i P(t, T_{i+1}), \quad S(t) = \frac{P(t, T_0) - P(t, T_N)}{A(t)}.$$

Let Q^A be the measure induced by using $A(t)$ as the numeraire, such that

$$V_{\text{swaption}}(0) = A(0)E^A((S(T_0) - c)^+), \quad (12.45)$$

where E^A denotes expectation in measure Q^A . We know that $S(t)$ is a martingale in Q^A and, due to our Markov setting, a deterministic function of $x(t)$, i.e. $S(t) = S(t, x(t))$. It follows from Ito's lemma that

$$dS(t) = q(t, x(t))^\top \sigma_x(t)^\top dW(t), \quad (12.46)$$

where $q(t, x)$ is a d -dimensional column vector with elements

$$q_j(t, x) = \frac{\partial S(t)}{\partial x_j} = \frac{\partial}{\partial x_j} \frac{P(t, T_0, x) - P(t, T_N, x)}{\sum_{i=0}^{N-1} \tau_i P(t, T_{i+1}, x)}, \quad j = 1, \dots, d.$$

From the reconstitution formula in Corollary 12.1.3 we can evaluate the partial derivatives explicitly, yielding

$$\begin{aligned} q_j(t, x) &= -\frac{P(t, T_0, x)G_j(t, T_0) - P(t, T_N, x)G_j(t, T_N)}{A(t, x)} \\ &\quad - \frac{S(t, x) \sum_{i=0}^{N-1} \tau_i P(t, T_{i+1}, x)G_j(t, T_{i+1})}{A(t, x)}, \end{aligned}$$

where we recall that

$$G_j(t, T) = \int_t^T e^{-\int_t^u \kappa_j(s) ds} du.$$

The functions q_j can be experimentally verified to be close to a constant⁹ so, as a good approximation, we can write

$$q_j(t, x(t)) \approx q_j(t, \bar{x}(t)), \quad j = 1, \dots, N, \quad (12.47)$$

where $\bar{x}(t)$ is some *deterministic* proxy for the random vector $x(t)$. A reasonable approach is to set $\bar{x}(t) = 0$, but see Chapter 13 for refinements. In any case, with the approximation (12.47), the following swaption pricing formula is easily proven.

Lemma 12.1.19. *Let $\bar{x}(t)$ be a deterministic function of time, and assume that (12.47) holds. Then*

$$V_{\text{swaption}}(0) = A(0) [(S(0) - c) \Phi(d) + \sqrt{v} \phi(d)],$$

where

$$d = \frac{S(0) - c}{\sqrt{v}}, \quad v = \int_0^{T_0} \|q(t, \bar{x}(t))^{\top} \sigma_x(t)^{\top}\|^2 dt.$$

Proof. Follows directly from the Bachelier pricing formula (7.16), expression for the swap rate volatility (12.46), and approximation (12.47). \square

12.1.7 Calibration and Parameterization via Benchmark Rates

With the swaption formulas developed in the previous sections, we have reached a point where we can entertain ideas of how to calibrate a multi-factor Gaussian model. As one would expect, the basic principles of such a calibration are, for the most part, identical for all multi-factor models, and a detailed discussion is best postponed to later chapters when our model catalog is more complete. Nevertheless, it is useful to present here a few ideas that will align the Gaussian multi-factor model with later material, particular that in Chapter 14.

Using the fundamental model setup from Section 12.1.1, we first observe that the parameters $g(t)$ and $h(t)$ of the Gaussian model (in the parameterization from Proposition 12.1.1) affect both volatilities and correlations of market rates. Thus, should we desire to recover $g(t)$ and $h(t)$ via calibration to market instruments, we would need both caplet/swaptions (for overall level of volatility) and spread options¹⁰ (for correlation) as calibration targets. The former are much more liquid than the latter, so it would be

⁹Reflecting a rather intuitive fact that in a Gaussian model a swap rate is approximately Gaussian.

¹⁰Or other derivatives with first-order correlation dependence. See Chapter 17 for details on yield curve spread options.

beneficial to be able to separate volatility and correlation calibration; such separation of different model parameters is a good idea anyway, for any model. The parameterization of Proposition 12.1.1 is somewhat awkward for this purpose, so let us attempt to reparameterize the model in quantities that are more closely related to market observations.

It is intuitively obvious that a model with d stochastic factors is fully specified by volatilities and correlations of d (different) forward rates. Let us select d *benchmark tenors* $\delta_1 < \dots < \delta_d$, and define d *benchmark rates* $f_i(t) = f(t, t + \delta_i)$, $i = 1, \dots, d$. Notice that these forward rates are defined with “sliding” tenors, to encourage time-homogeneity. Let $\lambda_i(t)$ be the instantaneous volatility of the rate $f_i(t)$, $i = 1, \dots, d$, and let $\chi_{i,j}(t)$ be the instantaneous correlation between rates $f_i(t)$ and $f_j(t)$. To recover model parameters from this data, we first observe that the instantaneous covariance matrix of the vector

$$f(t) = (f_1(t), \dots, f_d(t))^\top$$

is given by $R^f(t) = \{\lambda_i(t)\lambda_j(t)\chi_{i,j}(t)\}_{i,j=1}^d$. On the other hand, in the model parameterized as in Proposition 12.1.1, the instantaneous covariance matrix is given by (see (12.5))

$$H^f(t)g(t)^\top g(t)H^f(t)^\top,$$

where the $d \times d$ matrix $H^f(t)$ is obtained by “stacking” vectors $h(t + \delta_i)$ together,

$$H^f(t) = \begin{pmatrix} h(t + \delta_1)^\top \\ \vdots \\ h(t + \delta_d)^\top \end{pmatrix}.$$

Let us assume that the vector $h(t)$ is directly parameterized, for instance by the specification of d different (and typically constant) mean reversion parameters κ_i , $i = 1, \dots, d$, to be applied in (12.12). It follows that the matrix $g(t)$ can then be recovered by solving

$$H^f(t)g(t)^\top = C^f(t)^\top, \quad (12.48)$$

where $C^f(t)$ is such that

$$R^f(t) = C^f(t)^\top C^f(t).$$

While this completely determines the model parameterization, (12.48) is not the only way to specify $g(t)$. For computational reasons we could, for example, decide to determine $g(t)$ by fitting *correlation* rather than *covariance*, as suggested in Andreasen [2005]. Defining $X^f(t) = \{\chi_{i,j}(t)\}$ and $D^f(t)$ by $X^f(t) = D^f(t)^\top D^f(t)$, we then obtain the matrix $g(t)$ by solving

$$H^f(t)g(t)^\top = \text{diag}((\lambda_1(t), \dots, \lambda_d(t))^\top) D^f(t)^\top. \quad (12.49)$$

The computational advantage of (12.49) over (12.48) is that the matrix $D^f(t)$ does not depend on $\lambda_i(t)$'s and, hence, does not need to be recomputed after each update of the volatilities $\lambda_i(t)$. (Such frequent updates happen, for example, during volatility calibration to quoted option prices.) The disadvantage of (12.49) is that the *true* instantaneous correlation matrix of the vector $f(t)$ is not going to be exactly equal to $X^f(t)$. (Similar issues are discussed at length in Section 14.3.4.) This is less of a concern if $X^f(t)$ itself is fit to the observed prices of correlation-dependent instruments such as spread options.

The ruminations above can be used to design a relatively straightforward calibration routine, a sketch of which is listed below. For full details, we again refer to Chapter 14.

1. Specify $h(t)$ via the mean reversion parameterization of Section 12.1.1.1, using d different constant mean reverptions. The choice of mean reverptions defines the interpolation of volatilities and correlations, i.e. how the volatilities/correlations of non-benchmark rates are obtained from those of the benchmark ones. As such, mean reverptions have rather limited impact on prices of exotic derivatives, as volatilities and correlations of benchmark rates — presumably chosen not at random but to represent the primary risk factors for a given pricing problem — are controlled directly in our setup. With that in mind, we advocate choosing mean reverptions in such a way as to improve the numerical properties of the algorithm. In particular, as an inversion of the matrix $H^f(t)$ is implicit in (12.48) (or (12.49)), we suggest using mean reverptions that are sufficiently different from each other so that the matrix $H^f(t)$ has a better-behaved inverse. Besides stable numerics, a good choice of mean reverptions will generate volatility factors that are fundamentally consistent with observed swaption quotes, in the sense that calibrated volatilities $\lambda(t)$ will be close to time-stationary. Working with a four-factor model, Andreasen [2005] suggests the following values:

$$\begin{aligned}\varkappa(t) &= \text{diag}((0.015, 0.15, 0.30, 1.20)^\top), \\ \{\delta_1, \delta_2, \delta_3, \delta_4\} &= \{6m, 2y, 10y, 30y\},\end{aligned}$$

which gives us a good example to follow.

2. Populate the correlation matrix $\{\chi_{i,j}(t)\}$ of benchmark rates. This may be done through a smooth functional form, as described in Chapter 14. Most often a time-homogeneous specification should be used. The parameters of the functional form may be found from historical analysis or through calibration to market prices of spread options.
3. Calibrate benchmark rate volatilities against swaptions, using the results of Section 12.1.6. For a discussion of optimization techniques and relevant calibration norms, see Chapter 14.
4. Recover $g(t)$, the diffusion matrix of factors, using (12.48) or (12.49).

With the function $h(t)$ and the correlation matrix $\{\chi_{i,j}(t)\}$ pre-specified, the model has enough parameters to calibrate d swaption strips (see Section 10.1.4). In most applications, we recommend choosing d strips with constant swap tenors matching benchmark tenors $\delta_1, \dots, \delta_d$. An alternative would be to do a best-fit calibration to all available options (i.e., a global calibration).

12.1.8 Monte Carlo Simulation

Monte Carlo methods for the d -dimensional Gaussian model are straightforward, as all state variables are jointly Gaussian. To demonstrate, we adopt the parameterization in Section 12.1.1.1 and consider pricing a security that pays an amount $V(T)$ at time T , where $V(T)$ may be a function of the entire path of the discount curve on $[0, T]$. Working in the risk-neutral measure, we must compute

$$\begin{aligned} V(0) &= \mathbb{E}^Q \left(V(T) e^{-\int_0^T r(u) du} \right) \\ &= P(0, T) \mathbb{E}^Q \left(V(T; \{x(t) : 0 \leq t \leq T\}) e^{-\int_0^T \mathbf{1}^\top x(u) du} \right), \end{aligned} \quad (12.50)$$

where we have used the relation $r(t) = f(0, t) + \mathbf{1}^\top x(t)$, and also have emphasized the dependence of $V(T)$ on the entire path of $x(\cdot)$. We assume that the determination of $V(T)$ involves observations of the yield curve on a discrete schedule $\{t_i\}_{i=0}^N$ with $t_0 = 0$ and $t_N = T$.

We recall from Proposition 12.1.2 that the risk-neutral dynamics for $x(t) = (x_1(t), \dots, x_d(t))^\top$ are

$$dx(t) = (y(t)\mathbf{1} - \kappa(t)x(t)) dt + \sigma_x(t)^\top dW(t), \quad x(0) = 0,$$

for deterministic vectors/matrices $y(t)$, $\kappa(t)$, $\sigma_x(t)$. Observe that $x(t_{i+1})$, conditional on $x(t_i)$, is d -dimensional Gaussian with mean

$$\mathbb{E}^Q(x(t_{i+1})|x(t_i)) = e^{-\int_{t_i}^{t_{i+1}} \kappa(u) du} x(t_i) + \int_{t_i}^{t_{i+1}} e^{-\int_s^{t_{i+1}} \kappa(u) du} y(s)\mathbf{1} ds,$$

and covariance matrix

$$\text{Var}(x(t_{i+1})|x(t_i)) = \int_{t_i}^{t_{i+1}} e^{-\int_s^{t_{i+1}} \kappa(u) du} \sigma_x(s) \sigma_x(s)^\top e^{-\int_s^{t_{i+1}} \kappa(u)^\top du} ds. \quad (12.51)$$

Let the square root of the covariance matrix be denoted C .

Advancement of $x(\cdot)$ from t_i to t_{i+1} can now proceed in an obvious fashion:

1. Draw d independent Gaussian samples $Z = (Z_1, Z_2, \dots, Z_d)^\top$.
2. Compute $Z^* = CZ$.

3. Set $x(t_{i+1}) = \mathbb{E}^Q(x(t_{i+1})|x(t_i)) + Z^*$.

At each time on the simulated path $x(t_0), \dots, x(t_N)$, we can use the reconstitution formula in Corollary 12.1.3 to reconstruct the entire discount curve, in turn allowing us to compute $V(T)$ on the path. To evaluate (12.50), it remains to simulate the quantity

$$I(T) = - \int_0^T \mathbf{1}^\top x(u) du$$

on the path. Clearly $I(t)$ is a Gaussian process, so we can work out the moments of $I(t_{i+1})$ conditional on $I(t_i)$ and $x(t_i)$, allowing for bias-free joint time-stepping of $x(t)$ and $I(t)$ on the schedule $\{t_i\}_{i=0}^N$. The analysis proceeds as in Section 10.1.6.1, and is omitted for brevity. In practice, we find that it is often more convenient to change the measure (as in Section 10.1.6.3) or to compute $I(T)$ by numerical integration, i.e.

$$I(T) \approx -\mathbf{1}^\top \sum_{i=1}^N x(t_i).$$

While the last method obviously introduces some amount of bias, this is generally of little concern and can be controlled as needed through the insertion of extra dates in the schedule.

Finally, we remind the reader once again that all time integrals involved in the Monte Carlo discretization scheme above should be pre-computed before actual path simulations commence.

12.1.9 Finite Difference Methods

We finish our treatment of the multi-factor Gaussian short rate model with a brief discussion of finite difference applications. For this, let us consider a claim V the terminal payout of which can be computed solely from knowledge of the yield curve at time T . We assume that the yield curve is driven by a multi-factor Gaussian model, in the form described in Section 12.1.1.1. In this model, the discount curve at time T can be reconstituted solely from knowledge of the Markovian state variable vector $x(T)$, so we may write $V(T) = V(T, x(T))$. By standard results from Chapter 1, for $t < T$ we then have $V(t) = V(t, x(t))$ where $V(t, x)$ satisfies a d -dimensional PDE:

$$\mathcal{L}V - (f(0, t) + \mathbf{1}^\top x) V = 0, \quad (12.52)$$

where \mathcal{L} is a partial differential operator,

$$\mathcal{L}V = V_t + \boldsymbol{\nu}(t) (\mathbf{y}(t)\mathbf{1} - x(t))^\top V_x + \frac{1}{2} \text{tr} (V_{xx} \sigma_x(t)^\top \sigma_x(t)).$$

In the definition of \mathcal{L} , $\text{tr}(\cdot)$ is the matrix trace operator, V_x is a d -dimensional vector of first-order spatial derivatives, and V_{xx} a $d \times d$ matrix of second-order spatial derivatives. The d -dimensional PDE (12.52) is subject to a given terminal condition $V(T, x)$, the computation of which typically would involve usage of the discount bond reconstitution formulas in Corollary 12.1.3.

Numerical solution of (12.52) is practically feasible provided that d does not exceed 3 or 4, say. We recommend the Craig-Sneyd scheme in Section 2.12. As this is a splitting scheme with only one dimension being computed non-explicitly in each split step, applications of the one-dimensional side-boundary conditions of Section 10.1.5 carry over in straightforward fashion. A similar comment holds for the application of upwinding in the edges of the finite difference grid.

12.2 The Affine Model

12.2.1 Introduction

In Section 10.2, we introduced the one-factor affine short rate model through the short rate SDE

$$\begin{aligned} dr(t) &= \mu(t, r(t)) dt + \sigma(t, r(t)) dW(t) \\ &= \kappa(t)(\vartheta(t) - r(t)) dt + \sqrt{\alpha + \beta r(t)} dW(t), \end{aligned}$$

for deterministic parameters $\kappa(t)$, $\vartheta(t)$, α , β , subject to certain regularity conditions. The fact that both $\mu(t, r)$ and $\sigma(t, r)^2$ were linear (that is, affine) in r allowed for a discount bond price formula that was exponentially affine in r . We are now interested in examining how we can extend this one-factor setup to a multi-dimensional one. Our discussion of this extension shall be quite brief though, as the general multi-factor affine class is of fairly limited practical relevance in securities pricing applications. A subset of affine models — those that can be rewritten as linear-quadratic interest rate models — have some uses in practice and shall receive a fuller treatment in Section 12.3.

12.2.2 Basic Model

Let us consider a time-homogeneous¹¹ Markov system of state variables

$$dx(t) = \mu(x(t)) dt + \sigma(x(t)) dW(t), \quad (12.53)$$

¹¹As is standard practice in much of the literature on multi-factor affine models, for notational simplicity we here assume that μ and σ do not depend on time t . See the comments at the end of Section 12.2.4 for extensions to time-dependent parameters.

where $x(t) = (x_1(t), \dots, x_d(t))^\top$ has state space $D \subseteq \mathbb{R}^d$; W is a d -dimensional Brownian motion in the risk-neutral probability measure Q ; and where $\mu : D \rightarrow \mathbb{R}^d$ and $\sigma : D \rightarrow \mathbb{R}^d \times \mathbb{R}^d$ have sufficient regularity for (12.53) to have a unique solution. We further write

$$r(t) = F(x(t)) \quad (12.54)$$

for some deterministic function $F : D \rightarrow \mathbb{R}$.

For the one-dimensional affine model in Section 10.2, we concluded that discount bonds were *exponential affine* in the underlying state variables. It is of interest to establish the circumstances for which something similar holds for (12.53)–(12.54). That is, we wish to determine the form that μ , σ , and F must take such that $P(t, T)$ can be written¹²

$$P(t, T) = \mathbb{E} \left(\exp \left(\int_t^T F(x(u)) du \right) \right) = \exp(A(T-t) - B(T-t)^\top x(t)),$$

for deterministic functions $A : \mathbb{R} \rightarrow \mathbb{R}$ and $B : \mathbb{R} \rightarrow \mathbb{R}^d$. The following result is shown in Duffie and Kan [1996].

Proposition 12.2.1. *Suppose that in (12.53)–(12.54) μ , $\sigma\sigma^\top$ and F are affine functions of x . Then discount bond prices are exponential affine.*

Remark 12.2.2. Under additional non-degeneracy assumptions, Duffie and Kan [1996] show that the converse of Proposition 12.2.1 holds. However, many interesting models, including those of Chapter 13, violate these conditions.

The proof of Proposition 12.2.1 follows the line of attack of Proposition 10.2.2 and is omitted for brevity.

Duffie and Kan [1996] demonstrate that if $\sigma\sigma^\top$ is affine one may, under some mild non-degeneracy conditions, rearrange the dynamics for $x(t)$ such that $\sigma(t)$ is *diagonal*. That is, we may write

$$dx(t) = (a - bx(t)) dt + \Sigma \begin{pmatrix} \sqrt{v_1(x(t))} & 0 & \ddots & 0 \\ 0 & \sqrt{v_2(x(t))} & \ddots & \ddots \\ \ddots & \ddots & \ddots & 0 \\ 0 & \ddots & 0 & \sqrt{v_d(x(t))} \end{pmatrix} dW(t), \quad (12.55)$$

where $a \in \mathbb{R}^d$ and $b, \Sigma \in \mathbb{R}^{d \times d}$ and

$$v_i(x) = \alpha_i + \beta_i^\top x, \quad i = 1, \dots, d, \quad (12.56)$$

¹²Throughout this section, \mathbb{E} denotes expectation in the risk-neutral measure Q .

with α_i scalar and $\beta_i \in \mathbb{R}^d$. The representation (12.55) is convenient in practical work and going forward shall, together with an affine short rate specification

$$r(t) = \xi_0 + \xi^\top x(t) = \xi_0 + \sum_{i=1}^d \xi_i x_i(t), \quad (12.57)$$

constitute our working definition for a multi-factor affine model.

12.2.3 Regularity Issues

As was the case for the one-factor affine model, there are strong restrictions on which values of $\Sigma, \alpha_i, \beta_i, a, b$ allow for valid solutions to (12.55). To state these restrictions, we first establish the process for $v_i(x(t))$ to be

$$d(v_i(x(t))) = \beta_i^\top dx(t) = \beta_i^\top (a - bx(t)) dt + \beta_i^\top \Sigma v(x(t)) dW(t), \quad (12.58)$$

where $v(x(t))$ is the diagonal matrix in (12.55). It is clear that the volatility process $v_i(x(t))$ must be non-negative for all i and t , i.e. the valid domain for $x(t)$ is

$$D = \{x \in \mathbb{R}^d : v_i(x) \geq 0, \quad i = 1, \dots, d\}. \quad (12.59)$$

Clearly, when $x(t)$ is on the boundary $\partial_i D = \{x \in D : v_i(x) = 0\}$, we must ensure that i) the drift of $v_i(x(t))$ is positive; and ii) the instantaneous variance is zero. The first condition implies that, for all i ,

$$\beta_i^\top (a - bx) \geq 0, \quad \forall x \in \partial_i D.$$

The second condition requires that, for all i ,

$$\beta_i^\top \Sigma v(x) = 0, \quad \forall x \in \partial_i D.$$

This evidently requires that for $j = 1, \dots, d$ either $v_j(x) = k v_i(x)$ for some constant k ; or the j -th element of the row vector $\beta_i^\top \Sigma$ is zero, i.e. $(\beta_i^\top \Sigma)_j = 0$. As the constant k can be absorbed into the definition of Σ , we may simplify the first condition to $v_j(x) = v_i(x)$. These results motivate the following theorem, the detailed proof of which can be found in Duffie and Kan [1996].

Theorem 12.2.3. Consider the SDE (12.55) with domain D given in (12.59) and assume that for all $i = 1, \dots, d$,

1. $\beta_i^\top (a - bx) \geq 0$ for all x such that $v_i(x) = 0$.
2. For all $j = 1, \dots, d$, if $(\beta_i^\top \Sigma)_j \neq 0$ then $v_i(x) = v_j(x)$.

Then there exists a unique strong solution to (12.55) with $x(t) \in D$. If additionally $x(0)$ is such that $v_i(x(0)) > 0$ for all i , then also $v_i(x(t)) > 0$ provided that we replace Condition 1 with the stronger criterion

- 1*. $\beta_i^\top (a - bx) \geq \beta_i^\top \Sigma \Sigma^\top \beta_i / 2$ for all x such that $v_i(x) = 0$.

Remark 12.2.4. We recognize the criterion 1* above as a multi-variate generalization of the *Feller condition*, first encountered in Section 8.3.

We should emphasize that the regularity conditions outlined in Theorem 12.2.3 for the affine state variable SDE to be well-defined are quite strong and rule out many seemingly reasonable model specifications. Section 12.2.5 list some concrete models that satisfy the conditions of the theorem.

12.2.4 Discount Bond Prices

As advertised earlier, the main advantage of affine multi-factor models is the existence of discount bond reconstitution formulas that involve only the solution of ordinary Riccati ODEs. To develop this result, let

$$P(t, T, x) = \mathbb{E} \left(e^{-\int_t^T r(u) du} \middle| x(t) = x \right),$$

with $r(t)$ computed from (12.57) and the dynamics of the state variable vector $x(t)$ given in (12.55). From the Feynman-Kac result, we must have

$$\mathcal{L}P - (\xi_0 + \xi^\top x) P = 0, \quad (12.60)$$

where, using the same notation as in Section 12.1.9, \mathcal{L} is a partial differential operator

$$\mathcal{L}P = P_t + (a - bx)^\top P_x + \frac{1}{2} \text{tr} (P_{xx} \Sigma v(x) v(x)^\top \Sigma^\top).$$

Earlier results from Section 10.2.2 strongly hint that we should look for a solution of the form

$$P(t, T, x) = e^{A(T-t) - B(T-t)^\top x}, \quad (12.61)$$

where $A : \mathbb{R} \rightarrow \mathbb{R}$ and $B : \mathbb{R} \rightarrow \mathbb{R}^d$ are unknown deterministic functions. Inserting this solution into (12.60) and using the “matching principle”¹³ we get the following result.

Proposition 12.2.5. *The solution to (12.60) is*

$$P(t, T, x) = \exp (A(T-t) - B(T-t)^\top x),$$

where the real-valued function $A(\tau)$ and the vector-valued function $B(\tau)$ satisfy the system of Riccati ODE equations

$$\begin{aligned} \frac{d}{d\tau} B(\tau) &= -b^\top B(\tau) - \frac{1}{2} \beta^\top \text{diag}(\Sigma^\top B(\tau)) \Sigma^\top B(\tau) + \xi, \\ \frac{d}{d\tau} A(\tau) &= -a^\top B(\tau) + \frac{1}{2} \alpha^\top \text{diag}(\Sigma^\top B(\tau)) \Sigma^\top B(\tau) - \xi_0, \end{aligned}$$

¹³If $a + b^\top x = c + d^\top x$ holds for a non-empty open set, then $a = c$ and $b = d$.

with initial conditions $A(0) = 0$, $B(0) = 0$, where $\alpha = (\alpha_1, \dots, \alpha_d)^\top$ and the i -th row of matrix β is given by β_i^\top , $i = 1, \dots, d$. The ODE system can be written component-wise as

$$\begin{aligned} \frac{d}{d\tau} B_i(\tau) &= - \sum_{j=1}^d b_{j,i} B_j(\tau) \\ &\quad - \frac{1}{2} \sum_{k=1}^d \beta_{k,i} \left(\sum_{j=1}^d \Sigma_{j,k} B_j(\tau) \right)^2 + \xi_i, \quad i = 1, \dots, d, \\ \frac{d}{d\tau} A(\tau) &= - \sum_{j=1}^d a_j B_j(\tau) + \frac{1}{2} \sum_{k=1}^d \alpha_k \left(\sum_{j=1}^d \Sigma_{j,k} B_j(\tau) \right)^2 - \xi_0. \end{aligned}$$

While analytical solution of the ODEs in Proposition 12.2.5 is sometimes possible (see Sections 12.2.5.2 and 12.2.5.3), in general one has to rely on numerical solution. The Runge-Kutta algorithm is a good choice for this.

An application of Proposition 12.2.5 reveals that discount bond price dynamics are

$$dP(t, T)/P(t, T) = r(t) dt - B(T-t)^\top \Sigma v(x(t)) dW(t),$$

and that forward rate dynamics are

$$df(t, T) = O(dt) + \frac{\partial B(T-t)^\top}{\partial T} \Sigma v(x(t)) dW(t),$$

where the drift term can be computed from the HJM results in Section 4.4. It follows that

$$\begin{aligned} \text{Corr}(df(t, T_1), df(t, T_2)) \\ = \frac{Y(T_1-t, x(t))^\top Y(T_2-t, x(t))}{\sqrt{Y(T_1-t, x(t))^\top Y(T_1-t, x(t))} \sqrt{Y(T_2-t, x(t))^\top Y(T_2-t, x(t))}}, \end{aligned}$$

where

$$Y(T-t, x) = v(x) \Sigma^\top \frac{\partial B(T-t)}{\partial T}.$$

Unlike the forward rate correlations computed earlier for the multi-factor Gaussian model, we notice that in the affine model $\text{Corr}(df(t, T_1), df(t, T_2))$ generally depends on the random variable $x(t)$ and hence is *stochastic*.

Before moving on to concrete model examples, let us note that it is possible to extend the affine model to have time-dependent coefficients. In this case, the bond price equation would be

$$P(t, T, x) = \exp \left(A(t, T) - \sum_{i=1}^d B_i(t, T) x_i \right),$$

where, for $i = 1, 2, \dots, d$,

$$\begin{aligned} \frac{d}{dt} B_i(t, T) &= \sum_{j=1}^d b_{j,i}(t) B_j(t, T) \\ &\quad + \frac{1}{2} \sum_{k=1}^d \beta_{k,i}(t) \left(\sum_{j=1}^d \Sigma_{j,k}(t) B_j(t, T) \right)^2 - \xi_i(t), \\ \frac{d}{dt} A(t, T) &= \sum_{j=1}^d a_j(t) B_j(t, T) \\ &\quad - \frac{1}{2} \sum_{k=1}^d \alpha_k(t) \left(\sum_{j=1}^d \Sigma_{j,k}(t) B_j(t, T) \right)^2 + \xi_0(t), \end{aligned}$$

subject to $A(T, T) = B_1(T, T) = \dots = B_d(T, T) = 0$. A certain amount of time-dependence would always be required in order to calibrate the model to the initial yield curve and to observed option prices.

12.2.5 Some Concrete Models

12.2.5.1 Fong-Vasicek Model

In Section 11.2.3 we encountered the two-factor model

$$\begin{aligned} dr(t) &= \kappa_r (\vartheta_r - r(t)) dt + \sqrt{z(t)} dW_1(t), \\ dz(t) &= \kappa_z (\vartheta_z - z(t)) dt + \eta \sqrt{z(t)} (\rho dW_1(t) + \sqrt{1 - \rho^2} dW_2(t)). \end{aligned}$$

This model can be folded into the framework in Section 12.2.2 by writing $r(t) = x_1(t)$ (i.e., $\xi_0 = \xi_2 = 0$, $\xi_1 = 1$ in (12.57)), $z(t) = x_2(t)$, and

$$\begin{aligned} d \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} &= \left(\begin{pmatrix} \kappa_r \vartheta_r \\ \kappa_z \vartheta_z \end{pmatrix} - \begin{pmatrix} \kappa_r & 0 \\ 0 & \kappa_z \end{pmatrix} \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} \right) dt \\ &\quad + \begin{pmatrix} 1 & 0 \\ \eta \rho & \eta \sqrt{1 - \rho^2} \end{pmatrix} \begin{pmatrix} \sqrt{x_2(t)} & 0 \\ 0 & \sqrt{x_2(t)} \end{pmatrix} d \begin{pmatrix} W_1(t) \\ W_2(t) \end{pmatrix}. \end{aligned}$$

Clearly this is of the form (12.55). It is easy to verify that the restrictions in Theorem 12.2.3 are all satisfied here.

12.2.5.2 Longstaff-Schwartz Model

Longstaff and Schwartz [1992] have proposed a two-factor extension of the CIR model we encountered in Section 10.2. In the language of (12.55)–(12.57),

the Longstaff-Schwartz (LS) model can be written as $r(t) = x_1(t) + x_2(t)$ ($\xi_0 = 0$, $\xi_1 = \xi_2 = 1$) with risk-neutral dynamics of the form

$$\begin{aligned} d\begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} &= \left(\begin{pmatrix} \kappa_1 \vartheta_1 \\ \kappa_2 \vartheta_2 \end{pmatrix} - \begin{pmatrix} \kappa_1 & 0 \\ 0 & \kappa_2 \end{pmatrix} \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} \right) dt \\ &\quad + \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix} \begin{pmatrix} \sqrt{x_1(t)} & 0 \\ 0 & \sqrt{x_2(t)} \end{pmatrix} d\begin{pmatrix} W_1(t) \\ W_2(t) \end{pmatrix}. \end{aligned} \quad (12.62)$$

Again, it is easy to verify that the regularity conditions of Theorem 12.2.3 hold for (12.62). We notice that here the two state variables $x_1(t)$ and $x_2(t)$ are independent and both are time-homogeneous CIR. The independence assumption ensures that the analytical results from Section 10.2.2.1 can be used to solve the Riccati ODEs in Proposition 10.2.4 analytically. For completeness, we list the result below.

Lemma 12.2.6. *For the LS model (12.62), discount bond prices can be computed by*

$$P(t, T) = \exp(A_1(T-t) + A_2(T-t) - B_1(T-t)x_1(t) - B_2(T-t)x_2(t)),$$

where, for $i = 1, 2$,

$$\begin{aligned} A_i(\tau) &= \kappa_i \vartheta_i \sigma_i^{-2} (\kappa_i + \gamma_i) \tau - 2\kappa_i \vartheta_i \sigma_i^{-2} \ln \left(1 + \frac{(\kappa_i + \gamma_i)(e^{\gamma_i \tau} - 1)}{2\gamma_i} \right), \\ B_i(\tau) &= \frac{2(1 - e^{-\gamma_i \tau})}{(\kappa_i + \gamma_i)(1 - e^{-\gamma_i \tau}) + 2\gamma_i e^{-\gamma_i \tau}}, \end{aligned}$$

with $\gamma_i = \sqrt{\kappa_i^2 + 2\sigma_i^2}$.

Proof. From independence

$$P(t, T) = \mathbb{E} \left(e^{- \int_t^T r(u) du} \right) = \mathbb{E} \left(e^{- \int_t^T x_1(u) du} \right) \mathbb{E} \left(e^{- \int_t^T x_2(u) du} \right).$$

Application of the result in Proposition 10.2.4 (with $c_1 = 0$ and $c_2 = 1$) then proves the lemma. \square

We should note that it is common to reparameterize the LS model in terms of $r(t)$ and $v(t) \triangleq d\text{Var}_t(dr(t))/dt$, particularly when performing time series estimation. To quickly demonstrate the basic idea, notice that

$$\begin{aligned} dr(t) &= \kappa_1 (\vartheta_1 - x_1(t)) dt + \kappa_2 (\vartheta_2 - x_2(t)) dt \\ &\quad + \sigma_1 \sqrt{x_1(t)} dW_1(t) + \sigma_2 \sqrt{x_2(t)} dW_2(t), \end{aligned}$$

such that

$$v(t) = \sigma_1^2 x_1(t) + \sigma_2^2 x_2(t). \quad (12.63)$$

Combining (12.63) with the equation $r(t) = x_1(t) + x_2(t)$ allows us to write, provided $\sigma_1 \neq \sigma_2$,

$$x_1(t) = \frac{\sigma_2^2 r(t) - v(t)}{\sigma_2^2 - \sigma_1^2}, \quad x_2(t) = \frac{v(t) - \sigma_1^2 r(t)}{\sigma_2^2 - \sigma_1^2}.$$

From this, it is possible to eliminate $x_1(t)$ and $x_2(t)$ from the SDEs for $r(t)$ and $v(t)$, resulting in a Markov model with $r(t)$ and $v(t)$ being the only state variables. We leave the details to the reader (or see Longstaff and Schwartz [1992]).

As shown in Longstaff and Schwartz [1992], the two-factor time-homogeneous specification (12.62) allows one to produce a substantially richer set of yield curve shapes than an ordinary one-factor CIR model. Of course, without the introduction of time-dependence in one or more parameters¹⁴ (or application of the Dybvig “trick” from Section 11.3.2.4), the model will still never be able to perfectly fit the market-observed yield curve. The question of how to make an LS model acceptable for derivatives pricing purposes (which would necessarily involve further time-dependence and a scheme to allow for calibration to option prices) is of limited interest to us here, so we skip it and just point to Sections 10.2 and 12.3 for some general ideas. See also Clewlow and Strickland [1994] where some practical issues in parameter estimation for the LS model are discussed.

As a final remark, let us mention that the time-homogeneous LS model allows for an analytical pricing formula for European options on discount bonds. As both $x_1(t)$ and $x_2(t)$ are non-central chi-square random variables (see Section 8.3), the pricing formulas involve a two-dimensional non-central chi-square distribution, the practical computation of which is discussed in Chen and Scott [1992]. As time-homogeneous specifications are of little interest to the applications in this book, we do not list the pricing formulas, but simply refer to Longstaff and Schwartz [1992] for the details.

12.2.5.3 Multi-Factor CIR Models

A d -factor extension of the two-factor model in Section 12.2.5.2 would involve a system of decoupled SDEs of the form

$$dx_i(t) = \varkappa_i (\vartheta_i - x_i(t)) dt + \sigma_i \sqrt{x_i(t)} dW_i(t), \quad i = 1, \dots, d, \quad (12.64)$$

with $r(t) = \sum_i x_i(t)$ and all Brownian motions $W_1(t), \dots, W_d(t)$ mutually independent. The model satisfies the regularity conditions in Theorem 12.2.3 and it is clear from results in Section 8.3 that the resulting model will imply a non-negative short rate process. In fact, the short rate process will be strictly positive provided that there exists at least one $i \in \{1, \dots, d\}$ for which the Feller condition $2\varkappa_i \vartheta_i \geq \sigma_i^2$ is satisfied. Discount bond pricing in the model (12.64) can be done analytically by re-using the one-factor results in Section 10.2.2.1, in the same manner as in Lemma 12.2.6. We leave the details to the reader. The uncoupled nature of the SDEs for the

¹⁴Longstaff and Schwartz [1993] suggest making \varkappa_2 a function of time.

various $x_i(t)$ in (12.64) is rather convenient as it allows us to reuse analytical and numerical techniques from Section 10.2. For instance, Monte Carlo simulation of (12.64) can proceed by simply simulating each $x_i(t)$ according to the scheme discussed in Section 10.2.7.

12.2.6 Brief Notes on Option Pricing

Pricing of contingent claims with no path-dependence can be done via the PDE (12.60), the solution of which would often proceed by finite difference methods, at least if the dimension d is modest. See Section 12.1.9 for further details. When model dimension is high or the payout is path-dependent, Monte Carlo methods are required. In some cases (as in Section 12.2.5.3 above), Monte Carlo discretization of d -dimensional affine models is a straightforward application of one-dimensional schemes.

To calibrate affine models to market option data, it is, as always, important to have fast schemes for swaption pricing. Without going into details, we note that the ideas laid out earlier in Section 10.2 for the one-factor models may be applied to the multi-factor affine models as well. As our treatment of (and interest in) affine models is rather cursory, we just refer the reader to material in Section 12.3 and Chapter 13. We should also note the existence of several dedicated swaption approximations in the literature for affine models; see Collin-Dufresne and Goldstein [2002a] for an example and further references.

12.3 The Quadratic Gaussian Model

A sub-class of affine models called *quadratic Gaussian* (QG) models is particularly attractive for practical applications. While currently more familiar to academics than to practitioners (see Chen et al. [2004], Ahn et al. [2002], Assefa [2007]) the quadratic Gaussian models have several appealing properties: they are Markovian in a finite number of state variables, the state variables are Gaussian facilitating fast simulation, and the models in this class are capable of generating volatility smiles that can be parameterized in an intuitive way.

A QG model is obtained by generalizing (12.4) to include a quadratic term:

$$r(t) = z(t)^\top \gamma(t) z(t) + h(t)^\top z(t) + a(t), \quad (12.65)$$

where $\gamma(t)$ is a $d \times d$ symmetric matrix and, as before, $h(t)$ is a d -dimensional vector. The scalar function $a(t)$ is used to fit the initial yield curve. The state variable vector $z(t)$ follows (12.3), i.e.

$$dz(t) = g(t)^\top dW(t), \quad z(0) = 0, \quad (12.66)$$

with $W(t)$ a Brownian motion under the risk-neutral measure. Just like a linear Gaussian model, a quadratic model can be expressed in terms of mean-reverting state variables:

$$r(t) = x(t)^\top \tilde{\gamma}(t)x(t) + \mathbf{1}^\top x(t) + a(t), \quad (12.67)$$

where

$$\tilde{\gamma}(t) = H(t)^{-1}\gamma(t)H(t)^{-1},$$

and the transformed state variables $x(t)$, defined by $x(t) = H(t)z(t)$, follow

$$dx(t) = -\varkappa(t)x(t) dt + (g(t)H(t))^\top dW(t),$$

with $H(t)$ defined by (12.7) and $\varkappa(t)$ defined by (12.8). Other linear transformations, along the lines of Section 12.1.1.2, are also possible. Such representations may provide certain advantages for model implementation and numerical methods, as we often prefer to keep the diffusion term as constant in time as possible. Nevertheless, to reduce notational clutter we here stick to the driftless form (12.65)–(12.66).

We briefly saw a one-dimensional quadratic Gaussian model in Section 10.2.6; here, following Piterbarg [2009a], we study the multi-dimensional version.

12.3.1 Quadratic Gaussian Models are Affine

To show that the model defined by (12.65)–(12.66) is indeed affine, we introduce a vector of extra state variables $u(t)$ of length d^2 , whose elements are pairwise products of the coordinates of $z(t)$, i.e.

$$u(t) = (z_1(t)z_1(t), z_1(t)z_2(t), z_1(t)z_3(t), \dots, z_d(t)z_d(t))^\top.$$

Then, clearly, $r(t)$ is a linear function of $(z(t), u(t))$, and the coefficients of the SDE for the matrix $z(t)z(t)^\top$ are linear in $z(t)$:

$$d(z(t)z(t)^\top) = g(t)^\top dW(t) z(t)^\top + z(t) dW(t)^\top g(t) + (g(t)^\top g(t)) dt.$$

As we can write $u(t)$ by “unwrapping” the rows of $z(t)z(t)^\top$ into a vector, it follows that the coefficients of the SDE for $du(t)$ are linear in $z(t)$.

The analysis above makes it clear that the combined state variable vector $(z(t), u(t))$ has multi-factor affine dynamics. Represented as a standard affine model, the quadratic model would evidently require a total of $d(d + 1)$ state variables (rather than just d), so there is often good reason *not* to use such a representation explicitly. We note in passing that the quadratic parameterization satisfy the regularity constraints from Section 12.2.3 by construction.

12.3.2 The Basics

Since the QG model (12.65)–(12.66) is affine, it should come as no surprise that bond reconstruction formulas are available.

Proposition 12.3.1. *In the quadratic Gaussian model (12.65)–(12.66), zero-coupon discount bonds are exponentials of quadratic forms,*

$$\begin{aligned} & -\ln P(t, T) \\ &= z(t)^\top \gamma(t, T) z(t) + h(t, T)^\top z(t) + a(t, T) - \ln(P(0, T)/P(0, t)), \end{aligned}$$

with $\gamma(t, T)$, $h(t, T)$ satisfying Riccati equations

$$\begin{aligned} & -\frac{d}{dt} \gamma(t, T) + 2\gamma(t, T) g(t) g(t)^\top \gamma(t, T) = \gamma(t), \\ & -\frac{d}{dt} h(t, T) + 2\gamma(t, T) g(t) g(t)^\top h(t, T) = h(t), \end{aligned} \quad (12.68)$$

with terminal conditions $\gamma(T, T) = 0$, $h(T, T) = 0$.

Proof. By the same arguments as Proposition 12.2.5. \square

Remark 12.3.2. The function $a(t, T)$ also satisfies a Riccati equation; however, we find that it is better to determine it from the no-arbitrage condition $P(0, t)E^t(P(t, T)) = P(0, T)$, where E^t is the expected value operator under the t -forward measure Q^t , so that

$$a(t, T) = \ln E^t \left(\exp \left(-z(t)^\top \gamma(t, T) z(t) - h(t, T)^\top z(t) \right) \right). \quad (12.69)$$

To calculate $a(\cdot, T)$ in (12.69), we need to know the distribution of $z(t)$ under the t -forward measure Q^t ; this distribution will also be of general use in option pricing. From Proposition 12.3.1,

$$\frac{dP(t, T)}{P(t, T)} = O(dt) - \left(2z(t)^\top \gamma(t, T) + h(t, T)^\top \right) g(t)^\top dW(t),$$

so by Girsanov's theorem (Theorem 1.5.1),

$$dW^T(t) = dW(t) + g(t) (2\gamma(t, T) z(t) + h(t, T)) dt \quad (12.70)$$

is a Brownian motion under the T -forward measure Q^T . We use this fact to obtain the following result.

Proposition 12.3.3. *In the quadratic Gaussian model (12.65)–(12.66), the dynamics of the state process $z(t)$ in the T -forward measure Q^T are given by*

$$dz(t) = (m^T(t) - k^T(t)z(t)) dt + g(t)^\top dW^T(t), \quad (12.71)$$

where

$$k^T(t) = 2g(t)^\top g(t)\gamma(t, T), \quad m^T(t) = -g(t)^\top g(t)h(t, T),$$

and $W^T(t)$ is a \mathbb{Q}^T -Brownian motion. In particular, $z(t)$ is a Gaussian process under any T -forward measure, and is given in the integrated form by

$$\begin{aligned} z(s) &= J_{k^T}(s) \left[J_{k^T}(t)^{-1} z(t) + \int_t^s J_{k^T}(u)^{-1} m^T(u) du \right. \\ &\quad \left. + \int_t^s J_{k^T}(u)^{-1} g(u)^\top dW^T(u) \right], \end{aligned} \quad (12.72)$$

where the matrix-valued function $J_{k^T}(t)$ is defined by (12.15), i.e. satisfies an ODE

$$\frac{d}{dt} J_{k^T}(t) = -2g(t)^\top g(t)\gamma(t, T) J_{k^T}(t), \quad J_{k^T}(0) = I. \quad (12.73)$$

Proof. The equation (12.71) follows from (12.70). Integrating a linear Gaussian SDE (12.71) we obtain (12.72), see Lemma 12.1.6 or Karatzas and Shreve [1997]. \square

As $z(t)$ is Gaussian under any forward measure, its distribution is fully specified by its first and second moments.

Proposition 12.3.4. *In the quadratic Gaussian model (12.65)–(12.66), the conditional moments of the Gaussian state process $z(t)$ under the T -forward measure \mathbb{Q}^T ,*

$$m^T(t, s, z) \triangleq \mathbb{E}^T(z(s) | z(t) = z), \quad (12.74)$$

$$\begin{aligned} \nu^T(t, s, z) &\triangleq \text{Var}^T(z(s) | z(t) = z) \\ &= \mathbb{E}^T \left((z(s) - m^T(t, s, z)) (z(s) - m^T(t, s, z))^\top \middle| z(t) = z \right), \end{aligned} \quad (12.75)$$

are given by

$$m^T(t, s, z) = J_{k^T}(s) J_{k^T}(t)^{-1} z - J_{k^T}(s) \int_t^s J_{k^T}(u)^{-1} g(u)^\top g(u) h(u, T) du, \quad (12.76)$$

$$\nu^T(t, s, z) = J_{k^T}(s) \left(\int_t^s J_{k^T}(u)^{-1} g(u)^\top g(u) (J_{k^T}(u)^{-1})^\top du \right) J_{k^T}(s)^\top, \quad (12.77)$$

where $J_{k^T}(t)$ is defined by (12.73).

Proof. Follows immediately from (12.72). \square

To compute the function $a(t, T)$ via (12.69) we need to know the moment-generating function of a quadratic form of a Gaussian vector.

Proposition 12.3.5. Let Z be a K -dimensional Gaussian vector with mean m and variance ν . Let Q be a symmetric $K \times K$ matrix and u a K -dimensional vector. Define

$$\Psi(u, Q; m, \nu) \triangleq \ln E(\exp(Z^\top QZ + u^\top Z)).$$

If $\det(I - 2Q\nu) > 0$, then

$$\begin{aligned} \Psi(u, Q; m, \nu) = & \frac{1}{2} (2m^\top Q + u^\top) \nu (I - 2Q\nu)^{-1} (2Qm + u) \\ & + m^\top Qm + u^\top m - \frac{1}{2} \ln(\det(I - 2Q\nu)). \end{aligned}$$

Proof. In Appendix 12.A. \square

In addition to the proof of Proposition 12.3.5, Appendix 12.A contains a number of technical results useful for working with quadratic forms of Gaussian vectors.

The QG model is Markovian in d state variables (the vector z) and it should not be surprising that, with the help of the quadratic term, it can generate a genuine U-shaped volatility smile (see Figure 12.3). The state vector follows a Gaussian process, and it can be simulated at minimal computational cost (see Section 12.3.7). While these properties make the quadratic model both flexible and numerically efficient, its practical usefulness ultimately hinges on our ability to parameterize it in a sensible and intuitive way. Such a task is not trivial given that the generic time-dependent quadratic term $\gamma(t)$ is, essentially, unconstrained. While the richness of the model allows for a potentially large number of parameterization strategies, we here have little interest in exhaustive classification and content ourselves with presenting just one possible — and quite reasonable, we think — approach.

12.3.3 Parameterization

12.3.3.1 Smile Generation

To devise a parameterization strategy for the QG model, it is useful to understand the mechanism by which it generates a volatility smile. As the one-factor case is somewhat degenerate, we first look at the two-factor case for inspiration. We find it convenient to parameterize the quadratic term in such a way that we can identify one state variable as a “curve” factor and the other as a “volatility” factor (see Tezier [2005]). With that in mind, we set $d = 2$, $g_{1,2}(t) = g_{2,1}(t) = 0$, $h_1(t) = e^{-\kappa t}$, $h_2(t) = 0$,

$$\gamma_{1,1}(t) = \eta \bar{\omega} h_1(t)^2, \quad \gamma_{1,2}(t) = \frac{\eta}{2} \bar{\omega} h_1(t), \quad \gamma_{2,2}(t) = 0, \quad (12.78)$$

where $\bar{\omega} = \sqrt{1 - \omega^2}$. According to (12.65), the short rate is then given by

$$r(t) = (1 + \eta v(t)) h_1(t) z_1(t) + a(t), \quad (12.79)$$

where

$$v(t) = \varpi h_1(t) z_1(t) + \bar{\varpi} z_2(t).$$

If $\eta = 0$, the expression for the short rate reduces to

$$r(t) = h_1(t) z_1(t) + a(t),$$

and the model then becomes a one-factor (linear) Gaussian,

$$dr(t) = \nu(\vartheta(t) - r(t)) dt + e^{-\nu t} g_{1,1}(t) dW_1(t), \quad \vartheta(t) = a(t) + a'(t)/\nu.$$

Fittingly, we can identify $h_1(t) z_1(t)$ as a *curve factor*, i.e. as the factor that drives the state of the yield curve. If $\eta \neq 0$, the short rate is given by the curve factor times $1 + \eta v(t)$. As high values of $v(t)$ imply high volatility of $r(t)$, $\eta v(t)$ plays the role of¹⁵ “stochastic volatility”. Consequently, η may be interpreted as a volatility of volatility parameter. We notice that the volatility factor $v(t)$ is a linear combination of the curve factor $h_1(t) z_1(t)$ and a process $z_2(t)$ which is independent of the curve factor. The parameter ϖ therefore determines the correlation between the curve factor and the volatility factor.

As one would intuitively guess, the model outlined above is capable of producing volatility smiles that are similar to those of the stochastic volatility models we encountered in Chapter 8. Figure 12.3 shows a sample fit of the QG model to a market-implied volatility smile.

The parameterization (12.78) not only serves to identify one of the state variables as a curve factor and the other as a volatility factor, it also conceptually separates parameters that affect the volatility structure of the model ($h_1(t)$, $g_{1,1}(t)$), and those that affect the volatility smile (η , ϖ). Such separation is very convenient for building intuition for model dynamics and for the development of efficient European option approximations and practical calibration algorithms. Consequently, we seek to impose a similar structure as we build QG models of dimensions higher than two.

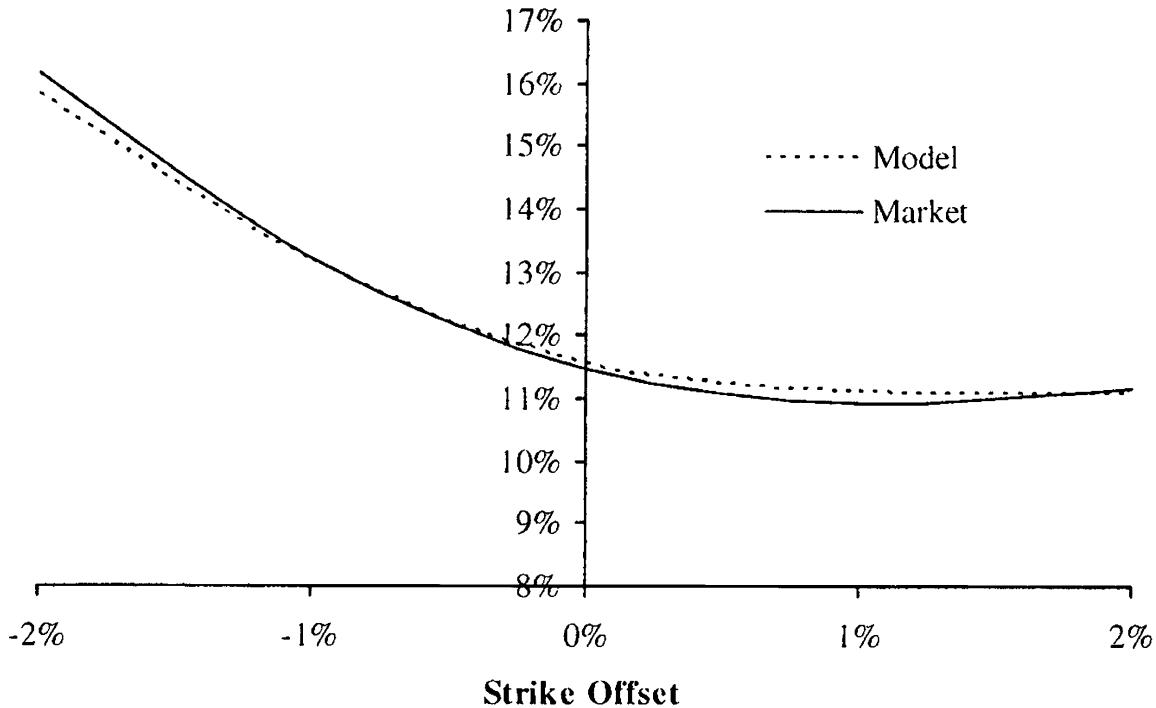
12.3.3.2 Quadratic Term

Given a budget of d curve factors and 1 volatility factor, we follow the example from the previous section and use a linear function of the d curve factors to define the yield curve dynamics, and the volatility factor to drive multiplicative scaling.

Let $z(t)$ be a $(d+1)$ -dimensional column vector of factors, with the first d coordinates, denoted by $z_{1:d}(t)$, being curve factors and $z_{d+1}(t)$ being the

¹⁵In the language of Section 11.2.3, the stochastic volatility driver is here evidently of the *spanned* type. See also Section 12.3.6.

Fig. 12.3. Implied Volatility Smile



Notes: Fit of a four-factor quadratic Gaussian model to the volatility smile of 10y \times 10y swaptions as observed in the summer of 2007. The swaption strike (“Strike Offset”) is set as an offset to the forward swap rate.

single volatility factor. Let the $(d + 1) \times (d + 1)$ matrix-valued function $g(t)$ in (12.66) be of the block form,

$$g(t) = \begin{pmatrix} g_{1:d}(t) & 0 \\ 0 & g_{d+1,d+1}(t) \end{pmatrix},$$

with $g_{1:d}(t)$ being a $d \times d$ diffusion matrix for curve factors, and $g_{d+1,d+1}(t)$ being a scalar diffusion coefficient for the volatility factor $z_{d+1}(t)$. We write (12.66) as

$$dz(t) = g(t)^\top dW(t),$$

where

$$W(t) = \begin{pmatrix} W_{1:d}(t) \\ W_{d+1}(t) \end{pmatrix}$$

is a $(d + 1)$ -dimensional Brownian motion in the risk-neutral measure Q . Notice that we assume that the volatility factor is independent of the curve factors. In (12.65), let the linear term $h(t)$, a $(d + 1)$ -dimensional column vector, have a last element 0 so that the volatility factor has no first-order effect on the short rate,

$$h(t) = \begin{pmatrix} h_{1:d}(t) \\ 0 \end{pmatrix}. \quad (12.80)$$

Finally, in (12.65) the quadratic term $\gamma(t)$ is specified to have the block form,

$$\gamma(t) = \begin{pmatrix} \gamma_{1:d}(t) & \gamma_{1:d,d+1}(t) \\ \gamma_{1:d,d+1}(t)^\top & 0 \end{pmatrix}, \quad (12.81)$$

where the $d \times d$ matrix $\gamma_{1:d}(t)$ is given by

$$\gamma_{1:d}(t) = \eta \varpi h_{1:d}(t) h_{1:d}(t)^\top,$$

and the $d \times 1$ vector $\gamma_{1:d,d+1}(t)$ is given by

$$\gamma_{1:d,d+1}(t) = \frac{1}{2} \eta \overline{\varpi} h_{1:d}(t), \quad \overline{\varpi} = \sqrt{1 - \varpi^2},$$

where $\varpi \in [-1, 1]$.

Combining everything above, our model for the short rate is given by

$$r(t) = (1 + \eta (\varpi \times (h(t)^\top z(t)) + \overline{\varpi} \times z_{d+1}(t))) \times (h(t)^\top z(t)) + a(t), \quad (12.82)$$

and we obtain a representation of the short rate as a curve factor times one plus a volatility process, similar to (12.79).

12.3.3.3 Linear Term

To understand the volatility structure of the QG model better, let us momentarily set $\eta = 0$ in (12.82). The model then reduces to a multi-factor (linear) Gaussian model,

$$r(t) = h(t)^\top z(t) + a(t), \quad (12.83)$$

for which we can use the tools and intuition we developed earlier in this chapter. In particular, the ideas of Section 12.1.7 could be fruitfully applied. Without repeating ourselves, we assume that d benchmark rates are specified, and the volatility structure is parameterized by their instantaneous volatilities $\lambda_i(t)$, $i = 1, \dots, d$, and instantaneous correlations $\{\chi_{i,j}(t)\}$, $i, j = 1, \dots, d$. Assuming, for concreteness, that the loadings vector $h(t)$ is specified through a series of d constant mean reverrections,

$$h(t) = (e^{-\kappa_1 t}, \dots, e^{-\kappa_d t}, 0)^\top, \quad (12.84)$$

we see that the ‘‘curve’’ part of the diffusion coefficient, i.e. the matrix $g_{1:d}(t)$, can be obtained by the algorithm from Section 12.1.7. Together with the quadratic form parameterization in (12.81), this completely specifies the model. Of course, it still remains to set the various model parameters in such a way that market prices for interest rate options are matched, a topic that we turn to next.

12.3.4 Swaption Pricing

12.3.4.1 State Vector Distribution Under the Annuity Measure

Adopting the notation employed in Section 12.1.6, we consider a swaption that fixes at $T_0 > 0$ and has fixed payments at $T_1 < \dots < T_N$, $\tau_i = T_{i+1} - T_i$. We remind the reader that the relevant swap rate and the annuity are given by

$$S(t) = \frac{P(t, T_0) - P(t, T_N)}{A(t)}, \quad A(t) = \sum_{i=0}^{N-1} \tau_i P(t, T_{i+1}). \quad (12.85)$$

The corresponding annuity measure is the measure Q^A for which $A(t)$ is the numeraire. As exploited on numerous occasions already, in Q^A the swap rate $S(t)$ is a martingale and the swaption price may be obtained as a European option on $S(T_0)$.

The term distribution of the state vector on the fixing date T_0 of the swap rate, $z(T_0)$, is easily characterized.

Lemma 12.3.6. *The distribution of $z(T_0)$ in the annuity measure Q^A is a Gaussian mixture, with the density $\psi^A(z)$ of $z(T_0)$ given by*

$$\begin{aligned} \psi^A(z) &= \sum_{i=0}^{N-1} w_i^A \phi\left(z; m^{T_{i+1}}(0, T_0, 0), \nu^{T_{i+1}}(0, T_0, 0)\right), \\ w_i^A &= \frac{\tau_i P(0, T_{i+1})}{A(0)}, \quad i = 0, \dots, N-1, \end{aligned} \quad (12.86)$$

where $\phi(z; m, \nu)$ is a $(d+1)$ -dimensional Gaussian density with mean m and covariance matrix ν . The mean $m^{T_{i+1}}(0, T_0, 0)$ and the covariance matrix $\nu^{T_{i+1}}(0, T_0, 0)$ are given in Proposition 12.3.4.

Proof. For an arbitrary scalar function $g(z)$ we have, from standard measure change arguments,

$$\begin{aligned} \mathbb{E}^A(g(z(T_0))) &= \frac{P(0, T_0)}{A(0)} \mathbb{E}^T(g(z(T_0)) A(T_0)) \\ &= \frac{P(0, T_0)}{A(0)} \sum_{i=0}^{N-1} \mathbb{E}^T(g(z(T_0)) \tau_i P(T_0, T_{i+1})) \\ &= \sum_{i=0}^{N-1} \frac{\tau_i P(0, T_{i+1})}{A(0)} \mathbb{E}^{T_{i+1}}(g(z(T_0))). \end{aligned}$$

The lemma follows directly from Proposition 12.3.4. \square

It follows from Lemma 12.3.6 that the mean and the covariance matrix of $z(T_0)$ in the annuity measure Q^A are given by

$$\begin{aligned} \mathbb{E}^A(z(T_0)) &= \sum_{i=0}^{N-1} w_i^A m^{T_{i+1}}(0, T_0, 0), \\ \text{Var}^A(z(T_0)) &= \sum_{i=0}^{N-1} w_i^A (\nu^{T_{i+1}}(0, T_0, 0) + m^{T_{i+1}}(0, T_0, 0)^2) - (\mathbb{E}^A(z(T_0)))^2. \end{aligned} \quad (12.87)$$

We emphasize that $z(T_0)$ is not Gaussian under the annuity measure, although it is tempting to use the approximation

$$z(T_0) \stackrel{d}{\approx} \mathcal{N}(\mathbb{E}^A(z(T_0)), \text{Var}^A(z(T_0))). \quad (12.88)$$

In practical uses of either (12.86) or (12.87), we need an efficient way to compute the moments $m^{T_{i+1}}(0, T_0, 0)$ and $\nu^{T_{i+1}}(0, T_0, 0)$ of $z(T_0)$. The results in Proposition 12.3.4 require an integration in the time domain for each $i = 0, \dots, N-1$, and, with N being potentially large, the computational effort would therefore often be quite high. Fortunately, a much faster alternative is available. Once $m^{T_0}(0, T_0, 0)$ and $\nu^{T_0}(0, T_0, 0)$ are on hand — and they are always known as both are required for yield curve fitting via (12.69) — the moments of the factors observed on the same date, but under all other forward measures, can be calculated by a measure-change formula:

$$m^{T_{i+1}}(0, T_0, 0) = \mathbb{E}^{T_{i+1}}(z(T_0)) = \frac{\mathbb{E}^{T_0}(z(T_0)P(T_0, T_{i+1}))}{\mathbb{E}^{T_0}(P(T_0, T_{i+1}))}. \quad (12.89)$$

The expression on the right-hand side is obtained in closed form in Corollary 12.A.3 of Appendix 12.A, utilizing the fact that the discount bond is an exponential of a quadratic form of a Gaussian vector.

12.3.4.2 Exact Pricing of European Swaptions

The result of Lemma 12.3.6 that shows that the distribution of the state vector in the annuity measure is a mixture of Gaussian distributions leads us to one possible European swaption pricing method. To elaborate, let us define $S(T_0, z)$ to be the value of the swap rate $S(T_0)$ when $z(T_0) = z$. It follows from (12.86) that the value of a European swaption with strike c can be represented as

$$V_{\text{swaption}}(0) = A(0)\mathbb{E}^A(S(T_0, \mu_\xi + \sigma_\xi X) - c)^+, \quad (12.90)$$

where ξ is an integer-valued random variable with distribution

$$\mathbb{Q}(\xi = i) = w_i^A, \quad i = 0, \dots, N-1,$$

X is a standard Gaussian $(d+1)$ -dimensional random variable, and

$$\mu_i = m^{T_{i+1}}(0, T_0, 0), \quad \sigma_i = \sqrt{\nu^{T_{i+1}}(0, T_0, 0)}, \quad i = 0, \dots, N-1.$$

Evaluation of (12.90) may, for instance, proceed by a simple one-step Monte Carlo simulation

$$\mathbb{E}^A (S(T_0) - c)^+ \approx \frac{1}{L} \sum_{l=1}^L (S(T_0, \mu_{\xi_l} + \sigma_{\xi_l} X_l) - c)^+,$$

where $\{\xi_l\}_{l=1}^L$ is an i.i.d. sample from the distribution of ξ , and $\{X_l\}_{l=1}^L$ is an i.i.d. sample from the standard $(d+1)$ -dimensional Gaussian distribution.

While this scheme still requires a Monte Carlo simulation to compute a swaption value, we emphasize that it is simple and fast: one only needs to draw a sample of the variable ξ (which essentially defines what mean/variance to use) and one sample of a standard Gaussian vector, in order to sample *directly* the terminal distribution of the swap rate. When combined with quasi-random numbers as in Section 3.2.10.1, the method could be seen as a type of outright $(d+1)$ -dimensional numerical integration.

12.3.4.3 Approximations for European Swaptions

As the short rate and all continuously compounded forward rates are quadratic forms of the state vector, it seems reasonable to assume that swap rates are approximately of this form as well. We can use this observation to develop analytically tractable approximations for swaptions. Let us define

$$h_S = \nabla^\top S(T_0, z^*), \quad \gamma_S = \frac{1}{2} \nabla^\top \nabla S(T_0, z^*), \quad (12.91)$$

where ∇ is the gradient operator, $\nabla = (\partial/\partial z_1, \dots, \partial/\partial z_{d+1})$ (row vector). In essence, h_S is the first, and γ_S is the (half of the) second-order derivative of the swap rate function $S(T_0, \cdot)$ at a specific point z^* . Both are easily computed by a numerical finite difference algorithm. The expansion point z^* could be 0 or, for a slightly more accurate approximation,

$$z^* = \mathbb{E}^A(z(T_0))$$

as computed by (12.87). Applying Taylor's expansion,

$$S(T_0, z) \approx (z - z^*)^\top \gamma_S (z - z^*) + h_S^\top (z - z^*) + s(z^*).$$

To ensure that the forward swap rate is repriced correctly under this approximation, we adjust the constant term accordingly, and define the *quadratic approximation to the swap rate* by

$$\begin{aligned} S(T_0, z) &\approx S_q(T_0, z), \\ S_q(T_0, z) &= z^\top \gamma_S z + h_S^\top z - \mathbb{E}^A (z(T_0)^\top \gamma_S z(T_0) + h_S^\top z(T_0)) + S(0), \end{aligned} \quad (12.92)$$

where the required expected value is calculated in Corollary 12.A.1.

Under the quadratic approximation to the swap rate and Gaussian approximation to the distribution of $z(T_0)$ (see (12.88)), it becomes possible to price options on the swap rate using Fourier integration methods. For this, we need the moment-generating function

$$q(u) \triangleq \tilde{E} \left(e^{uS_q(T_0, z(T_0))} \right),$$

where \tilde{E} is the expected value operator under the assumption that $z(T_0)$ is Gaussian. This expression is indeed available in closed form, thanks to Proposition 12.3.5:

$$q(u) = \exp \left(\Psi \left(uh_S, u\gamma_S; E^A(z(T_0)), \text{Var}^A(z(T_0)) \right) \right).$$

Given $q(u)$, we can compute the option price

$$\tilde{E} \left((S_q(T_0, z(T_0)) - c)^+ \right)$$

by Theorem 8.4.3. A suitable control variate as in Theorem 8.4.4 is essential for improving numerical performance.

It should be noted here that the application of Fourier methods to swaption pricing does not hinge on the Gaussian approximation (12.88), as the true moment-generating function of $S_q(T_0, z(T_0))$ is readily available under Q^A . Indeed, from the mixing formula (12.86) we get

$$\begin{aligned} E^A \left(e^{uS_q(T_0, z(T_0))} \right) &= \sum_{i=0}^{N-1} w_i^A E^{T_{i+1}} \left(e^{uS_q(T_0, z(T_0))} \right) \\ &= \sum_{i=0}^{N-1} w_i^A \exp \left(\Psi \left(uh_S, u\gamma_S; m^{T_{i+1}}(0, T_0, 0), \nu^{T_{i+1}}(0, T_0, 0) \right) \right). \end{aligned}$$

We can use this formula in option pricing instead of the Gaussian approximation;; however, we find that the resulting increase in computational cost — for each value of u we now require N evaluations of the function Ψ , instead of just one — is rarely justified by the (slight) improvements of accuracy.

While we find the Fourier integration method to be robust and efficient¹⁶, we can further explore the specifics of our parameterization to design even faster valuation algorithms. In particular, from (12.82) we notice that the quadratic form defining the short rate is not of full rank, as the short rate is a quadratic function of only two “aggregate” quantities, $h(t)^\top z(t)$ and $z_{d+1}(t)$. We can expect that this rank-2 structure is preserved, at least approximately, in swap rates. The linear term in the quadratic approximation for the swap rate, $h_S^\top z$, will be one of the two aggregate factors to use in re-parameterizing

¹⁶Contrary to some claims in the literature, see e.g. Boyarchenko and Levendorski [2007].

the quadratic part $z^\top \gamma_S z$. The other one, naturally, will be the stochastic volatility factor. In summary, we seek to approximate

$$z^\top \gamma_S z \approx (h_S^\top z, z_{d+1}) \hat{\gamma}_S (h_S^\top z, z_{d+1})^\top, \quad (12.93)$$

where $\hat{\gamma}_S$ is a 2×2 (symmetric) matrix.

To formalize the idea outlined above, let us define a two-dimensional stochastic vector $\hat{z}(T_0) = (\hat{z}_1(T_0), \hat{z}_2(T_0))^\top$ by

$$\hat{z}_1(T_0) = h_S^\top z(T_0), \quad \hat{z}_2(T_0) = z_{d+1}(T_0),$$

or, in matrix notation,

$$\begin{pmatrix} \hat{z}_1(T_0) \\ \hat{z}_2(T_0) \end{pmatrix} = Rz(T_0), \quad R = \begin{pmatrix} h_{S,1} & \dots & h_{S,d} & h_{S,d+1} \\ 0 & \dots & 0 & 1 \end{pmatrix}.$$

Then (12.93) can be re-written as

$$z^\top \gamma_S z \approx z^\top R^\top \hat{\gamma}_S R z.$$

Formally, we set $\hat{\gamma}_S$ to be a solution to the following minimization problem

$$\text{Var}^A (z(T_0)^\top (\gamma_S - R^\top \hat{\gamma}_S R) z(T_0)) \rightarrow \min; \quad (12.94)$$

we then call this $\hat{\gamma}_S$ a *rank-2 quadratic approximation*. The problem (12.94) can be solved explicitly (if rather tediously) using Corollary 12.A.1, resulting in the following approximation to the value of an option on the swap rate.

Theorem 12.3.7. *Under the rank-2 quadratic approximation (12.94) to the swap rate defined by (12.85), and Gaussian approximation to the distribution of the state vector $z(T_0)$ under Q^A , the value of a European swaption with strike c is approximately given by*

$$V_{\text{swaption}}(0) \approx A(0) \int_{\mathbb{R}^2} (\hat{z}^\top \hat{\gamma}_S \hat{z} + \hat{z}_1 + \hat{\alpha}_S + S(0) - c)^+ \phi(\hat{z}; \hat{m}, \hat{\nu}) d\hat{z}, \quad (12.95)$$

where $\hat{z} = (\hat{z}_1, \hat{z}_2)^\top$ and $\phi(\hat{z}; \hat{m}, \hat{\nu})$ a two-dimensional Gaussian density with mean \hat{m} and covariance matrix $\hat{\nu}$. Also, $\hat{\alpha}_S$ is defined by

$$\hat{\alpha}_S = -E^A (z(T_0)^\top R^\top \hat{\gamma}_S R z(T_0) + h_S^\top z(T_0)), \quad (12.96)$$

with

$$\hat{\gamma}_S = (2\hat{m}\hat{m}^\top + \hat{\nu})^{-1} R (2m m^\top + \nu) \gamma_S \nu R^\top \hat{\nu}^{-1},$$

where

$$m = E^A (z(T_0)), \quad \nu = \text{Var}^A (z(T_0)), \quad \hat{m} = Rm, \quad \hat{\nu} = R\nu R^\top,$$

As in Section 12.1.6.1, the two-dimensional integral in (12.95) can be computed efficiently by conditioning on one of the integration variables, evaluating the resulting sub-expression in closed form, and performing the outer integration using, say, *Gauss-Hermite* quadrature (see Press et al. [1992]). We omit straightforward details. Table 12.1 demonstrates typical quality of various approximations. Data in the table represent $10y \times 10y$ swaption volatilities, computed from the same model settings as those used to construct Figure 12.3.

Strike	ATM-2%	ATM-1%	ATM	ATM+1%	ATM+2%
Model exact	15.84	13.17	11.54	11.12	11.09
Gauss approx	15.84	13.17	11.55	11.12	11.09
Gauss+Quadratic	15.89	13.17	11.51	11.06	11.00
Gauss+Quadratic+Rank 2	15.89	13.17	11.51	11.07	10.99

Table 12.1. Implied Black Volatilities in a Quadratic Gaussian Model. “Exact” is defined by (12.90). “Gauss approx” is defined by (12.88). “Quadratic” is defined by (12.92). “Rank 2” is defined by (12.95). Results for a $10y \times 10y$ swaption in %.

12.3.5 Calibration

While a number of viable approaches to calibration exist, we recommend organizing it as a multi-pass bootstrap algorithm, an approach that should be familiar to the reader by now (see e.g. Section 10.2.5.2). First, the parameters ϖ and η are fixed to the desired shape of volatility smile. Next, the correlation matrix of the benchmark rates $\{\chi_{i,j}(t)\}$ is parameterized by a convenient functional form (see the discussion in Section 14.3.2), generally to either match historical correlations of the relevant rates or to fit market-implied prices of CMS spread options. After that, the calibration problem is reduced to the problem of matching at-the-money swaption volatilities by manipulating the benchmark rate volatilities $\lambda_i(t)$, $i = 1, \dots, d$ (the reader will recall from Sections 12.3.3.3 and 12.1.7 that they are used to construct the state vector diffusion matrix $g(t)$). Having d time-dependent volatilities allows us to calibrate to d swaption strips. While not strictly necessary, we find it convenient to choose the swaption strips to be of constant tenor, with the tenors matching those of the benchmark rates. Denoting $t_1 < \dots < t_K$ to be the expiry dates of the swaptions in the calibration set, we break the calibration into K subproblems, where in the j -th sub-calibration we match the j -th row of the swaption matrix by tweaking $\lambda_i(t_j)$, $i = 1, \dots, d$. In the linear case, i.e. when the quadratic term is zero, only one pass for $j = 1, \dots, K$ is required, as swaption prices with expiry t_j depend on $\lambda(s)$ for $s \in [0, t_j]$ only. In the general quadratic case, this is no longer the case, and prices of swaptions with expiry t_j depend on $\lambda(s)$ for all $s > t_j$.

through bond reconstruction formulas). However, this “tail” dependence is minor, and we can still calibrate sequentially by performing multiple passes (typically two or three). For a more performant algorithm, we could use fast swaption approximations for initial pass(es), saving a more accurate one for the final pass. The specifics of such a multi-pass calibration should follow closely the ideas discussed in more details in the context of affine models in Section 10.2.5.

12.3.6 Spanned Stochastic Volatility

While we have used the term “stochastic volatility” throughout to describe our parameterization of the QG model, the model clearly does not involve true unspanned stochastic volatility, of the type defined in Section 11.2.3. In particular, the discount bond reconstitution formulas in Proposition 12.3.1 depend on the full vector of state variables $z(t)$. However, in parameterizing the model we were careful to assign zero weight to z_{d+1} , the “volatility” factor, in the linear part of the quadratic form for the short rate (see (12.80)), ensuring that discount bonds have rather small (second order) dependence on it. Hence, we expect the model to exhibit some traits of stochastic volatility models. Lending some credibility to this observation, Piterbarg [2009a] (and, with more details, Piterbarg [2008]) analyzes the dynamics of volatility smiles in two-factor quadratic Gaussian models and concludes that these models lie somewhere between local volatility and true stochastic volatility models (which we introduce in Chapter 13 below).

12.3.7 Numerical Methods

We round out our discussion of quadratic Gaussian models with a quick review of numerical methods available for derivatives pricing. The discussion is mercifully brief because the state variables in a quadratic Gaussian model follow the same process as the state variables in a linear Gaussian model, making the material of Section 12.1.8 directly applicable. This fortunate circumstance is, in fact, one of the key attractions of the quadratic Gaussian models, as we mentioned earlier. For instance, PDE methods for the quadratic Gaussian model carry over unchanged from the linear Gaussian case, as the state variables in both classes of models follow essentially identical processes. We refer to Section 12.1.9 for details.

As for Monte Carlo simulation, we can reuse results from Section 12.1.8, and emphasize that the state vector can be simulated at low cost and bias-free over a period of time of any length without adding any intermediate dates. As a result, the performance of the Monte Carlo method for the quadratic Gaussian model is on par with the linear Gaussian model, and far ahead of any alternative multi-factor model with volatility smile, such as the Libor market model (Chapter 14) or even the multi-factor quasi-Gaussian model (Chapter 13).

12.A Appendix: Quadratic Forms of Gaussian Vectors

First, we prove Proposition 12.3.5. We have,

$$\begin{aligned}\Psi(u, Q; m, \nu) &= \frac{1}{(2\pi)^{K/2} \sqrt{\det(\nu)}} \\ &\times \int \exp(z^\top Qz + u^\top z) \exp\left(-\frac{1}{2}(z - m)^\top \nu^{-1}(z - m)\right) dz.\end{aligned}$$

We have,

$$z^\top Qz + u^\top z = \left((z - m)^\top Q(z - m) + 2m^\top Qz - m^\top Qm + u^\top z\right),$$

so the integrand is equal to

$$\exp\left(-\frac{1}{2}(z - m)^\top (\nu^{-1} - 2Q)(z - m)\right) \exp(2m^\top Qz - m^\top Qm + u^\top z).$$

Define

$$\nu_Q \triangleq (\nu^{-1} - 2Q)^{-1} = \nu(I - 2Q\nu)^{-1}.$$

Let Q^{ν_Q} be a measure under which Z has mean m and variance ν_Q , and E^{ν_Q} the corresponding expected value operator. Then

$$\Psi(u, Q; m, \nu) = \exp(-m^\top Qm) \frac{\sqrt{\det(\nu_Q)}}{\sqrt{\det(\nu)}} E^{\nu_Q}(\exp(2m^\top QZ + u^\top Z)).$$

By the standard results for exponents of Gaussian linear forms (see e.g. Kotz et al. [2000]),

$$\begin{aligned}E^{\nu_Q}(\exp((2m^\top Q + u^\top)Z)) \\ = \exp\left((2m^\top Q + u^\top)m + \frac{1}{2}(2m^\top Q + u^\top)\nu(I - 2Q\nu)^{-1}(2Qm + u)\right).\end{aligned}$$

Thus we get

$$\begin{aligned}\ln E(\exp(Z^\top QZ + u^\top Z)) &= \frac{1}{2}(2m^\top Q + u^\top)\nu(I - 2Q\nu)^{-1}(2Qm + u) \\ &\quad + m^\top Qm + u^\top m + \frac{1}{2} \ln \det(\nu^{-1}\nu_Q).\end{aligned}$$

The proposition has been proven.

Once the moment-generating function is available, other characteristics of the distribution follow. For example, we can easily calculate the mean and the variance of a quadratic form of a Gaussian vector.

Corollary 12.A.1. Let Z be a K -dimensional Gaussian vector with mean m and variance ν . Let Q be a symmetric $K \times K$ matrix and u a K -dimensional vector. Then

$$\begin{aligned}\mathbb{E}(Z^\top QZ + u^\top Z) &= (m^\top Qm + u^\top m) + \text{tr}(Q\nu), \\ \text{Var}(Z^\top QZ + u^\top Z) &= (2m^\top Q + u^\top) \nu (2Qm + u) + 2\text{tr}(Q\nu Q\nu).\end{aligned}$$

Proof. Clearly,

$$\begin{aligned}\mathbb{E}(Z^\top QZ + u^\top Z) &= \frac{d}{d\epsilon} \Psi(\epsilon u, \epsilon Q; m, \nu) \Big|_{\epsilon=0}, \\ \text{Var}(Z^\top QZ + u^\top Z) &= \frac{d^2}{d\epsilon^2} \Psi(\epsilon u, \epsilon Q; m, \nu) \Big|_{\epsilon=0}.\end{aligned}$$

Recall

$$\begin{aligned}\Psi(\epsilon u, \epsilon Q; m, \nu) &= \frac{\epsilon^2}{2} (2m^\top Q + u^\top) (\nu^{-1} - 2\epsilon Q)^{-1} (2Qm + u) \\ &\quad + \epsilon (m^\top Qm + u^\top m) - \frac{1}{2} \ln \det(I - 2\epsilon Q\nu).\end{aligned}$$

By Jacobi's formula,

$$\frac{d}{d\epsilon} \det(I - 2\epsilon Q\nu) = -\det(I - 2\epsilon Q\nu) \text{tr}\left(2(I - 2\epsilon Q\nu)^{-1} Q\nu\right),$$

so

$$\frac{d}{d\epsilon} \ln \det(I - 2\epsilon Q\nu) = -\text{tr}\left(2(I - 2\epsilon Q\nu)^{-1} Q\nu\right).$$

Then

$$\begin{aligned}\frac{d}{d\epsilon} \Psi(\epsilon u, \epsilon Q; m, \nu) &= \epsilon (2m^\top Q + u^\top) (\nu^{-1} - 2Q\epsilon)^{-1} (2Qm + u) \\ &\quad + \frac{1}{2} \epsilon^2 \times (\dots) \\ &\quad + (m^\top Qm + u^\top m) + \text{tr}\left((I - 2\epsilon Q\nu)^{-1} Q\nu\right),\end{aligned}$$

so that

$$\frac{d}{d\epsilon} \Psi(\epsilon u, \epsilon Q; m, \nu) \Big|_{\epsilon=0} = (m^\top Qm + u^\top m) + \text{tr}(Q\nu).$$

Furthermore,

$$\begin{aligned}\frac{d}{d\epsilon} \text{tr}\left((I - 2\epsilon Q\nu)^{-1} Q\nu\right) &= \text{tr}\left(\frac{d}{d\epsilon} \left((I - 2\epsilon Q\nu)^{-1} Q\nu\right)\right) \\ &= 2\text{tr}\left((I - 2\epsilon Q\nu)^{-1} Q\nu (I - 2\epsilon Q\nu)^{-1} Q\nu\right),\end{aligned}$$

So,

$$\begin{aligned} \frac{d^2}{d\epsilon^2} \Psi(\epsilon u, \epsilon Q; m, \nu) &= (2m^\top Q + u^\top) (\nu^{-1} - 2\epsilon Q)^{-1} (2Qm + u) \\ &\quad + \epsilon \times (\dots) \\ &\quad + 2\text{tr} \left((I - 2\epsilon Q\nu)^{-1} Q\nu (I - 2\epsilon Q\nu)^{-1} Q\nu \right), \end{aligned}$$

thus

$$\left. \frac{d^2}{d\epsilon^2} \Psi(\epsilon u, \epsilon Q; m, \nu) \right|_{\epsilon=0} = (2m^\top Q + u^\top) \nu (2Qm + u) + 2\text{tr}(Q\nu Q\nu).$$

□

Interestingly, we can obtain covariances or, indeed, any cross-moments of multiple quadratic forms of the same vector Z using the same idea as in the previous corollary.

Corollary 12.A.2. *Let Z be a K -dimensional Gaussian vector with mean m and variance ν . Let Q_1, Q_2 be symmetric $K \times K$ matrices and u_1, u_2 be K -dimensional vectors. Then*

$$\begin{aligned} \mathbb{E} \left((Z^\top Q_1 Z + u_1^\top Z)^n (Z^\top Q_2 Z + u_2^\top Z)^m \right) \\ = \left. \frac{\partial^{n+m}}{\partial \epsilon_1^n \partial \epsilon_2^m} \exp(\Psi(\epsilon_1 u_1 + \epsilon_2 u_2, \epsilon_1 Q_1 + \epsilon_2 Q_2; m, \nu)) \right|_{\epsilon_1=\epsilon_2=0}. \end{aligned}$$

In particular,

$$\begin{aligned} \text{Cov}(Z^\top Q_1 Z + u_1^\top Z, Z^\top Q_2 Z + u_2^\top Z) \\ = (2m^\top Q_1 + u_1^\top) \nu (2Q_2 m + u_2) + 2\text{tr}(Q_1 \nu Q_2 \nu). \end{aligned}$$

The next corollary helps with calculating moments of the state vector under different forward measures in quadratic Gaussian models, see (12.89).

Corollary 12.A.3. *Let Z be a K -dimensional Gaussian vector with mean m and variance ν . Let Q be a symmetric $K \times K$ matrix and u a K -dimensional vector. Denote*

$$\begin{aligned} \hat{m} &= \frac{\mathbb{E}(Z \exp(-(Z^\top Q Z + u^\top Z)))}{\mathbb{E}(\exp(-(Z^\top Q Z + u^\top Z)))}, \\ \hat{\nu} &= \frac{\mathbb{E}(Z Z^\top \exp(-(Z^\top Q Z + u^\top Z)))}{\mathbb{E}(\exp(-(Z^\top Q Z + u^\top Z)))} - \hat{m} \hat{m}^\top. \end{aligned}$$

Then

$$\hat{m} = m - \nu(I + 2Q\nu)^{-1}(2Qm + u), \quad (12.97)$$

$$\hat{\nu} = \nu(I + 2Q\nu)^{-1}. \quad (12.98)$$

Proof. First, we note that

$$\begin{aligned}\widehat{m} &= -\frac{d}{du} \Psi(-u, -Q; m, \nu), \\ \widehat{\nu} &= -\frac{d}{dQ} \Psi(-u, -Q; m, \nu) - \widehat{m} \widehat{m}^\top.\end{aligned}$$

From Proposition 12.3.5,

$$\begin{aligned}\Psi(-u, -Q; m, \nu) &= \frac{1}{2} (2m^\top Q + u^\top) (\nu^{-1} + 2Q)^{-1} (2Qm + u) \\ &\quad - m^\top Qm - u^\top m - \frac{1}{2} \ln \det(I + 2Q\nu).\end{aligned}$$

Then

$$\begin{aligned}\frac{d}{du} \left(\frac{1}{2} (2m^\top Q + u^\top) \nu (I + 2Q\nu)^{-1} (2Qm + u) - m^\top Qm - u^\top m \right) \\ = \nu (I + 2Q\nu)^{-1} (2Qm + u) - m,\end{aligned}$$

and (12.97) follows.

The proof of (12.98) proceeds along similar lines, using the fact that for any matrix A

$$\frac{d}{dA} \det A = (\det A) (A^{-1})^\top$$

and, in particular,

$$\frac{d}{dQ} \det(I + 2Q\nu) = 2 \det(I + 2Q\nu) \nu (I + 2Q\nu)^{-1}.$$

□

The Quasi-Gaussian Model with Local and Stochastic Volatility

In this chapter we consider extensions to one- and multi-factor Gaussian short rate models (Chapters 10, 11 and 12) with local and stochastic volatility. The extensions come at additional computational cost, as extra state variables are required to preserve the Markovian structure of the model. Following the pioneering work of Jamshidian [1991b], we use the term *quasi-Gaussian*¹ for the models in this chapter; their development for practical applications was undertaken in Andreasen [2001], Andersen and Andreasen [2002] and Andreasen [2005], building on early work by Jamshidian [1991b], Babbs [1990], Cheyette [1991] and Ritchken and Sankarasubramanian [1995]. Low-dimensional versions of quasi-Gaussian models are, in our opinion, among the best — if not *the* best — low-factor short rate models, as they combine flexibility of volatility smile specification, relative ease of calibration, and efficient numerical implementation. Higher-dimensional quasi-Gaussian models, while not yet mainstream, provide an alluring alternative to the better-established Libor market models (see Chapter 14).

We start this chapter by developing a one-factor quasi-Gaussian model with a local volatility function. The problems of volatility and mean reversion calibration are given considerable attention, and are followed by a discussion of various numerical methods used for model implementation. A straightforward extension to stochastic volatility is presented next, followed by development of multi-factor quasi-Gaussian models.

13.1 One-Factor Quasi-Gaussian Model

13.1.1 Definition

Recall that any HJM model is defined by a volatility structure of instantaneous forward rates. In particular, for any “reasonable” random function $\sigma_f(t, T) = \sigma_f(t, T, \omega)$, the following SDE defines a valid HJM model,

¹Also known as *pseudo-Gaussian* or *Cheyette* models.

$$df(t, T) = \sigma_f(t, T) \left(\left(\int_t^T \sigma_f(t, u) du \right) dt + dW(t) \right), \quad 0 \leq t \leq T < \infty. \quad (13.1)$$

Here $W(t)$ is a one-dimensional Brownian motion in the risk-neutral measure, and $\{f(t, T)\}_{T \geq t}$ is a collection of instantaneous forward rates.

It is shown in Section 4.5.2 that a one-factor Markovian Gaussian model is obtained by imposing a separability condition on the deterministic volatility structure of instantaneous forward rates, see (4.44). A general class of one-factor *quasi-Gaussian* (qG)² models is obtained by retaining the separability condition, but relaxing the deterministic requirement in a specific way. In particular, the component of the volatility structure that is a function of calendar time (the function g), is now allowed to be stochastic:

$$\sigma_f(t, T, \omega) = g(t, \omega) h(T). \quad (13.2)$$

In line with the notations of Section 10.1.2.2, we define

$$\begin{aligned} \kappa(t) &= -\frac{h'(t)}{h(t)}, \\ G(t, T) &= \frac{\int_t^T h(s) ds}{h(t)}, \\ \sigma_r(t, \omega) &= \sigma_f(t, t, \omega) = g(t, \omega) h(t). \end{aligned} \quad (13.3)$$

The proof of Proposition 10.1.7 carries through unchanged even for stochastic $g(t, \omega)$, and we obtain the following result.

Proposition 13.1.1. *Consider a general one-factor qG model, i.e. the HJM model (13.1) with the separable volatility condition (13.2). Define stochastic processes $x(t)$ and $y(t)$ by*

$$\begin{aligned} dx(t) &= (y(t) - \kappa(t)x(t)) dt + \sigma_r(t, \omega) dW(t), \\ dy(t) &= (\sigma_r(t, \omega)^2 - 2\kappa(t)y(t)) dt, \\ x(0) &= y(0) = 0. \end{aligned} \quad (13.4)$$

In the general qG model all zero-coupon discount bonds are deterministic functions of the processes $x(t)$ and $y(t)$,

$$P(t, T) = P(t, T, x(t), y(t)),$$

where

$$P(t, T, x, y) = \frac{P(0, T)}{P(0, t)} \exp \left(-G(t, T)x - \frac{1}{2}G(t, T)^2 y \right), \quad (13.5)$$

²We use a small q in the abbreviation qG to avoid notational conflict with the quadratic Gaussian (QG) models of Chapter 12.

the instantaneous forward rates are given by

$$f(t, T) = f(0, T) + \frac{h(T)}{h(t)} (x(t) + G(t, T) y(t)), \quad (13.6)$$

and the short rate is

$$r(t) = f(t, t) = f(0, t) + x(t). \quad (13.7)$$

The proposition demonstrates that the evolution of the whole interest rate curve, as parameterized by either forward rates or discount bonds, in the model can be reduced to the evolution of just two *state variables* $x(t)$ and $y(t)$, with dynamics given by (13.4). Unlike many of the models in Chapter 11, the qG model has a closed-form bond reconstitution formula for arbitrary choices of $g(t, \omega)$; this tractability comes at the cost at requiring *two* state variables (x and y), even though the Brownian motion $W(t)$ is only one-dimensional. Observe that in general, the function $y(t)$ is *not* deterministic, except in the case of pure Gaussian dynamics, i.e. when $\sigma_r(t, \omega)$ is a deterministic function of t . However, even when it is not deterministic, $y(t)$ does not have the diffusion term; we call such processes *locally deterministic*.

The roles of the two state variables $x(t)$ and $y(t)$ in the qG model are rather different. The variable $x(t)$ constitutes the main yield curve driver, as evidenced in (13.7), whereas $y(t)$ is an auxiliary “convexity” variable required to uphold the no-arbitrage condition; in general, it is convenient to think of the model as having “one and a half” factors.

13.1.2 Local Volatility

A one-factor qG model with *local volatility* is obtained by requiring $g(\cdot)$ to be a deterministic, time-dependent function of the state variables,

$$g(t) = g(t, x(t), y(t)). \quad (13.8)$$

Then, the short rate volatility $\sigma_r(\cdot)$ is also a function of the state variables,

$$\sigma_r(t) = \sigma_r(t, x(t), y(t)) = g(t, x(t), y(t)) h(t), \quad (13.9)$$

and the dynamics of the state variables in the local volatility qG model are given by

$$dx(t) = (y(t) - \kappa(t)x(t)) dt + \sigma_r(t, x(t), y(t)) dW(t), \quad (13.10)$$

$$dy(t) = (\sigma_r(t, x(t), y(t))^2 - 2\kappa(t)y(t)) dt. \quad (13.11)$$

Clearly, (13.10)–(13.11) define a two-dimensional Markovian process. As all zero-coupon discount bonds are functions of these two state variables by Proposition 13.1.1, the local volatility qG model is Markovian in two state variables.

For future use, we denote

$$\sigma_r^0(t) \triangleq \sigma_r(t, 0, 0). \quad (13.12)$$

If $\sigma_r(t, x, y)$ is independent of x, y , the model reduces to a purely Gaussian model with the deterministic short rate volatility $\sigma_r^0(t)$.

13.1.3 Swap Rate Dynamics

For the purpose of European swaption pricing in the qG model, we shall need to establish swap rate dynamics in an annuity measure. For this purpose, let us fix a tenor structure

$$0 < T_0 < T_1 < T_2 < \dots < T_N,$$

with

$$\tau_n = T_{n+1} - T_n.$$

Consider a forward swap rate $S(t)$ with the first fixing T_0 and the last payment T_N (see Section 4.1.3), i.e.

$$S(t) \triangleq S_{0,N}(t) = \frac{P(t, T_0) - P(t, T_N)}{A(t)}, \quad (13.13)$$

$$A(t) \triangleq A_{0,N}(t) = \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}). \quad (13.14)$$

It follows from Proposition 13.1.1 that all zero-coupon bonds $P(t, \cdot)$ are deterministic functions of $x(t)$ and $y(t)$, and hence so is the forward swap rate; accordingly we define

$$S(t, x, y) \triangleq \frac{P(t, T_0, x, y) - P(t, T_N, x, y)}{\sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}, x, y)}. \quad (13.15)$$

The following proposition, a simple extension of the results from Section 10.1.3.2, determines the dynamics of the swap rate under its corresponding annuity measure, i.e. the measure Q^A for which $A(t)$ is a numeraire.

Proposition 13.1.2. *We have*

$$dS(t) = \left(\frac{\partial S}{\partial x}(t, x(t), y(t)) \right) \sigma_r(t, x(t), y(t)) dW^A(t), \quad (13.16)$$

where $W^A(t)$ is a Brownian motion in measure Q^A . Here

$$\begin{aligned} \frac{\partial S}{\partial x}(t, x, y) &= -\frac{1}{A(t, x, y)} (P(t, T_0, x, y) G(t, T_0) - P(t, T_N, x, y) G(t, T_N)) \\ &\quad + \frac{S(t, x, y)}{A(t, x, y)} \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}, x, y) G(t, T_{n+1}). \end{aligned} \quad (13.17)$$

Proof. By definition

$$S(t) = S(t, x(t), y(t)).$$

The statement of the proposition follows by applying Ito's lemma and dropping dt terms from the expression, as justified by the fact that the swap rate is a martingale in the annuity measure, per Lemma 4.2.4. \square

13.1.4 Approximate Local Volatility Dynamics for Swap Rate

The SDE (13.16) shows that a swap rate follows a local volatility process where the volatility is a function of the short rate state x and the auxiliary variable y . Since the model is essentially one-factor, it is reasonable to assume that there is a strong linkage between a swap rate and the state of the short rate. Hence, it seems plausible that the dynamics for a swap rate could be written with a diffusion term that is just a function of the swap rate itself. Such a simplification would be convenient in many applications, as methods from Chapter 7 could be called upon to solve the resulting SDE or to price options. The following proposition proven by the methods of *Markovian projection* (see Appendix A) makes matters precise.

Lemma 13.1.3. *The values of all European options on the swap rate $S(t)$ in the model (13.16) are identical to values computed in a vanilla model with time-dependent local volatility function:*

$$dS(t) = \varphi(t, S(t)) dW^A(t), \quad (13.18)$$

where

$$\varphi(t, s)^2 = E^A \left(\left(\frac{\partial S}{\partial x}(t, x(t), y(t)) \sigma_r(t, x(t), y(t)) \right)^2 \middle| S(t) = s \right). \quad (13.19)$$

While evaluating conditional expectations such as (13.19) is often rather difficult, here we are aided considerably by the essentially one-dimensional structure of the problem. If we assume that $y(t)$ is well approximated by a deterministic function $\bar{y}(t)$ — an approximation we shall use repeatedly in this chapter — then $S(t)$ would just be a deterministic function of $x(t)$ and time t . The opposite would also be true, i.e. $x(t)$ would be a deterministic function of $S(t)$, $x(t) = X(t, S(t))$. If this function were available, the evaluation of the conditional expected value in (13.19) would boil down to evaluating $\frac{\partial S}{\partial x}(t, x(t), y(t))\sigma_r(t, x(t), y(t))$ at $x(t) = X(t, s)$ and $y(t) = \bar{y}(t)$,

$$\varphi(t, s) \approx \frac{\partial S}{\partial x}(t, X(t, s), \bar{y}(t)) \sigma_r(t, X(t, s), \bar{y}(t)),$$

where $\partial S / \partial x$ is given by (13.17).

Before presenting various methods for approximating \bar{y} and the function X , let us emphasize again that once the volatility function $\varphi(t, s)$ is determined, swaption values can be computed from (13.18) by methods developed for local volatility vanilla models in Chapter 7.

13.1.4.1 Simple Approximation

A simple approximation for the function φ is obtained by slightly extending the idea from Section 10.1.3.2. Setting

$$\bar{y}(t) = 0,$$

and applying a linear approximation

$$\begin{aligned} S(t, x, 0) &\approx S(t, 0, 0) + \frac{\partial S}{\partial x}(t, 0, 0)x, \\ \frac{\partial S}{\partial x}(t, x, 0) &\approx \frac{\partial S}{\partial x}(t, 0, 0), \end{aligned} \quad (13.20)$$

we obtain

$$x \approx \frac{S(t, x, 0) - S(t, 0, 0)}{\partial S(t, 0, 0)/\partial x}. \quad (13.21)$$

Hence we arrive at the approximation

$$\varphi(t, s) \approx \frac{\partial S}{\partial x}(t, 0, 0) \sigma_r \left(t, \frac{s - S(t, 0, 0)}{\partial S(t, 0, 0)/\partial x}, 0 \right). \quad (13.22)$$

13.1.4.2 Advanced Approximation

The approximation (13.22) is quite accurate provided that volatility is low or moderate. For an approximation with a greater range of applicability, we can consider improving (13.20) by using a higher-order expansion and by taking greater care in selecting the expansion point. Starting with the latter, we note that the conditional expectation in (13.19) is taken in the annuity measure, suggesting that the expected values of $x(t)$ and $y(t)$ under Q^A provide a good expansion point.

Proposition 13.1.4. *Let*

$$\bar{y}(t) = E^A(y(t)).$$

Then, approximately,

$$\bar{y}(t) \approx h(t)^2 \int_0^t \sigma_r^0(s)^2 h(s)^{-2} ds, \quad t \in [0, T_0], \quad (13.23)$$

where, per (13.3),

$$h(t) = \exp \left(- \int_0^t \kappa(u) du \right).$$

Proof. Recall the dynamics of $y(t)$ in the quasi-Gaussian model (Proposition 13.1.1),

$$dy(t) = \left(\sigma_r(t, x(t), y(t))^2 - 2\kappa(t)y(t) \right) dt, \quad y(0) = 0.$$

Taking expected values, we obtain

$$dE^A(y(t)) = \left(E^A(\sigma_r(t, x(t), y(t))^2) - 2\kappa(t)E^A(y(t)) \right) dt, \quad (13.24)$$

subject to $E^A(y(0)) = 0$. Approximating, for the purposes of this calculation,

$$E^A(\sigma_r(t, x(t), y(t))^2) \approx \sigma_r^0(t)^2,$$

the equation (13.24) yields

$$dE^A(y(t)) \approx (\sigma_r^0(t)^2 - 2\kappa(t)E^A(y(t))) dt, \quad E^A(y(0)) = 0.$$

Solving this ODE leads to (13.23). \square

As above, let $X(t, s)$ be the function inverse, in x , to $S(t, x, \bar{y}(t))$, i.e.

$$S(t, X(t, s), \bar{y}(t)) \equiv s, \quad (13.25)$$

and let $x_0(t)$ be given as the solution of

$$S(t, x_0(t), \bar{y}(t)) = S(0), \quad (13.26)$$

where $S(0)$ is the forward swap rate at time 0.

Remark 13.1.5. The function $S(t, x, \bar{y}(t))$ is known in closed form from (13.15) and is smooth and monotonic in x . As such, (13.26) can be solved for $x_0(t)$ in just a few iterations of the Newton algorithm. A good starting point for the search is $x = 0$.

It turns out that $x_0(t)$ in (13.26) is a good expansion point itself, and can also be used to calculate an even better one:

Lemma 13.1.6. *The function $x_0(t)$ is a first-order approximation to $E^A(x(t))$. An approximation to second order is given by*

$$\bar{x}(t) = x_0(t) + \frac{\partial^2 X}{\partial s^2}(t, S(0)) \text{Var}^A(S(t)), \quad t \in [0, T_0].$$

Proof. Expanding $X(t, s)$ around $s = S(0)$ to first order, we obtain

$$x(t) - x_0(t) = \left. \frac{\partial X}{\partial s}(t, s) \right|_{s=S(0)} (S(t) - S(0)) + O((S(t) - S(0))^2). \quad (13.27)$$

Taking expected values and using the fact that $S(t)$ is a martingale in measure Q^A , we get

$$\mathbb{E}^A(x(t)) - x_0(t) = O(\mathbb{E}^A((S(t) - S(0))^2)),$$

i.e. $x_0(t)$ is an approximation to $\mathbb{E}^A(x(t))$ to first order. The approximation $\bar{x}(t)$ is obtained by expanding (13.27) to second order,

$$\begin{aligned} x(t) - x_0(t) &= \frac{\partial X}{\partial s}(t, S(0))(S(t) - S(0)) \\ &\quad + \frac{1}{2} \frac{\partial^2 X}{\partial s^2}(t, S(0))(S(t) - S(0))^2 + O((S(t) - S(0))^3). \end{aligned}$$

Then

$$\mathbb{E}^A(x(t)) = x_0(t) + \frac{\partial^2 X}{\partial s^2}(t, S(0)) \text{Var}^A(S(t)) + O(\mathbb{E}^A((S(t) - S(0))^3)).$$

□

Remark 13.1.7. As high precision is not required when calculating the variance $\text{Var}^A(S(t))$ in Lemma 13.1.6, it can be evaluated by considering a simple Gaussian approximation to the dynamics of $S(t)$, i.e. using $x(t) = 0$, $y(t) = 0$ in (13.16),

$$dS(t) \approx \frac{\partial S}{\partial x}(t, 0, 0)\sigma_r^0(t) dW^A(t),$$

which would yield

$$\text{Var}^A(S(t)) \approx \int_0^t \left(\frac{\partial S}{\partial x}(s, 0, 0)\sigma_r^0(s) \right)^2 ds.$$

The second derivative $\partial^2 X(t, s)/\partial s^2$ can be computed by differentiating the implicit definition (13.25) twice.

Having now established an expansion point, let us proceed to determine an approximation to $X(t, s)$ with higher accuracy than the linear one in (13.21). Empirically, it can be observed that $S(t, x)$ is closely approximated as a quadratic function of x across a wide range of the argument x . This suggests a second-order expansion of (13.25) around $\bar{x}(t)$, and approximating $X(t, s)$ with $\xi = \xi(t, s)$, the solution of the following quadratic equation in ξ :

$$\begin{aligned} S(t, \bar{x}(t), \bar{y}(t)) + \frac{\partial S}{\partial x}(t, \bar{x}(t), \bar{y}(t))(\xi - \bar{x}(t)) \\ + \frac{1}{2} \frac{\partial^2 S}{\partial x^2}(t, \bar{x}(t), \bar{y}(t))(\xi - \bar{x}(t))^2 = s. \end{aligned} \quad (13.28)$$

With $S(t, \bar{x}(t), \bar{y}(t))$, $\partial S(t, \bar{x}(t), \bar{y}(t))/\partial x$ and $\partial^2 S(t, \bar{x}(t), \bar{y}(t))/\partial x^2$ pre-computed, the evaluation of $\xi(t, s)$ for any s is essentially instantaneous, and we obtain the following efficient approximation for $\varphi(t, s)$.

Proposition 13.1.8. *An approximation to $\varphi(t, s)$ in (13.19) is given by*

$$\varphi(t, s) \approx \frac{\partial S}{\partial x}(t, \xi(t, s), \bar{y}(t)) \sigma_r(t, \xi(t, s), \bar{y}(t)),$$

where $\xi(t, s)$ is the solution to the quadratic equation (13.28), with $\bar{x}(t)$ given by Lemma 13.1.6 and $\bar{y}(t)$ given by Proposition 13.1.4.

13.1.5 Linear Local Volatility

As demonstrated above, the local volatility qG model (13.9) has the flexibility to generate essentially arbitrary local volatility dynamics for swap rates. Using the results of Lemma 13.1.3 and Proposition 13.1.8, the function $\sigma_r(t, x, y)$ could therefore, in principle, be calibrated non-parametrically (see Dupire [1994] and the discussion in Section 7.1.3) to the implied volatilities of a collection of swaptions across all strikes. However, as explained in Section 7.1.3, we recommend the volatility function $\sigma_r(t, x, y)$ to be chosen from a parametric family of monotone, downward sloping functions of state variable(s). While power functions that give rise to models with CEV-type³ dynamics could be used, as explained in Remark 7.2.14 linear functions provide a less-involved alternative capable of producing essentially the same range of volatility smiles as CEV models.

With the above in mind, let us consider the following short rate local volatility function

$$\sigma_r(t, x, y) = \lambda_r(t)(\alpha_r(t) + b_r(t)x). \quad (13.29)$$

The scale function $\alpha_r(t)$ is redundant (as it can be absorbed in $\lambda_r(t)$) and may be set exogenously; Section 13.1.6 discusses a convenient choice. The functions $\lambda_r(t)$ (volatility) and $b_r(t)$ (skew) are calibrated to the market. Under (13.29), the local volatility of the swap rate S is given, approximately, by

$$\varphi(t, s) \approx \lambda_r(t) \frac{\partial S}{\partial x}(t, \xi(t, s), \bar{y}(t)) (\alpha_r(t) + b_r(t)\xi(t, s)), \quad (13.30)$$

with $\xi(t, s)$ as in Proposition 13.1.8. As the local volatility of the short rate is linear in x , it seems reasonable that the local volatility of the swap rate would be well approximated by a linear function as well. To exploit this, let φ be given as in (13.30) and notice that

$$\begin{aligned} \varphi(t, S(0)) &\approx \lambda_r(t) \frac{\partial S}{\partial x}(t, \xi(t, S(0)), \bar{y}(t)) (\alpha_r(t) + b_r(t)\xi(t, S(0))), \\ \frac{\partial \varphi}{\partial s}(t, S(0)) &\approx \lambda_r(t) \frac{\partial S}{\partial x}(t, \xi(t, S(0)), \bar{y}(t)) \frac{\partial \xi}{\partial s}(t, S(0)) \\ &\quad \times \left[\frac{\frac{\partial^2 S}{\partial x^2}(t, \xi(t, S(0)), \bar{y}(t))}{\frac{\partial S}{\partial x}(t, \xi(t, S(0)), \bar{y}(t))} (\alpha_r(t) + b_r(t)\xi(t, S(0))) + b_r(t) \right]. \end{aligned}$$

³Constant Elasticity of Variance, see Section 7.2.

Clearly

$$\xi(t, S(0)) \approx \bar{x}(t),$$

and, with $\xi(t, \cdot)$ being an approximate inverse to $S(t, \cdot, \bar{y}(t))$,

$$\frac{\partial \xi}{\partial s}(t, S(0)) \approx \frac{1}{\frac{\partial S}{\partial x}(t, \bar{x}(t), \bar{y}(t))}.$$

It follows that

$$\varphi(t, S(0)) \approx \lambda_r(t) \frac{\partial S}{\partial x}(t, \bar{x}(t), \bar{y}(t)) (\alpha_r(t) + b_r(t) \bar{x}(t)), \quad (13.31)$$

$$\frac{\partial \varphi}{\partial s}(t, S(0)) \approx \lambda_r(t) \left[\frac{\frac{\partial^2 S}{\partial x^2}(t, \bar{x}(t), \bar{y}(t))}{\frac{\partial S}{\partial x}(t, \bar{x}(t), \bar{y}(t))} (\alpha_r(t) + b_r(t) \bar{x}(t)) + b_r(t) \right]. \quad (13.32)$$

The following corollary to Proposition 13.1.8 holds.

Corollary 13.1.9. *Under the assumption of linear local short rate volatility (13.29) for the quasi-Gaussian model (13.10), the dynamics of the swap rate $S(t)$ are approximated by*

$$dS(t) \approx \lambda_S(t) (b_S(t)S(t) + (1 - b_S(t))S(0)) dW^A(t), \quad (13.33)$$

where

$$\lambda_S(t) = \lambda_r(t) \frac{1}{S(0)} \frac{\partial S}{\partial x}(t, \bar{x}(t), \bar{y}(t)) (\alpha_r(t) + b_r(t) \bar{x}(t)), \quad (13.34)$$

$$b_S(t) = \frac{S(0)}{(\alpha_r(t) + b_r(t) \bar{x}(t))} \frac{b_r(t)}{\frac{\partial S}{\partial x}(t, \bar{x}(t), \bar{y}(t))} + \frac{S(0) \frac{\partial^2 S}{\partial x^2}(t, \bar{x}(t), \bar{y}(t))}{\left(\frac{\partial S}{\partial x}(t, \bar{x}(t), \bar{y}(t)) \right)^2}. \quad (13.35)$$

Proof. Under a linear approximation to the local volatility function of the swap rate we have, with φ defined in (13.30),

$$dS(t) \approx \left(\varphi(t, S(0)) + \frac{\partial \varphi}{\partial s}(t, S(0)) (S(t) - S(0)) \right) dW^A(t)$$

which, after rearranging the terms, yields

$$dS(t) \approx \frac{\varphi(t, S(0))}{S(0)} \left(S(0) + S(0) \frac{\partial \varphi(t, S(0)) / \partial s}{\varphi(t, S(0))} (S(t) - S(0)) \right) dW^A(t).$$

Defining

$$\lambda_S(t) \triangleq \frac{\varphi(t, S(0))}{S(0)}, \quad b_S(t) \triangleq S(0) \frac{\partial \varphi(t, S(0)) / \partial s}{\varphi(t, S(0))},$$

the result follows from (13.31)–(13.32). \square

We recognize (13.33) as a displaced log-normal SDE with time-dependent volatility $\lambda_S(t)$ and skew $b_S(t)$. Using averaging techniques from Section 7.6.2, we can convert it into a displaced log-normal SDE with time-constant parameters $\bar{\lambda}_S$ and \bar{b}_S , see Proposition 7.2.12. For convenience, we list the resulting swaption pricing formula below.

Proposition 13.1.10. *Consider a payer swaption with strike (i.e., coupon) c and expiry T_0 on the swap rate $S(t)$ defined in (13.13). In the quasi-Gaussian model (13.10) with linear short rate volatility (13.29), the swaption price can be approximated by the displaced log-normal option formula*

$$V_{\text{swaption}}(0) \approx A(0) \left[(S(0) + S(0)(1 - \bar{b}_S)/\bar{b}_S) \Phi(d_+) - (c + S(0)(1 - \bar{b}_S)/\bar{b}_S) \Phi(d_-) \right],$$

$$d_{\pm} = \frac{\ln \left(\frac{S(0) + S(0)(1 - \bar{b}_S)/\bar{b}_S}{c + S(0)(1 - \bar{b}_S)/\bar{b}_S} \right) \pm \frac{1}{2} \bar{b}_S^2 \bar{\lambda}_S^2 T_0}{\bar{b}_S \bar{\lambda}_S \sqrt{T_0}},$$

where

$$\bar{\lambda}_S = \left(\frac{1}{T_0} \int_0^{T_0} \lambda_S(t)^2 dt \right)^{1/2}, \quad (13.36)$$

$$\bar{b}_S = \int_0^{T_0} b_S(t) w_S(t) dt, \quad (13.37)$$

$$w_S(t) = \frac{\lambda_S(t)^2 \int_0^t \lambda_S(s)^2 ds}{\int_0^{T_0} (\lambda_S(u)^2 \int_0^u \lambda_S(s)^2 ds) du},$$

with $\lambda_S(t)$, $b_S(t)$ given by Corollary 13.1.9.

13.1.6 Linear Local Volatility for a Swaption Strip

The quasi-Gaussian model (13.10) is typically calibrated to a swaption strip (recall the definition in Section 10.1.4) on a maturity grid $0 = T_0 < \dots < T_N$, i.e. a collection of $N - 1$ swaptions with the n -th swaption expiring on T_n , $n = 1, \dots, N - 1$. Let us suppose a swaption strip is specified, and that the n -th swaption has an underlying swap with $\mu(n)$ periods. For each n , we denote the corresponding swap rate and the annuity by

$$S_n(t) \triangleq S_{n,\mu(n)}(t), \quad A_n(t) \triangleq A_{n,\mu(n)}(t), \quad n = 1, \dots, N - 1.$$

For the model with the linear volatility specification (13.29), Proposition 13.1.10 will, after proper adjustment of maturities (see footnote 6 in Chapter 10), allow us to value all swaptions of all strikes in the strip by using the

displaced log-normal model with the effective parameters computed from the local volatility function of the model.

As mentioned earlier, in the specification (13.29) the function $\alpha_r(t)$ is only included for convenient scaling. To find a good value for it, let us consider the relationship between the local swap rate skews $b_{S_n}(t)$ and the local short rate skew $b_r(t)$ in (13.35). Ignoring small terms,

$$b_{S_n}(t) \approx \frac{S_n(0)}{\alpha_r(t)} \frac{b_r(t)}{\partial S_n(t, \bar{x}(t), \bar{y}(t)) / \partial x}. \quad (13.38)$$

It is often convenient to parametrize the model in such a way that the values of model parameters (here $b_r(t)$) are roughly of the same order of magnitude as the output parameters (here $b_{S_n}(t)$). This allows one to quickly check whether model parameters are sensible, and may well lead to better numerical properties of the calibration algorithm. Based on (13.38), we elect to set $\alpha_r(t)$ equal to $S_n(0)$ (of course we need to account for different values of n), and rescale $b_r(t)$ to incorporate the term $\partial S_n / \partial x$ (again, for different n). In summary, we specialize the definition (13.29) to be

$$\sigma_r(t, x, y) = \sum_{n=1}^{N-1} \lambda_n (S_n(0) + b_n D_n x) \mathbf{1}_{\{t \in (T_{n-1}, T_n]\}}, \quad (13.39)$$

where the *skew scalings* D_n are given by

$$D_n = \frac{\partial S_n}{\partial x}(t, 0, 0).$$

This definition recognizes the fact that the behavior of $\lambda_r(t)$ and $b_r(t)$ between the knot dates $\{T_n\}$ is of no consequence, wherefore these functions can be taken to be piecewise constant,

$$\begin{aligned} \lambda_r(t) &= \sum_{n=1}^{N-1} \lambda_n \mathbf{1}_{\{t \in (T_{n-1}, T_n]\}}, \\ \alpha_r(t) &= \sum_{n=1}^{N-1} S_n(0) \mathbf{1}_{\{t \in (T_{n-1}, T_n]\}}, \\ b_r(t) &= \sum_{n=1}^{N-1} b_n D_n \mathbf{1}_{\{t \in (T_{n-1}, T_n]\}}. \end{aligned}$$

13.1.7 Volatility Calibration

Let us assume that a swaption strip is given, and the model is parameterized with the local volatility of the form (13.39). The model parameters λ_n and b_n , $n = 1, \dots, N - 1$, need to be determined by calibrating the model to

market prices of swaptions in the swaption strip. For now, let us suppose that the mean reversion function $\varkappa(t)$ in (13.10) is specified externally — we will return to its calibration later in the chapter.

As each swap rate has an approximately displaced log-normal distribution in the model, the calibration objective could be expressed as the problem of matching *displaced log-normal parameters*, as given by the model for each swap rate, to a similar set of market-implied parameters. We already saw a similar approach in Section 9.3.4 and recall that performing the calibration in model parameter space, rather than in the space of calibration instrument *values*, avoids the expense of invoking option pricing formulas within the calibration loop.

Accordingly, assume that a collection of market parameters, i.e. displaced log-normal volatilities and skews $(\hat{\lambda}_{S_n}, \hat{b}_{S_n})$, $n = 1, \dots, N - 1$, is given. In practice, these parameters are obtained by fitting a series of constant-parameter displaced log-normal vanilla models to the observed swaption volatility smiles at all expiries T_1, \dots, T_{N-1} . A best-fit model calibration across swaption strikes is possible (at the expense of using a numerical optimizer), or we may simply set the volatility $\hat{\lambda}_{S_n}$ to match at-the-money swaption volatilities, while the skew \hat{b}_{S_n} is fit to the slope of the volatility smile at-the-money, or to the volatility at some relevant non-ATM strike.

It is clear from the swaption pricing formula in Proposition 13.1.10 that the value of a swaption with expiry T_n depends on model parameters (λ_i, b_i) for $i = 1, \dots, n$ only. Hence, the qG model can be calibrated by a bootstrap method, similarly to the pure Gaussian case from Section 10.1.4. In the bootstrap method, the equations (13.36)–(13.37) are solved sequentially for $n = 1, \dots, N - 1$, with the two equations on step n used to determine two unknown model parameters (λ_n, b_n) . For example, the following algorithm could be used.

1. Set (λ_n, b_n) , $n = 1, \dots, N - 1$, to some reasonable starting values, e.g. set λ_n 's to (properly scaled) volatilities obtained by calibrating a pure Gaussian model as in Section 10.1.4, and $b_n = \hat{b}_{S_n}$.
2. Set $n = 1$.
3. For given n , (λ_i^*, b_i^*) are known for $i = 1, \dots, n - 1$. Note we use a star to denote *calibrated* values of the model parameters.
4. Calculate $\bar{x}(t)$ (Lemma 13.1.6) and $\bar{y}(t)$ (Proposition 13.1.4) for $t \in [0, T_n]$ using (λ_i^*, b_i^*) , $i = 1, \dots, n - 1$, and the initial guess for (λ_n, b_n) from Step 1. Note that $\bar{x}(t), \bar{y}(t)$ implicitly depend on n as their definition depends on the swap rate/annuity measure used.
5. Calculate $\lambda_{S_n}(t), b_{S_n}(t)$ for $t \in [0, T_{n-1}]$ from (λ_i^*, b_i^*) , $i = 1, \dots, n - 1$, using (13.34)–(13.35).
6. Make another guess for (λ_n, b_n) .
7. Update $\lambda_{S_n}(t), b_{S_n}(t)$ for $t \in (T_{n-1}, T_n]$ from (λ_i^*, b_i^*) , $i = 1, \dots, n - 1$, using (13.34)–(13.35).
8. Calculate $\bar{\lambda}_{S_n}, \bar{b}_{S_n}$ using Proposition 13.1.10.

9. Compare $(\bar{\lambda}_{S_n}, \bar{b}_{S_n})$ to $(\hat{\lambda}_{S_n}, \hat{b}_{S_n})$. If not equal within given tolerance, go to Step 6. Otherwise, proceed to Step 10.
10. As we have reached acceptable convergence between $(\bar{\lambda}_{S_n}, \bar{b}_{S_n})$ and $(\hat{\lambda}_{S_n}, \hat{b}_{S_n})$, set the calibrated model parameter values to the latest trial values, $(\lambda_n^*, b_n^*) = (\lambda_n, b_n)$.
11. Update $n \rightarrow n + 1$. If $n \leq N - 1$ go to Step 3. Otherwise, conclude.

It may appear more accurate to make Step 4 a part of the calibration loop (for each n), with $\bar{x}(t)$, $\bar{y}(t)$, $t \in (T_{n-1}, T_n]$, and dependent quantities such as $\partial S_n(t, \bar{x}(t), \bar{y}(t))/\partial x$, etc. computed using the current guess for (λ_n, b_n) in each loop iteration (together with already-calibrated values (λ_i^*, b_i^*) , $i = 1, \dots, n - 1$). While such extensions are relatively straightforward, our experience shows that the quality of the calibration is rarely improved enough to justify the additional complexity.

It is possible to add regularity terms to Step 9 of the algorithm above, to ensure that the resulting model parameters do not behave irregularly as functions of t . See Section 13.1.8.2 below for an example of this.

13.1.8 Mean Reversion Calibration

With the volatility calibration out of the way, let us discuss what to do with the remaining model parameter, the mean reversion function $\varkappa(t)$ in (13.10). We start with a short review of the effects of mean reversion.

13.1.8.1 Effects of Mean Reversion

Let us first consider a few simple examples as a way of building intuition about the effects of mean reversion on market values of various securities. For simplicity, we use continuously compounded rates as convenient proxies for Libor and swap rates, and consider a pure Gaussian model with constant volatility and mean reversion,

$$\sigma_r(t) \equiv \sigma_r, \quad \varkappa(t) \equiv \varkappa.$$

A continuously compounded forward yield over a period $[T, M]$, observed at time t , is given by⁴

$$F(t, T, M) = -\frac{1}{M - T} \ln \frac{P(t, M)}{P(t, T)}.$$

⁴We normally use $y(t, T, M)$ for continuously compounded forward yield, see Section 4.1.1, and $F(t, T, M)$ for a futures rate, see Section 4.1.2, but for the remainder of this chapter only we allow ourselves a slight notational inconsistency to avoid the possibility of confusion between y , the state variable, and y , the forward yield.

According to Proposition 10.1.7, the forward yield can be expressed as a function of the state variables and parameters of the model,

$$F(t, T, M) = \frac{G(t, M) - G(t, T)}{M - T} x(t) + \frac{1}{2} \frac{G(t, M)^2 - G(t, T)^2}{M - T} y(t),$$

$$G(t, u) = \frac{1 - e^{-\kappa(u-t)}}{\kappa}.$$

Recall that $y(t)$ is deterministic in the Gaussian case, so the standard deviation of $F(T, T, M)$ is equal to

$$\text{Stdev}(F(T, T, M)) = \frac{G(T, M) - G(T, T)}{M - T} (\text{Var}(x(T)))^{1/2}$$

$$= \frac{1 - e^{-\kappa(M-T)}}{\kappa(M - T)} (\text{Var}(x(T)))^{1/2}.$$

For two maturities M_1, M_2 , $T \leq M_1 \leq M_2$, we therefore observe that

$$\frac{\text{Stdev}(F(T, T, M_2))}{\text{Stdev}(F(T, T, M_1))} = \left(\frac{1 - e^{-\kappa(M_2-T)}}{\kappa(M_2 - T)} \right) / \left(\frac{1 - e^{-\kappa(M_1-T)}}{\kappa(M_1 - T)} \right),$$

i.e. the ratio of standard deviations of two forward yields with the same expiry T but different tenors is independent of the volatility parameter σ_r , and solely determined by the mean reversion parameter κ . Since standard deviations of forward yields can loosely be thought of as proxies for implied swaption volatilities, we observe that the mean reversion parameter changes the relative levels of implied volatilities of swaptions with the same expiry but different underlying swap tenors. Specifically, for a fixed level of volatility σ_r , an increase in mean reversion makes the volatility of a longer-tenor swaption decrease relative to the volatility of a shorter-tenor swaption, assuming both have the same expiry date.

With the mean reversion effect above in mind, consider (say) a caplet and a swaption with the same expiry, and imagine an experiment in which the mean reversion is changed, but the volatility σ_r is adjusted to keep the implied volatility of the caplet unchanged. It should be clear from the discussion above that as mean reversion increases, the swaption volatility will *decrease*. If, instead, the market volatility of the swaption is kept constant by adjusting σ_r for each level of mean reversion, then the caplet volatility will *increase* when the mean reversion increases. Such complementarity of effects of the mean reversion κ and volatility σ_r allows us, in principle, to set both in such a way that we match the market-implied volatilities of both the caplet and the swaption.

For an alternative look at the effect of mean reversion, let us consider two forward rates with different fixing dates, $F(T_1, T_1, M_1)$ and $F(T_2, T_2, M_2)$ with $T_1 \leq T_2$. Observing that

$$x(t) = \sigma_r \int_0^t e^{-\kappa(t-u)} dW(u),$$

it is easy to establish that

$$\begin{aligned} \text{Corr}(F(T_1, T_1, M_1), F(T_2, T_2, M_2)) \\ = \text{Corr}(x(T_1), x(T_2)) = e^{-\kappa(T_2-T_1)} \left(\frac{1 - e^{-2\kappa T_1}}{1 - e^{-2\kappa T_2}} \right)^{1/2}. \end{aligned}$$

This correlation, which we can call *inter-temporal correlation* as the forward yields are observed at different times, depends only on κ , i.e. on the auto-correlation properties of the process $x(t)$. At a level of $\kappa = 0$, the inter-temporal correlation is $(T_1/T_2)^{1/2}$ and decreases to 0 as $\kappa \rightarrow \infty$.

The dependence of inter-temporal correlation on mean reversion is potentially useful for calibration purposes. To see this, consider a hypothetical security, a basket option on a set of rates with different expiry dates, with a payoff

$$\max \{F(T_n, T_n, M_n), n = 1, \dots, N-1\}.$$

It is well-known (and intuitively obvious) that prices of basket options are decreasing functions of correlation, hence we expect the price of the contract above to be increasing in the mean reversion κ , *ceteris paribus*. If such a basket option were traded in the market, the mean reversion could in principle be implied from its price.

While basket options on forward yields are, of course, not traded outright, the example above is not as far-fetched as it might seem since a Bermudan swaption (see Section 5.12) gives the holder a right to exercise into one of several different swaps that, critically, are observed on different exercise dates. The Bermudan swaption is therefore conceptually similar to a basket option on a strip of swap rates, with each rate fixing on its own fixing date. The implications of this analogy, and the effect of mean reversion on inter-temporal correlations, will be exploited later in developing the local projection method for Bermudan swaptions (see Chapter 19).

13.1.8.2 Calibrating Mean Reversion to Volatility Ratios

We consider the qG model with the volatility function (13.39). In Section 13.1.7 we developed the volatility calibration algorithm for this model under the assumption that the mean reversion function $\kappa(t)$ was already available. Since it is often beneficial to calibrate different model parameters separately, the calibration of mean reversion should then ideally not require the knowledge of the volatility $\sigma_r(t, x, y)$ of the model.

As presented earlier, the qG model volatility calibration involves a strip of swap rates $\{S_n(\cdot)\}_{n=1}^{N-1}$, with the rate $S_n(\cdot)$ fixing on T_n and having $\mu(n)$ periods; these swap rates define the volatility function $\sigma_r(t, x, y)$ in (13.39)

and serve as volatility calibration targets. As indicated in Section 13.1.8.1, the ratio of volatilities of two swaptions with the same expiry date is more-or-less independent of volatility, suggesting the usage of a *second strip*⁵ of swap rates to define targets for mean reversion calibration as the ratios between these rates and the original ones.

We assume that a second strip of swap rates is given, with the n -th rate fixing on T_n and having $\nu(n)$ periods, where $\nu(n) \neq \mu(n)$. We use extended notations to distinguish the two strips, with

$$\{S_{n,\mu(n)}(\cdot)\}_{n=1}^{N-1}$$

used for the original strip and

$$\{S_{n,\nu(n)}(\cdot)\}_{n=1}^{N-1}$$

used for the additional one. The mean reversion is then calibrated to the pairwise ratios of implied volatilities of $S_{n,\nu(n)}(T_n)$ and $S_{n,\mu(n)}(T_n)$, $n = 1, \dots, N - 1$. The following result forms the basis of mean reversion calibration.

Proposition 13.1.11. *In the quasi-Gaussian model with local short rate volatility function (13.9), the ratio of variances of two swap rates fixing on T_n with m_1 and m_2 periods, respectively, is approximately given by either*

$$\frac{\text{Var}(S_{n,m_1}(T_n))}{\text{Var}(S_{n,m_2}(T_n))} \approx \frac{\int_0^{T_n} \left(\sigma_r^0(t) \frac{\partial S_{n,m_1}}{\partial x}(t, 0, 0) \right)^2 dt}{\int_0^{T_n} \left(\sigma_r^0(t) \frac{\partial S_{n,m_2}}{\partial x}(t, 0, 0) \right)^2 dt}, \quad (13.40)$$

where $\sigma_r^0(t)$ is defined in (13.12), or

$$\frac{\text{Var}(S_{n,m_1}(T_n))}{\text{Var}(S_{n,m_2}(T_n))} \approx \frac{\int_0^{T_n} \left(\frac{\partial S_{n,m_1}}{\partial x}(t, 0, 0) \right)^2 dt}{\int_0^{T_n} \left(\frac{\partial S_{n,m_2}}{\partial x}(t, 0, 0) \right)^2 dt}. \quad (13.41)$$

Proof. The first formula is obtained by using $x(t) = 0$, $y(t) = 0$ in (13.16); see also Remark 13.1.7. The second one follows from the first under the approximation of the time-dependent volatility $\sigma_r^0(t)$ by a constant,

$$\sigma_r^0(t) \approx \sigma_r^0(0). \quad (13.42)$$

□

⁵Note that it is also possible to calibrate mean reversion to a whole *collection* of European swaptions (i.e. more than two) sharing the same expiry; the motivation for such calibration and an outline of the algorithm are presented in Section 19.4.4.

Remark 13.1.12. In practice, using the simpler approximation (13.41) instead of (13.40) does not reduce the accuracy of mean reversion calibration much. However, if the more accurate formula (13.40) is preferred, one can use an estimate of $\sigma_r^0(t)$ obtained from, for example, a pre-calibration of a pure Gaussian model.

Proposition 13.1.11 suggests an algorithm for mean reversion calibration. Recalling formulas (13.17), (13.3) we notice that the ratios

$$\left\{ \frac{\text{Var}(S_{n,\nu(n)}(T_n))}{\text{Var}(S_{n,\mu(n)}(T_n))} \right\}_{n=1}^{N-1}, \quad (13.43)$$

as computed in (13.41) depend on the mean reversion parameter $\kappa(t)$ only, allowing us to set up an optimization problem in which $\kappa(t)$ is chosen to match the ratios (13.43) to their market-implied values.

For a possible calibration algorithm, suppose that we discretize $\kappa(t)$ on the time grid,

$$\kappa(t) = \sum_{n=1}^{N-1} \kappa_n \times 1_{\{t \in (T_{n-1}, T_n]\}} + \kappa_N \times 1_{\{t \in (T_{N-1}, \infty)\}}.$$

It is generally not advisable to find these mean reversions only by best-fitting to variance ratios, as the function $\kappa(t)$ would likely end up being quite irregular. To prevent this, we suggest the inclusion of a regularization term penalizing non-stationary behavior. While there are multiple ways of doing this, one simple approach would be to set optimal mean reversion levels $\{\kappa_n^*\}_{n=1}^N$ as the solution to the following optimization problem,

$$\{\kappa_n^*\} = \underset{\{\kappa_n\}}{\text{argmin}} \left\{ \sum_{n=1}^{N-1} \left(\frac{\text{Var}(S_{n,\nu(n)}(T_n))}{\text{Var}(S_{n,\mu(n)}(T_n))} (\{\kappa_n\}) - \widehat{\text{Var}}(S_{n,\nu(n)}(T_n)) \right)^2 + w \sum_{n=1}^{N-1} (\kappa_{n+1} - \kappa_n)^2 \right\}, \quad (13.44)$$

where

$$\frac{\text{Var}(S_{n,\nu(n)}(T_n))}{\text{Var}(S_{n,\mu(n)}(T_n))} (\{\kappa_n\}) = \frac{\int_0^{T_n} \left(\frac{\partial S_{n,\nu(n)}}{\partial x}(t, 0, 0) \right)^2 dt}{\int_0^{T_n} \left(\frac{\partial S_{n,\mu(n)}}{\partial x}(t, 0, 0) \right)^2 dt}.$$

Here $w > 0$ is a user-specified regularization weight, and the $\widehat{\text{Var}}(S_{n,m})$'s are market-implied variances of swap rates. This minimization problem is easily solved by standard non-linear optimization methods, to be discussed in some detail later (in particular in Section 14.5).

At this point the reader might, of course, wonder whether target variances of swap rates, as required for the mean reversion calibration, could indeed be observed in the market. The answer is yes: in general, a market-implied variance of a swap rate can be calculated directly from values of options on the swap rate (i.e. swaptions) across a collection of strikes, as discussed in detail in Chapter 16. In the most commonly used linear local short rate volatility (13.29) case, we may simplify matters further. Indeed, in this model the market volatility parameters of swap rates are the displaced log-normal volatilities $\hat{\lambda}_{S_{n,m}}$ (see (13.33) and Section 13.1.7), in which case the market-implied variance of a swap rate can be approximated by

$$\widehat{\text{Var}}(S_{n,m}(T_n)) \approx \left(S_{n,m}(0)\hat{\lambda}_{S_{n,m}}\right)^2 T_n. \quad (13.45)$$

Hence, in this case the mean reversion calibration targets could be specified directly in terms of market-implied displaced log-normal volatilities. In any case, when paired up with the volatility calibration procedure from Section 13.1.7, numerical solution of the optimization problem (13.44) ultimately allows us to calibrate the model to market-implied values of two separate swaption strips.

13.1.8.3 Calibrating Mean Reversion to Inter-Temporal Correlations

In the previous section, we took advantage of the observation from Section 13.1.8.1 that a ratio of standard deviations of two rates with the same expiry in a one-factor constant-volatility Gaussian model is independent of volatility. In this section, we develop a different calibration method for mean reversion, based on our earlier observation that inter-temporal correlations between forward yields are also nearly independent of volatility.

We focus on correlations (for the rest of this section, we omit the “inter-temporal” qualifier for brevity) of the original strip of swap rates $\{S_n(\cdot)\}_{n=1}^{N-1}$. While a different swap rate strip could be used, the practical importance of such a generalization is limited.

Proposition 13.1.13. *Let $\sigma_r^0(t)$ be given by (13.12). The correlation between two rates $S_{n_1}(T_{n_1})$ and $S_{n_2}(T_{n_2})$, $n_1 \leq n_2$, in the quasi-Gaussian model with the general local volatility function (13.9) can be approximated by either*

$$\begin{aligned}
& \text{Corr}(S_{n_1}(T_{n_1}), S_{n_2}(T_{n_2})) \\
& \approx \int_0^{T_{n_1}} \left(\frac{\partial S_{n_1}}{\partial x}(t, 0, 0) \right) \left(\frac{\partial S_{n_2}}{\partial x}(t, 0, 0) \right) \sigma_r^0(t)^2 dt \\
& \quad \times \left(\int_0^{T_{n_1}} \left(\frac{\partial S_{n_1}}{\partial x}(t, 0, 0) \right)^2 \sigma_r^0(t)^2 dt \right)^{-1/2} \\
& \quad \times \left(\int_0^{T_{n_2}} \left(\frac{\partial S_{n_2}}{\partial x}(t, 0, 0) \right)^2 \sigma_r^0(t)^2 dt \right)^{-1/2}, \tag{13.46}
\end{aligned}$$

or

$$\begin{aligned}
& \text{Corr}(S_{n_1}(T_{n_1}), S_{n_2}(T_{n_2})) \approx \int_0^{T_{n_1}} \left(\frac{\partial S_{n_1}}{\partial x}(t, 0, 0) \right) \left(\frac{\partial S_{n_2}}{\partial x}(t, 0, 0) \right) dt \\
& \times \left(\int_0^{T_{n_1}} \left(\frac{\partial S_{n_1}}{\partial x}(t, 0, 0) \right)^2 dt \right)^{-1/2} \left(\int_0^{T_{n_2}} \left(\frac{\partial S_{n_2}}{\partial x}(t, 0, 0) \right)^2 dt \right)^{-1/2}. \tag{13.47}
\end{aligned}$$

Proof. By Proposition 13.1.2,

$$dS_{n_i}(t) = \frac{\partial S_{n_i}}{\partial x}(t, x(t), y(t)) \sigma_r(t, x(t), y(t)) dW^{A_{n_i}}(t)$$

in $Q^{A_{n_i}}$, $i = 1, 2$. In the risk-neutral measure, the SDE for $S_{n_i}(t)$ will have a stochastic drift; however, for the purposes of calculating the correlations in question, we ignore drift contributions. Thus, in the risk-neutral measure Q we have approximately that

$$\begin{aligned}
& E(S_{n_1}(T_{n_1}) S_{n_2}(T_{n_2})) - S_{n_1}(0) S_{n_2}(0) \\
& = E(S_{n_1}(T_{n_1}) S_{n_2}(T_{n_1})) - S_{n_1}(0) S_{n_2}(0) \\
& = E \int_0^{T_{n_1}} \frac{\partial S_{n_1}}{\partial x}(t, x(t), y(t)) \frac{\partial S_{n_2}}{\partial x}(t, x(t), y(t)) \sigma_r(t, x(t), y(t))^2 dt.
\end{aligned}$$

Using (13.12) and approximating all functions of state variables with their values at $x = y = 0$, we obtain

$$\begin{aligned}
& E(S_{n_1}(T_{n_1}) S_{n_2}(T_{n_2})) - S_{n_1}(0) S_{n_2}(0) \\
& \approx \int_0^{T_{n_1}} \left(\frac{\partial S_{n_1}}{\partial x}(t, 0, 0) \right) \left(\frac{\partial S_{n_2}}{\partial x}(t, 0, 0) \right) \sigma_r^0(t)^2 dt.
\end{aligned}$$

Similar formulas hold for $E(S_{n_i}(T_{n_i})^2)$, $i = 1, 2$, and (13.46) follows. The result (13.47) follows from (13.46) by approximating the deterministic volatility function $\sigma_r^0(t)$ with a constant,

$$\sigma_r^0(t) \approx \sigma_r^0(0).$$

□

If a “market-implied” correlation matrix

$$\hat{\chi}_{n_1, n_2} = \text{Corr}(S_{n_1}(T_{n_1}), S_{n_2}(T_{n_2})), \quad 1 \leq n_1 \leq n_2 < N - 1,$$

is somehow known — extracted, for example, from Bermudan swaption prices that, per discussion in Section 13.1.8.1, depend strongly on such correlations — then Proposition 13.1.13 can be used to calibrate the mean reversion function $\varkappa(t)$ to this matrix. The formula (13.47) is independent of the volatility term $\sigma_r(t, x, y)$, and the mean reversion calibration can precede volatility calibration. Alternatively, the more accurate formula (13.46) could be used with the volatility obtained from a pre-calibration of a pure Gaussian model.

13.1.8.4 Final Comments on Mean Reversion Calibration

The reader has undoubtedly noticed that the formulas for mean reversion calibration were derived using rather crude approximations. Refinements are certainly possible, but achieving a high level of accuracy for mean reversion calibration was never our objective. Indeed, while it is possible to execute a global calibration to vanilla option in which the mean reversion and the volatility function are calibrated together using numerical valuation methods (such as the PDE method, see later in the chapter), we believe that in setting mean reversion, market information should be a rough guide rather than a “hard” calibration target. For this reason, we recommend using mean reversion calibration to match ratios of volatilities or inter-temporal correlations in an approximate sense only; subsequently, we can apply much more precise volatility calibration to recover our main targets, namely the implied volatilities of the primary swaption strip.

Calibrating a model with relatively limited set of parameters to market inevitably leads to time-dependent parameters and, subsequently, questions about the stationarity of the resulting volatility structure. In most applications, a completely time-stationary model would yield a clearly unacceptable fit to market data, and a certain degree of non-stationarity is unavoidable. However, as far as mean reversion calibration is concerned, we often advocate using a *constant* mean reversion function $\varkappa(t) \equiv \varkappa$ in the calibration routines developed previously. For example, we can use a constant mean reversion \varkappa to roughly match the volatility ratios of caplets and swaptions, and then calibrate a time-dependent volatility function to match the volatilities of swaptions in the swaption strip exactly. To encourage even more time stationarity in the model, it is also possible to set (through an optimizer) the mean reversion \varkappa in such a way that the volatility parameters of the calibrated model (i.e. the λ_n 's) are as close to constant as possible.

Let us note that there is one instance of the application of the model where we recommend time-dependent mean reversion, namely when the

model is used as part of the *local projection method*. The method is developed in detail later in the book (see e.g. Sections 18.4, 19.2, 20.1.3, 20.2.1), but it can loosely be described as using a “small” model, such as the qG model, as a local or instrument-specific proxy for a “big” model, such as a Libor market model (see Chapter 14). In this case the dynamics of the volatility structure are defined by the “big” and, hopefully, realistic model, and the local model is effectively just a mechanism to reduce numerical complexity of valuation.

Finally, we note that it is, of course, also possible to use the qG model without explicitly calibrating the mean reversion parameter. For example, it could be exposed as an “exotic risk” parameter and set exogenously by a trader to reflect his estimation of the market prices of Bermudan swaptions or other exotic securities. Such practice is, we believe, quite common.

13.1.9 Numerical Methods

13.1.9.1 Direct Integration

We start the discussion of numerical methods for pricing derivatives in the quasi-Gaussian model by deriving an approximation to the density of the state variables $x(T)$ and $y(T)$, $T > 0$. Our approximation is constructed to be suitable for small T ; we find that it has good accuracy for T around 1–2 years or less, depending on the level of volatility. While usable for valuing European-style derivatives by direct integration — a method generally preferable to PDE or Monte Carlo methods when available — the real utility of having a probability density comes in improving the accuracy of the PDE method, as described in Sections 2.8.2 and 2.8.3.

Consider a contract with a payoff $V(x(T), y(T))$ at time T . Its value at time 0 is given by

$$V_0 = P(0, T) \mathbb{E}^T (V(x(T), y(T))),$$

where \mathbb{E}^T (as always) denotes expectation in the T -forward measure Q^T . Accordingly, we seek the density of $x(T), y(T)$ in the T -forward measure. As we focus on short times to maturity, it suffices to replace $y(T)$ with a deterministic approximation $\bar{y}(T)$ from Proposition 13.1.4. The local volatility (13.9) and mean reversion $\kappa(t)$ are, as a rule, piecewise constant in time; thus, in the small- T regime, they can be assumed to be independent of t . It is also safe to ignore the dependence of the volatility function on y . With this in mind, we define

$$\begin{aligned} \kappa &= \kappa(0), \\ v(x) &= \sigma_r(0, x, 0) / \sigma_r^0(0). \end{aligned} \tag{13.48}$$

The following result holds.

Theorem 13.1.14. Let us define $\pi(x)$ by the ODE

$$v(x)(\pi'(x))^2 + 2\kappa G_2(T)\pi(x)(x/\pi(x))' - 1 = 0, \quad x \in \mathbb{R}, \quad \pi(0) = 0,$$

where $G_2(T)$ is defined by (see (13.3) for the definition of $h(t)$)

$$G_2(T) = \int_0^T h(s)^2 ds, \quad (13.49)$$

and the prime denotes differentiation with respect to x . Set

$$\varpi(x) = x/\pi(x),$$

and let $\Psi(T, x)$ be the CDF of $x(T)$ in the T -forward measure. Then

$$\begin{aligned} \Psi(T, x) \approx & \Phi \left(\frac{x}{\sigma_r^0(0)\varpi(x)\sqrt{G_2(T)}} \right) \\ & + \sigma_r^0(0)\sqrt{G_2(T)}\varpi'(x)\phi \left(-\frac{x}{\sigma_r^0(0)\varpi(x)\sqrt{G_2(T)}} \right), \end{aligned}$$

where $\Phi(z)$, $\phi(z)$ are the standard Gaussian CDF and PDF, respectively.

Proof. While lengthy and somewhat technical, the proof is instructive as it shows a general approach to deriving short-time densities for local volatility term structure models. Full details are shown in Appendix 13.A of this chapter. \square

Remark 13.1.15. If $\sigma_r(0, x, 0) = \text{const}$ (pure Gaussian case), then $v(x) = 1$, and $\pi(x) = x$ is a solution to the ODE. Therefore, $\varpi(x) = 1$, and we recover the Gaussian CDF of $x(T)$ as expected. The function $\varpi(x) = x/\pi(x)$ measures the deviation of the model from the Gaussian case; $\sigma_r^0(0)\varpi(x)$ can be thought of as an “effective term volatility” at x .

Equipped with Theorem 13.1.14, we can recover an approximation for the density

$$\psi(T, x) = \frac{\partial \Psi}{\partial x}(T, x),$$

which allows us to value (short-dated) derivative contracts by numerical integration. Specifically, the value of a contract with a payoff $V(x(T), y(T))$ at time T in the quasi-Gaussian local volatility model is approximately equal to

$$V_0 \approx P(0, T) \int_{-\infty}^{\infty} V(x, \bar{y}(T)) \psi(T, x) dx,$$

where $\bar{y}(T)$ is given in Proposition 13.1.4. Again, the primary use of the results in Theorem 13.1.14 is to improve on the finite difference method, using the results of Sections 2.8.2 and 2.8.3.

13.1.9.2 Finite Difference Methods

We consider the model in its general local volatility form (13.10),

$$\begin{aligned} dx(t) &= (y(t) - \kappa(t)x(t)) dt + \sigma_r(t, x(t), y(t)) dW(t), \\ dy(t) &= \left(\sigma_r(t, x(t), y(t))^2 - 2\kappa(t)y(t) \right) dt. \end{aligned}$$

Using methods from Chapter 2, a PDE for the value of a security as a function of x, y can be easily derived from these SDEs. We find it more convenient, however, to transform the variables first to replace $y(t)$ with a locally deterministic process $u(t)$ that is drift-free on average.

Recall the definition (13.12), and define the deterministic function $\bar{y}(t)$ by (as in Proposition 13.1.4)

$$d\bar{y}(t) = (\sigma_r^0(t)^2 - 2\kappa(t)\bar{y}(t)) dt, \quad \bar{y}(0) = 0,$$

or

$$\bar{y}(t) = h(t)^2 \int_0^t \sigma_r^0(s)^2 h(s)^{-2} ds.$$

We define a new, normalized auxiliary variable $u(t)$ by

$$u(t) = y(t) - \bar{y}(t), \quad (13.50)$$

so that the state process $(x(t), u(t))$ satisfies

$$\begin{aligned} dx(t) &= (u(t) + \bar{y}(t) - \kappa(t)x(t)) dt + \sigma_r(t, x(t), u(t) + \bar{y}(t)) dW(t), \\ (13.51) \end{aligned}$$

$$du(t) = \left(\left(\sigma_r(t, x(t), u(t) + \bar{y}(t))^2 - \sigma_r^0(t)^2 \right) - 2\kappa(t)u(t) \right) dt, \quad (13.52)$$

subject to $x(0) = u(0) = 0$. Values of zero-coupon discount bonds can easily be re-expressed in terms of the new state variables,

$$P(t, T, x, u) = \frac{P(0, T)}{P(0, t)} \exp \left(-G(t, T)x - \frac{1}{2}G(t, T)^2 u - \frac{1}{2}G(t, T)^2 \bar{y}(t) \right).$$

The new parameterization of the qG model reduces nicely to the (Gaussian) case of deterministic volatility, in the following sense. If $\sigma_r(t, x, y)$ is independent of x and y , then

$$\sigma_r(t, x, y) \equiv \sigma_r^0(t), \quad (13.53)$$

and the SDE for the state variable $u(t)$ becomes

$$du(t) = -2\kappa(t)u(t) dt, \quad u(0) = 0,$$

with the unique solution

$$u(t) \equiv 0. \quad (13.54)$$

Thus, in the pure Gaussian case, the system of SDEs (13.51) reduces to a single SDE, in line with the way a one-factor Gaussian model was developed in Section 10.1.2.2.

Aesthetic reasons aside, the change of variable from $y(t)$ to $u(t)$ improves the numerical properties of a discretization of the PDE. Specifically, the variable y (or u) does not have a diffusion term so, at least in the y direction, the PDE is convection-dominated, in the sense described in Section 2.6. Removing most of the drift outside of the time-stepping scheme alleviates some of the numerical issues associated with such PDEs⁶.

The PDE associated with the dynamics (13.51) is derived in the standard way. Let $V(t, x, u)$ be the value, at time t , of a derivative with a payoff $V(x(T), u(T))$ at time T , given that $x(t) = x$, $u(t) = u$,

$$V(t, x, u) = \mathbb{E} \left(e^{-\int_t^T r(s) ds} V(x(T), u(T)) \middle| x(t) = x, u(t) = u \right).$$

Then the function $V(t, x, u)$ satisfies the PDE

$$\frac{\partial V}{\partial t}(t, x, u) + (\mathcal{L}V)(t, x, u) = (f(0, t) + x)V(t, x, u), \quad 0 \leq t < T, \quad (13.55)$$

$$V(T, x, u) = V(x, u),$$

where

$$\mathcal{L} = \mathcal{L}_x + \mathcal{L}_u,$$

and

$$\begin{aligned} \mathcal{L}_x &= (u + \bar{y}(t) - \varkappa(t)x) \frac{\partial}{\partial x} + \frac{1}{2} \sigma_r(t, x, u + \bar{y}(t))^2 \frac{\partial^2}{\partial x^2}, \\ \mathcal{L}_u &= \left((\sigma_r(t, x, u + \bar{y}(t))^2 - \sigma_r^0(t)^2) - 2\varkappa(t)u \right) \frac{\partial}{\partial u}. \end{aligned}$$

The PDE (13.55) has two space dimensions and no mixed derivatives, and can be solved numerically by the Douglas-Rachford ADI method outlined in Section 2.10. The fact that $u(t)$ has no diffusion term does not complicate the situation much; even without upwinding (which we nevertheless recommend), we have observed no noticeable deterioration in the stability or accuracy of the scheme. There is some evidence that a 5-point discretization in the u direction may improve precision slightly, as may semi-Lagrangian schemes (see e.g. Chen and Forsyth [2007]); as standard methods produce adequate results, we consider such improvements optional and leave them to the reader to explore.

⁶Removing deterministic drift components from variables before discretizing a PDE is a useful trick for any model.

The simple nature of the process for $u(t)$ typically allows one to use a rather coarse discretization in the u direction. For instance, a typical setting might involve $n_t = 100$ time steps, $n_x = 150$ steps in x direction, and $n_u = 10$ steps in the u direction. With these settings, most instruments are priced to basis point precision. It should be noted that n_u can be chosen to reflect the degree to which the model deviates from the pure Gaussian case. To elaborate, consider a volatility term of the form (13.29). If $b = 0$ the model is Gaussian, in which case only *one* discretization point ($u_0 = 0$) is required, as is clear from (13.54). As b is increased, the model becomes increasingly non-Gaussian, and an ever larger number of points in u direction are required to maintain adequate precision. In other words, a practical scheme would set u as an increasing function of the skew parameter b .

The choice of the domain for $(x(t), u(t))$ follows standard prescriptions from Chapter 2. Boundaries in the x dimension are most easily obtained under the Gaussian approximation to the short rate state dynamics

$$dx(t) \approx (\bar{y}(t) - \varkappa(t)x(t)) dt + \sigma_r^0(t) dW(t), \quad x(0) = 0, \quad (13.56)$$

using the formula (10.31) for the size of the grid while calculating $E(x(T))$, $\text{Var}(x(T))$ from (13.56) (see also footnote 7 in Chapter 10). For a slightly more refined approach, we note that with the linear local volatility (13.29) (or under a linear approximation in x to the general volatility function $\sigma_r(t, x, y)$) the distribution of $x(T)$ is closer to a displaced log-normal than to a Gaussian; we can then set the boundaries by moment-matching a displaced log-normal variable to the distribution of $x(T)$ and using appropriate quantiles. We leave it to the reader to fill in missing details.

A slightly more interesting question is how to dimension the grid in the u direction. As $u(t)$ does not have its own diffusion term, the randomness in $u(T)$ comes from the stochasticity of the drift in (13.52) which, to first order, is driven by $x(t)$. To extract this dependence, we apply a linear approximation to the volatility function $\sigma_r(t, x, y)$,

$$\sigma_r(t, x, u + \bar{y}(t))^2 - \sigma_r^0(t)^2 \approx 2\sigma_r^0(t) (\partial\sigma_r(t, 0, 0) / \partial x) x,$$

in the SDE (13.52), which allows us to integrate (13.52) in an approximate sense,

$$u(T) \approx 2h(T)^2 \int_0^T x(t) \sigma_r^0(t) (\partial\sigma_r(t, 0, 0) / \partial x) h(t)^{-2} dt.$$

Then, under the Gaussian approximation (13.56) to the dynamics of $x(t)$, we see that $u(T)$ is also Gaussian with an easily calculated variance, which allows us to set the boundaries in the u direction using an analog to the formula (10.31).

Once the grid boundaries are determined, appropriate spatial boundary conditions need to be specified. This is straightforward, with the ideas from

Section 10.1.5 easily transferable to the quasi-Gaussian model. For instance, a reasonable boundary specification is to assume linearity in V as a function of u at the u -boundaries (as in Section 2.2.2), while using the PDE itself (as in Section 10.1.5.2) to establish the boundary conditions at the x -boundaries.

13.1.9.3 Monte Carlo Simulation

Application of the Monte Carlo method to the quasi-Gaussian model is straightforward and can follow general guidelines from Chapter 3. As with the PDE method, there may be some accuracy gains associated with using the state variable $u(t)$ in (13.50) instead of $y(t)$.

Consider the problem of computing the value of a derivative with the payoff $V(x(T), u(T))$ at time T . Let

$$0 = t_0 < t_1 < \dots < t_N = T, \quad \Delta_n = t_n - t_{n-1},$$

be the discretization of the time domain. By applying a standard Euler discretization (see Section 3.2.3) to (13.51), the following stepping scheme is obtained,

$$\begin{aligned} \widehat{x}_n &= \widehat{x}_{n-1} + (\widehat{u}_{n-1} + \bar{y}(t_{n-1}) - \varkappa(t_{n-1}) \widehat{x}_{n-1}) \Delta_n \\ &\quad + \sigma_r(t_{n-1}, \widehat{x}_{n-1}, \widehat{u}_{n-1} + \bar{y}(t_{n-1})) Z_n \sqrt{\Delta_n}, \\ \widehat{u}_n &= \widehat{u}_{n-1} + \Delta_n \\ &\quad \times \left(\sigma_r(t_{n-1}, \widehat{x}_{n-1}, \widehat{u}_{n-1} + \bar{y}(t_{n-1}))^2 - \sigma_r^0(t_{n-1})^2 - 2\varkappa(t_{n-1}) \widehat{u}_{n-1} \right), \end{aligned}$$

where $\widehat{x}_0 = \widehat{u}_0 = 0$, $\{Z_n\}_{n=1}^N$ is a collection of i.i.d. standard Gaussian random variables, and $\{\widehat{x}_n, \widehat{u}_n\}_{n=0}^N$ is an approximation to $\{x(t_n), u(t_n)\}_{n=0}^N$. Of course, more advanced discretization schemes are possible, as explained in Chapter 3. Some of the ideas of Section 10.1.6 are also applicable here, including the observation that the model can be simulated under either the terminal measure or the spot measure to avoid the bias involved in time-discretizing the continuously compounded money market account $e^{\int_0^T r(s) ds}$ in the risk-neutral measure.

13.1.9.4 Single-State Approximations

The extra state variable, and the resultant requirement of a *two*-dimensional PDE scheme (see (13.55)) for an essentially one-factor model, is the price one has to pay for the flexibility of volatility specification in the qG model. This price is relatively modest in practice, but does make the model slightly slower than a classical one-factor short rate model, and also makes it somewhat more challenging to use as a building block for more complicated models, such as equity or FX-linked interest rate hybrid models. In this section we

briefly outline a few ideas for reducing the dimensionality of the model to one state only.

A very simple idea that can be traced to Hagan and Woodward [1999a] (in the multi-dimensional setting) is to force the $y(t)$ variable to be deterministic in the SDE for $x(t)$:

$$dx(t) = (y(t) - \kappa(t)x(t)) dt + \sigma_r(t, x(t)) dW(t), \quad x(0) = 0,$$

where now $y(t)$ is deterministic. Then, using generic machinery of Section 11.3, we can fit $y(t)$ to the initial yield curve via forward induction. Of course, tractability of bond reconstruction formulas is then largely lost, although much of the intuition behind the model is retained and we still maintain separate control over the at-the-money volatility structure (via mean reversion) and the volatility smile (via the local volatility function). Also, the bond reconstruction formulas from Proposition 13.1.1 can be considered as approximations in the new model, and potentially could be used to speed up volatility calibration.

A straight deterministic approximation is a rather blunt tool and would most likely not deliver the level of accuracy we require. A more refined approach for replacing the stochastic variable $y(t)$ in the qG model involves using its *projection* on the variable $x(t)$, as proposed in Kramin [2008]. To develop this idea in a bit more detail, consider a qG model with local volatility (13.10). Focusing first on calculating the following one-dimensional risk-neutral expectation

$$\mathbb{E}(V(x(T))),$$

with the payoff $V(x)$ is a function of $x(T)$ only, the ideas behind Markovian projection (Gyöngy's theorem, see Theorem A.1.1) give us the following exact result.

Proposition 13.1.16. *The undiscounted expected value of a payoff $V(x(T))$ in the model (13.10) is equal to*

$$\mathbb{E}(V(x(T))) = \mathbb{E}(V(\tilde{x}(T))), \quad (13.57)$$

where the process $\tilde{x}(t)$ satisfies

$$d\tilde{x}(t) = (\tilde{y}(t, \tilde{x}(t)) - \kappa(t)\tilde{x}(t)) dt + \tilde{\sigma}_r(t, \tilde{x}(t)) dW(t), \quad (13.58)$$

with

$$\tilde{y}(t, x) = \mathbb{E}(y(t) | x(t) = x), \quad \tilde{\sigma}_r(t, x)^2 = \mathbb{E}(\sigma_r(t, x(t), y(t))^2 | x(t) = x). \quad (13.59)$$

The equality (13.57) does not hold for the more realistic, and useful, case of calculating *discounted* expected values of a payoffs that depend on *both* x and y . Nor is there much theoretical justification for using this

projection when calculating expected values of payoffs that depend on values of the state variables at *multiple times*. Nevertheless, there is some empirical evidence that the approximations work reasonably well in practice. With that in mind, let us define a generic one-state approximate quasi-Gaussian local volatility model by

$$\begin{aligned} dx(t) &= (\tilde{y}(t, x(t)) - \kappa(t)x(t)) dt + \tilde{\sigma}_r(t, x(t)) dW(t), \\ P(t, T) &= \frac{P(0, T)}{P(0, t)} \exp \left(-G(t, T)x(t) - \frac{1}{2}G(t, T)^2 \tilde{y}(t, x(t)) \right), \\ r(t) &= f(0, t) + x(t), \end{aligned} \quad (13.60)$$

with $\tilde{y}(t, x)$, $\tilde{\sigma}_r(t, x)$ given by (13.59). Then, the value $V(t, x)$ of a given security at time t in state x satisfies the following one-dimensional PDE

$$\frac{\partial V}{\partial t} + (\tilde{y}(t, x) - \kappa(t)x) \frac{\partial V}{\partial x} + \frac{1}{2} \tilde{\sigma}_r(t, x)^2 \frac{\partial^2 V}{\partial x^2} = (f(0, t) + x)V. \quad (13.61)$$

Assuming that we can evaluate all terms efficiently, solving the one-dimensional PDE (13.61) is typically quicker than solving the PDE (13.55) for the real model. There is little, if any, benefit in applying the approximation to the Monte Carlo method.

Needless to say, the PDE (13.61) should only be considered for problems inside the domain of applicability of the approximation (13.60). In general, we would expect the approximation to work reasonably well for low to moderate volatilities and maturities (up to, say, 20 years), and deteriorate for longer maturities and/or for large volatilities. Kramin [2008] reports good results across a wide maturity spectrum.

To effectively use (13.60), we need to compute/approximate $\tilde{y}(t, x)$ and $\tilde{\sigma}_r(t, x)$. Typically the volatility term $\sigma_r(t, x, y)$ either does not depend on y at all (see e.g. (13.29)), or depends on y in a close-to-linear fashion, due to the low variance of $y(t)$ compared to $x(t)$. In both cases, the following simple approximation

$$\tilde{\sigma}_r(t, x) = \sigma_r(t, x, \tilde{y}(t, x))$$

appears to be justified.

To calculate $\tilde{y}(t, x)$, we recall the definition of $y(t)$,

$$y(t) = h(t)^2 \int_0^t \sigma_r(s, x(s), y(s))^2 h(s)^{-2} ds.$$

Conditioning on $x(t)$ and replacing $y(s)$ with $\bar{y}(s)$ in the argument of σ_r , where $\bar{y}(\cdot)$ is defined in (13.23), we obtain

$$\tilde{y}(t, x) \approx h(t)^2 \int_0^t E \left(\sigma_r(s, x(s), \bar{y}(s))^2 \middle| x(t) = x \right) h(s)^{-2} ds.$$

Invoking approximate linearity of $\sigma_r(s, x, y)^2$ in x we obtain

$$\tilde{y}(t, x) \approx h(t)^2 \int_0^t \sigma_r(s, \mathbb{E}(x(s)|x(t)=x), \bar{y}(s))^2 h(s)^{-2} ds.$$

Under the Gaussian approximation

$$\mathbb{E}(x(s)|x(t)=x) \approx \frac{\text{Var}(x(s))}{\text{Var}(x(t))} x \approx \frac{h(s)^2 \int_0^s \sigma_r^0(u)^2 h(u)^{-2} du}{h(t)^2 \int_0^t \sigma_r^0(u)^2 h(u)^{-2} du} x = \frac{\bar{y}(s)}{\bar{y}(t)} x,$$

so we obtain

$$\tilde{y}(t, x) \approx h(t)^2 \int_0^t \sigma_r(s, (\bar{y}(s)/\bar{y}(t)) x, \bar{y}(s))^2 h(s)^{-2} ds. \quad (13.62)$$

A direct application of (13.62) is rather costly: with n_t discretized points in t direction and n_x points in x direction, the cost of computing $\tilde{y}(\cdot, \cdot)$ for all t, x on the grid is $O(n_t^2 n_x)$, i.e. higher than for solving the PDE (13.61) itself. One remedy for this issue is to approximate $\sigma_r(s, x, \bar{y}(s))^2$ by a first- or second-order polynomial in x for each s , to obtain a polynomial approximation to $\tilde{y}(s, x)$ with the coefficients computed at $O(n_t)$ cost. Alternatively, we can derive a recursive update equation for $\tilde{y}(t, x)$ with the additional advantage that it does not rely on approximate linearity of $\sigma_r(s, x, y)^2$ in x (unlike (13.62)), a condition that, although generally desirable as we pointed out before, is not necessarily satisfied in all applications. We recall the equation (13.11) satisfied by $y(t)$, and discretize it for the time step $[t_n, t_{n+1}]$. Using short-hand notations $x_n = x(t_n)$, etc., we obtain

$$y_{n+1} = y_n + \left(\sigma_r(t_{n+1}, x_{n+1}, \bar{y}_{n+1})^2 - 2\kappa(t)y_n \right) \Delta_n, \quad , \Delta_n = t_{n+1} - t_n,$$

where $\sigma_r(t, x, \bar{y})^2$ is evaluated at the right point of the interval for reasons that will be clear momentarily. Next, conditioning on x_n, x_{n+1} , we obtain

$$\begin{aligned} \mathbb{E}(y_{n+1}|x_n, x_{n+1}) &= \mathbb{E}(y_n|x_n, x_{n+1}) \\ &\quad + \left(\sigma_r(t_{n+1}, x_{n+1}, \bar{y}_{n+1})^2 - 2\kappa(t)\mathbb{E}(y_n|x_n, x_{n+1}) \right) \Delta_n. \end{aligned}$$

By the Markov property

$$\mathbb{E}(y_n|x_n, x_{n+1}) = \mathbb{E}(y_n|x_n),$$

so

$$\begin{aligned} \mathbb{E}(y_{n+1}|x_n, x_{n+1}) &= \mathbb{E}(y_n|x_n) \\ &\quad + \left(\sigma_r(t_{n+1}, x_{n+1}, \bar{y}_{n+1})^2 - 2\kappa(t)\mathbb{E}(y_n|x_n) \right) \Delta_n, \end{aligned}$$

which gives us a way to obtain $E(y_{n+1}|x_n, x_{n+1})$ from $E(y_n|x_n)$. To get at the quantity that we want, namely $E(y_{n+1}|x_{n+1})$, we average over x_n ,

$$\begin{aligned} E(y_{n+1}|x_{n+1}) &= \int E(y_{n+1}|x_n = x, x_{n+1}) Q(x_n \in dx|x_{n+1}) \\ &= \int \left(E(y_n|x_n = x) \right. \\ &\quad \left. + \left(\sigma_r(t_{n+1}, x_{n+1}, \bar{y}_{n+1})^2 - 2\kappa(t)E(y_n|x_n = x) \right) \Delta_n \right) \\ &\quad \times Q(x_n \in dx|x_{n+1}) \end{aligned}$$

and, after rearranging some terms, we obtain a recursive formula for $E(y_{n+1}|x_{n+1})$,

$$\begin{aligned} E(y_{n+1}|x_{n+1}) &= (1 - 2\kappa(t)\Delta_n) \int E(y_n|x_n = x) Q(x_n \in dx|x_{n+1}) \\ &\quad + \sigma_r(t_{n+1}, x_{n+1}, \bar{y}_{n+1})^2 \Delta_n. \end{aligned}$$

For small Δ_n , the density

$$Q(x_n \in dx|x_{n+1})$$

is approximately Gaussian, and the required integral can be quickly computed numerically with just a few terms, giving us an algorithm of numerical complexity $O(n_t n_x)$.

Finally, we point out that the model (13.60) could be made exactly arbitrage-free by introducing a time-dependent deterministic component in its drift that is fit numerically to the initial yield curve, in line with the discussion at the beginning of this section — but of course at the cost of losing analytical tractability.

13.2 One-Factor Quasi-Gaussian Model with Stochastic Volatility

The most general one-factor quasi-Gaussian model specification allows for the short rate volatility to be a stochastic process, see (13.2). While so far we only considered the case of deterministic dependence of the volatility on state variables of the model, we now proceed to generalize the setup to include stochastic volatility.

13.2.1 Definition

Introduction of a stochastic variance process (see Chapter 8) in the specification of $g(\cdot)$ in (13.2) leads to a *stochastic volatility quasi-Gaussian* model. In particular, defining $z(t)$ to be the familiar CIR process,

$$dz(t) = \theta(z_0 - z(t)) dt + \eta(t)\sqrt{z(t)} dZ(t), \quad \langle dZ(t), dW(t) \rangle = 0,$$

we obtain a stochastic volatility qG model by specifying the volatility structure of the form

$$g(t, \omega) = \sqrt{z(t)} g(t, x(t), y(t)), \quad (13.63)$$

where $g(t, x, y)$ is a function of t, x, y only. With the standard definition

$$\sigma_r(t, x, y) = g(t, x, y) h(t),$$

the model is defined by the collection of SDEs

$$\begin{aligned} dx(t) &= (y(t) - \varkappa(t)x(t)) dt + \sqrt{z(t)}\sigma_r(t, x(t), y(t)) dW(t), \\ dy(t) &= \left(z(t)\sigma_r(t, x(t), y(t))^2 - 2\varkappa(t)y(t) \right) dt, \\ dz(t) &= \theta(z_0 - z(t)) dt + \eta(t)\sqrt{z(t)} dZ(t), \end{aligned} \quad (13.64)$$

subject to

$$x(0) = y(0) = 0, \quad z(0) = z_0 = 1, \quad \langle dZ(t), dW(t) \rangle = 0.$$

When specifying the local volatility function, it was natural to use piecewise constant functions for various parameters (see (13.39)), and we do the same with the volatility of variance,

$$\eta(t) = \sum_{n=1}^{N-1} \eta_n 1_{\{t \in (T_{n-1}, T_n]\}}. \quad (13.65)$$

The bond reconstitution formulas in the model (13.64) are the same as for the local volatility case; as follows from Proposition 13.1.1, they are the same for *any* quasi-Gaussian model. In particular, the zero-coupon discount bond formulas do not depend on the stochastic volatility process $z(t)$, and thus the model is a “true” stochastic volatility model, i.e. its stochastic volatility is unspanned and cannot be hedged by discount bonds. We remind the reader of the discussion of this topic in Section 11.2.3 and note that the model (13.64) has the lowest possible number of state variables — three — for an unspanned stochastic volatility term structure model, see Collin-Dufresne and Goldstein [2002b].

We note in passing that the assumption of zero correlation $\langle dZ(t), dW(t) \rangle = 0$ is a technical restriction helpful for developing efficient calibration formulas. It does not, however, restrict the range of available volatility smiles in the model, as the local volatility term can be used to control the slope of the smile. See also our discussion in Section 13.2.5.

13.2.2 Swap Rate Dynamics

Many results obtained in the local volatility case extend naturally to incorporate stochastic volatility, including Proposition 13.1.2, Lemma 13.1.3 and Proposition 13.1.8. The following analog to Corollary 13.1.9 is particularly useful for calibration.

Proposition 13.2.1. *Under the assumption of linear local short rate volatility (13.29), the dynamics of a swap rate $S(t)$ in the stochastic volatility quasi-Gaussian model (13.64) are given approximately by*

$$\begin{aligned} dS(t) &= \sqrt{z(t)} \lambda_S(t) (b_S(t)S(t) + (1 - b_S(t)) S(0)) dW^A(t), \\ dz(t) &= \theta(z_0 - z(t)) dt + \eta(t) \sqrt{z(t)} dZ(t), \\ z(0) = z_0 &= 1, \quad \langle dZ(t), dW^A(t) \rangle = 0, \end{aligned}$$

where $\lambda_S(t)$ and $b_S(t)$ are given by (13.34)–(13.35).

The dynamics in Proposition 13.2.1 are easily recognized to be those of a time-dependent SV model, see Chapter 9. As time averaging methods are available for stochastic volatility models (see Section 9.3), the following proposition, an analog to Proposition 13.1.10, should not come as a surprise. See also Theorem 9.3.1, Corollary 9.3.5, and Theorem 9.3.6.

Proposition 13.2.2. *In the setting of the stochastic volatility quasi-Gaussian model (13.64) with the linear local volatility (13.29), consider a T -maturity swaption on a swap rate $S(T)$. For the purpose of European option pricing, the dynamics of $S(t)$ in its annuity measure can be approximated by the following time-homogeneous stochastic volatility model,*

$$\begin{aligned} dS(t) &= \sqrt{z(t)} \bar{\lambda}_S (\bar{b}_S S(t) + (1 - \bar{b}_S(t)) S(0)) dW^A(t), \\ dz(t) &= \theta(z_0 - z(t)) dt + \bar{\eta}_S \sqrt{z(t)} dZ(t), \end{aligned}$$

where

- The effective volatility of variance $\bar{\eta}_S$ is given by

$$\bar{\eta}_S^2 = \frac{\int_0^T \eta(t)^2 \rho_S(t) dt}{\int_0^T \rho_S(t) dt}, \quad (13.66)$$

with the weight function $\rho_S(t)$ given by

$$\rho_S(r) = \int_r^T \int_s^T \lambda_S(t)^2 \lambda_S(s)^2 e^{-\theta(t-s)} e^{-2\theta(s-r)} dt ds.$$

- The effective skew \bar{b}_S is given by

$$\bar{b}_S = \int_0^T b_S(t) w_S(t) dt, \quad (13.67)$$

with the weight function $w_S(t)$ given by

$$w_S(t) = \frac{v_S(t)^2 \lambda_S(t)^2}{\int_0^T v_S(u)^2 \lambda_S(u)^2 du},$$

$$v_S(t)^2 = z_0^2 \int_0^t \lambda_S(s)^2 ds + z_0 e^{-\theta t} \int_0^t \lambda_S(s)^2 e^{-\theta s} \int_0^s \eta(u)^2 e^{2\theta u} du ds.$$

- The effective volatility $\bar{\lambda}_S$ is given by the solution to the equation

$$\Psi_{\bar{z}} \left(\frac{h''(\zeta_S)}{h'(\zeta_S)} \bar{\lambda}_S^2, 0; T \right) = \Psi_{\bar{z}\bar{\lambda}^2} \left(\frac{h''(\zeta_S)}{h'(\zeta_S)}, 0; T \right), \quad (13.68)$$

where (see Theorem 9.3.1)

$$\zeta_S = z_0 \int_0^T \lambda_S(t)^2 dt,$$

$$\Psi_{\bar{z}\bar{\lambda}^2}(v, 0; T) = E \left(\exp \left(v \int_0^T \lambda_S(t)^2 z(t) dt \right) \right),$$

$$\Psi_{\bar{z}}(v, 0; T) = E \left(\exp \left(v \int_0^T z(t) dt \right) \right),$$

$$h(x) = \frac{S(0)}{\bar{b}_S} (2\Phi(\bar{b}_S \sqrt{x}/2) - 1).$$

The functions $\lambda_S(t)$, $b_S(t)$ are given by (13.34)–(13.35).

13.2.3 Volatility Calibration

With a swaption strip $\{S_n(\cdot)\}_{n=1}^{N-1}$ given as in Sections 13.1.6 and 13.1.7, the volatility calibration algorithm can proceed along the same principles as in Section 13.1.7, where swap rate distributions in the stochastic volatility qG model can be found from the constant-parameter displaced SV SDEs in Proposition 13.2.2. As before, we assume that a collection of market parameters $(\hat{\lambda}_{S_n}, \hat{b}_{S_n}, \hat{\eta}_{S_n})$, $n = 1, \dots, N - 1$, is given; in practice, these parameters may be obtained by fitting a vanilla SV model to swaptions of a given expiry/tenor across strikes, a procedure described in more detail in Section 16.1.4.

While it is easy to modify the algorithm of Section 13.1.7 to introduce one more variable to solve for (η_n) in Step 6 for each n , the calibration algorithm

is typically more stable if we first solve for the volatility of variance function $\eta(t)$, i.e. find $\{\eta_n^*\}$ for all n , and then follow the algorithm from Section 13.1.7 to solve recursively in n for (λ_n^*, b_n^*) , using slightly modified formulas from Proposition 13.1.10 as given in Proposition 13.2.2. For completeness, we repeat the algorithm with the necessary modifications.

1. Set (λ_n, b_n) , $n = 1, \dots, N - 1$, to some reasonable starting values, e.g. set λ_n 's to (properly scaled) volatilities obtained by calibrating a pure Gaussian model as in Section 10.1.4, and $b_n = \hat{b}_{S_n}$.
2. Solve for η_n^* for $n = 1, \dots, N - 1$, using (13.66). For weights $\rho_{S_n}(t)$, use $\lambda_{S_n}(t)$ as computed from the first guess for λ_n 's (using (13.34)) obtained on Step 1.
3. Set $n = 1$.
4. For given n , (λ_i^*, b_i^*) are known for $i = 1, \dots, n - 1$. Note we use a star to denote *calibrated* values of the model parameters.
5. Calculate $\bar{x}(t)$ (Lemma 13.1.6) and $\bar{y}(t)$ (Proposition 13.1.4) for $t \in [0, T_n]$ using (λ_i^*, b_i^*) , $i = 1, \dots, n - 1$, and the initial guess for (λ_n, b_n) from Step 1. Note that $\bar{x}(t), \bar{y}(t)$ implicitly depend on n as their definition depends on the swap rate/annuity measure used.
6. Calculate $\lambda_{S_n}(t), b_{S_n}(t)$ for $t \in [0, T_{n-1}]$ from (λ_i^*, b_i^*) , $i = 1, \dots, n - 1$, using (13.34)–(13.35).
7. Make another guess for (λ_n, b_n) .
8. Update $\lambda_{S_n}(t), b_{S_n}(t)$ for $t \in (T_{n-1}, T_n]$ from (λ_i^*, b_i^*) , $i = 1, \dots, n - 1$, using (13.34)–(13.35).
9. Calculate $\bar{\lambda}_{S_n}, \bar{b}_{S_n}$ using Proposition 13.2.2.
10. Compare $(\bar{\lambda}_{S_n}, \bar{b}_{S_n})$ to $(\hat{\lambda}_{S_n}, \hat{b}_{S_n})$. If not equal within given tolerance, go to Step 7. Otherwise, proceed to Step 11.
11. As we have reached acceptable convergence between $(\bar{\lambda}_{S_n}, \bar{b}_{S_n})$ and $(\hat{\lambda}_{S_n}, \hat{b}_{S_n})$, set the calibrated model parameter values to the latest trial values, $(\lambda_n^*, b_n^*) = (\lambda_n, b_n)$.
12. Update $n \rightarrow n + 1$. If $n \leq N - 1$ go to Step 4. Otherwise, conclude.

As in Section 13.1.7, Step 5 can be performed inside of the calibration loop with some positive impact on the quality of calibration at a (moderate) cost of extra complexity.

13.2.4 Mean Reversion Calibration

The stochastic variance scaling $z(t)$ has some impact on the inter-temporal correlations or volatility ratios of swaptions of different maturities, an effect we discuss further in Section 20.2.4. Still, as mean reversion calibration is not meant to be overly precise, we can continue to use formulas developed in Section 13.1.8 that, in the context of a stochastic volatility qG model, imply that we roundly ignore any such effects.

While on the topic of mean reversions, let us not forget the parameter θ , the mean reversion of the variance process. The parameter θ fundamentally determines the decay of volatility smile curvature as a function of option maturity T ; a good setting of this parameter will help keep the parameter $\eta(t)$ from depending excessively on calendar time t . In general, we inherit the same θ as used for the vanilla SV model calibration to European swaptions, a subject we discuss in depth in Section 16.1.4.

13.2.5 Non-Zero Correlation

A question sometimes arises whether it is too restrictive to assume (as we do) that the correlation between the Brownian motions driving the curve factor and the stochastic volatility is zero. Empirical evidence from all major fixed income markets generally suggests that correlations between interest rates and their (short-dated) volatilities are small; see e.g. the analysis in Chen and Scott [2001]. Moreover, the assumption of zero correlation has little impact on the range of volatility smiles that the model can produce, as the skew term in the local volatility can produce the necessary tilting of the smile. Nor does it affect hedging implications of the model as long as minimum variance hedging is employed; see our discussion in Section 8.9. In our view, the zero-correlation constraint has little consequence in practice, but brings substantial technical benefits, particularly the ability to shift pricing numeraires without affecting the form of the stochastic variance process. If desired, non-zero correlation can still be accommodated, as most numerical schemes — including averaging formulas — are easy to adapt, as we briefly explain in Remark 9.3.7. One should be mindful, however, of the fact that under non-zero correlation, measure changes introduce a rate-dependent term in the drift of the stochastic variance (see Proposition 8.3.9), which requires additional considerations when deriving approximations for European swaptions, say. If interested, the reader can attack the problem along the same lines as in Chapter 15, where the case of non-zero correlation is considered in a context of a different model (the Libor market model).

13.2.6 PDE and Monte Carlo Methods

The PDE method of Section 13.1.9.2 can be extended to cover the stochastic volatility case, using the techniques of Section 9.4 for the stochastic volatility part. The resulting PDE will involve three spatial dimensions, which can be handled by the Craig-Sneyd scheme. The same is true for the Monte Carlo method, where a combination of ideas from Sections 13.1.9.3 and 9.5 cover most practical issues.

13.3 Multi-Factor Quasi-Gaussian Model

13.3.1 General Multi-Factor Model

Multi-factor quasi-Gaussian models combine the flexibility of volatility specification of multi-factor models (see Chapter 12) with the ability to generate volatility skews and smiles. Practical multi-factor quasi-Gaussian models are relatively new (see Andreasen [2005]), but could provide a compelling alternative to the Libor market model in Chapter 14⁷, the current *de-facto* market standard for multi-factor models.

Following the steps that lead to the one-factor quasi-Gaussian model, the multi-factor qG model is obtained by imposing a separability condition on the volatility structure of a multi-factor HJM model. Specifically, let us consider the forward rate process

$$df(t, T) = \sigma_f(t, T, \omega)^\top \left(\left(\int_t^T \sigma_f(t, u, \omega) du \right) dt + dW(t) \right), \quad (13.69)$$

where $\sigma_f(t, T, \omega)$ is a d -dimensional stochastic process, and $W(t)$ a d -dimensional Brownian motion in the risk-neutral measure. Let us assume that $\sigma_f(t, T, \omega)$ is separable, in the sense that it can be written as

$$\sigma_f(t, T, \omega) = g(t, \omega)h(T), \quad (13.70)$$

where $g(t, \omega)$ is a $d \times d$ stochastic matrix-valued process, and $h(t)$ is a d -dimensional deterministic vector-valued function of time. Then we define

$$H(t) = \text{diag}(h(t)) = \begin{pmatrix} h_1(t) & 0 & \ddots & 0 \\ 0 & h_2(t) & \ddots & \ddots \\ \ddots & \ddots & \ddots & 0 \\ 0 & \ddots & 0 & h_d(t) \end{pmatrix}.$$

Let us assume further that $h_i(t) \neq 0$, $i = 1, \dots, d$, for all t , whereby $H(t)$ is then invertible, and so we can define a diagonal $d \times d$ matrix $\varkappa(t)$ by

$$\varkappa(t) = -\frac{dH(t)}{dt} H(t)^{-1} \quad (13.71)$$

(this is the same as in (12.7)–(12.8)). Moreover, let us define

$$G(t, T) = \int_t^T H(u)H(t)^{-1}\mathbf{1} du, \quad \sigma_r(t, \omega) = g(t, \omega)H(t),$$

where we use the notation $\mathbf{1} = (1, 1, \dots, 1)^\top$ from Section 12.1.1.1.

⁷For readers not fully familiar with Libor market models, we recommend reading Chapters 14 and 15 before proceeding with this section.

Proposition 13.3.1. Consider a general multi-factor quasi-Gaussian model, i.e. an HJM model (13.69) with the separable volatility condition (13.70). Define stochastic processes $x(t)$, $y(t)$ by

$$\begin{aligned} dx(t) &= (y(t)\mathbf{1} - \varkappa(t)x(t)) dt + \sigma_r(t, \omega)^\top dW(t), \\ dy(t) &= (\sigma_r(t, \omega)^\top \sigma_r(t, \omega) - \varkappa(t)y(t) - y(t)\varkappa(t)) dt, \end{aligned} \quad (13.72)$$

where $x(t) \in \mathbb{R}^d$, $y(t) \in \mathbb{R}^{d \times d}$, and $x(0) = 0$, $y(0) = 0$. Then, zero-coupon discount bonds are given by

$$P(t, T) = P(t, T, x(t), y(t)),$$

with

$$P(t, T, x, y) = \frac{P(0, T)}{P(0, t)} \exp \left(-G(t, T)^\top x - \frac{1}{2} G(t, T)^\top y G(t, T) \right).$$

In addition, the instantaneous forward rates are given by

$$f(t, T) = f(0, T) + \mathbf{1}^\top H(T)H(t)^{-1}(x(t) + y(t)G(t, T)), \quad (13.73)$$

with the short rate

$$r(t) = f(0, t) + \mathbf{1}^\top x(t).$$

Proof. Follows closely that of Proposition 12.1.2. \square

13.3.2 Local and Stochastic Volatility Parameterization

While a pure local volatility specification of the multi-factor qG model is certainly possible, for brevity let us proceed directly to a more general setting where we have both local and stochastic volatility. Following Section 13.2, we start by specifying a one-dimensional process $z(t)$ by

$$dz(t) = \theta(z_0 - z(t)) dt + \eta(t)\sqrt{z(t)} dZ(t), \quad z(0) = z_0 = 1, \quad (13.74)$$

with $\langle dZ(t), dW(t) \rangle = 0$. Inspired by the one-dimensional case, we would like to specify a model with the volatility structure of the type

$$\sigma_r(t, \omega)^\top = \sqrt{z(t)} \sigma_x(t, x(t), y(t))^\top, \quad (13.75)$$

where $\sigma_x(t, x, y)$ is a multi-dimensional local volatility function responsible for inducing the skews in volatility smiles of swaptions. However, it is not entirely obvious how to parametrize $\sigma_x(t, x, y)$ sensibly, as the volatility function is not only responsible for skews but also for the general volatility structure of the model, including volatilities and correlations of all the rates.

Fortunately, the ideas of Section 12.1.7 could be fruitfully recycled and extended here, as suggested by Andreasen [2005] (see also Cheyette [1991]).

Recalling the definition of benchmark rates from Section 12.1.7, we take d benchmark tenors $\delta_1 < \dots < \delta_d$, and define d “rolling” benchmark rates $f_i(t) = f(t, t + \delta_i)$, $i = 1, \dots, d$. Ideally, it would be convenient if the qG model specification was such that the dynamics of the benchmark rates $f_i(t)$, $i = 1, \dots, d$, were of the familiar form

$$df_i(t) = \sqrt{z(t)} \lambda_i^f(t) \left(\alpha_i^f(t) + b_i^f(t) f_i(t) \right) dU_i(t) + O(dt), \quad i = 1, \dots, d, \quad (13.76)$$

where $\{U_i(t)\}_{i=1}^d$ is a d -dimensional vector of Brownian motions with the correlation matrix $X^f(t) = \{\chi_{i,j}(t)\}$. The following proposition shows how model parameters need to be set for these dynamics to hold.

Proposition 13.3.2. *Let us define the $d \times d$ matrix-valued process H^f by*

$$H^f(t) = \begin{pmatrix} h(t + \delta_1)^\top \\ \vdots \\ h(t + \delta_d)^\top \end{pmatrix},$$

and σ^f by

$$\begin{aligned} & \sigma^f(t, f(t)) \\ &= \text{diag} \left((\lambda_1^f(t)(\alpha_1^f(t) + b_1^f(t)f_1(t)), \dots, \lambda_d^f(t)(\alpha_d^f(t) + b_d^f(t)f_d(t)))^\top \right), \end{aligned}$$

where $f(t) = (f_1(t), \dots, f_d(t))^\top$. Also, let $D^f(t)$ be specified by $X^f(t) = D^f(t)^\top D^f(t)$. In (13.72), let us set

$$\begin{aligned} \sigma_r(t, \omega)^\top &= \sqrt{z(t)} \sigma_x(t, x(t), y(t))^\top, \\ \sigma_x(t, x(t), y(t))^\top &= H(t) H^f(t)^{-1} \sigma^f(t, f(t)) D^f(t)^\top, \end{aligned} \quad (13.77)$$

where σ_x is a function of $x(t)$, $y(t)$ because f_i 's are, see (13.73). Then the qG model in Proposition 13.3.1 is consistent with the benchmark rate dynamics of equation (13.76).

Proof. From (13.73) and using the definition of the vector $f(t)$, we obtain

$$f(t) = f(0) + H^f(t) H(t)^{-1} (x(t) + y(t) G(t, T)).$$

Then, with the help of (13.72),

$$\begin{aligned} df(t) &= O(dt) + H^f(t) H(t)^{-1} dx(t) \\ &= O(dt) + H^f(t) H(t)^{-1} \sigma_r(t, \omega)^\top dW(t). \end{aligned}$$

Using (13.77), we obtain

$$\begin{aligned} df(t) &= O(dt) + \sqrt{z(t)} H^f(t) H(t)^{-1} H(t) H^f(t)^{-1} \sigma^f(t, f(t)) D^f(t)^\top dW(t) \\ &= O(dt) + \sqrt{z(t)} \sigma^f(t, f(t)) D^f(t)^\top dW(t). \end{aligned}$$

The statement of the proposition follows once we set $dU(t) = D^f(t)^\top dW(t)$.

□

With the parameterization outlined in Proposition 13.3.2, the benchmark rates follow “SV-like” dynamics (13.76), and we can reasonably expect Libor and swap rates to follow similar dynamics, at least approximately. Besides reducing the generic qG model to a familiar class of dynamics, the parameterization in Proposition 13.3.2 also achieves a clear distinction between the effects of the various model parameters: we now have volatility parameters $\{\lambda_i^f(t)\}$, rate correlation parameters $\{\chi_{i,j}(t)\}$, skew parameters $\{b_i^f(t)\}$, and the volatility of variance $\eta(t)$.

As was the case for the model in Section 12.1.7, the qG model above has enough flexibility to calibrate to d swaption strips, if we assume that the mean reversions $\varkappa(t)$ and the correlation matrix $X^f(t)$ are specified prior to volatility calibration. If the swaption strips are of constant-tenor type — a sensible choice for $d > 1$ — then it is natural to set the tenors of the swaptions we want to use in calibration equal to the benchmark tenors δ_i , $i = 1, \dots, d$. With the stochastic volatility parameterization (13.74), (13.77), we can calibrate

1. At-the-money volatilities for d swaption strips.
2. Slopes of the volatility smiles for d swaption strips.
3. Curvature of the smile for one swaption strip.

The last point is not as restrictive as it may appear, since the curvatures of the smiles for swaptions of the same expiry but different tenors tend to be fairly similar.

To parametrize the model in a suitable way, let us choose a tenor structure $0 = T_0 < \dots < T_N$ and denote the swap rates/annuities for the i -th swaption strip, $i = 1, \dots, d$, in the calibration set by

$$S_{n,\mu_i(n)}(t, x(t), y(t)), \quad A_{n,\mu_i(n)}(t, x(t), y(t)), \quad n = 1, \dots, N - 1.$$

Extending (13.39), it is natural to define, for $i = 1, \dots, d$,

$$\begin{aligned} \lambda_i^f(t) &= \sum_{n=1}^{N-1} \lambda_{i,n} 1_{\{t \in (T_{n-1}, T_n]\}}, \quad \alpha_i^f(t) = \sum_{n=1}^{N-1} S_{n,\mu_i(n)}(0) 1_{\{t \in (T_{n-1}, T_n]\}}, \\ b_i^f(t) &= \sum_{n=1}^{N-1} b_{i,n} D_{i,n} 1_{\{t \in (T_{n-1}, T_n]\}}, \quad \eta(t) = \sum_{n=1}^{N-1} \eta_n 1_{\{t \in (T_{n-1}, T_n]\}}, \end{aligned}$$

where the skew scalings $D_{i,n}$ are given by an approximate (as we ignore ∂y terms) derivative of $S_{n,\mu_i(n)}$ “in the direction” of f_i ,

$$D_{i,n} = \mathbf{1}^\top H(t + \delta_i) H(t)^{-1} (\nabla S_{n,\mu_i(n)}) ,$$

with

$$\nabla S \triangleq (\partial S / \partial x_1, \dots, \partial S / \partial x_d).$$

In summary, the volatility smile parameters in the model are the $(2d + 1) \times (N - 1)$ parameters $\{\lambda_{i,n}\}$, $\{b_{i,n}\}$ and $\{\eta_n\}$.

13.3.3 Swap Rate Dynamics and Approximations

The swap rate dynamics in the multi-factor qG model can be derived and simplified by techniques similar to those we applied earlier in the one-factor case. First, we establish a multi-dimensional counterpart to Proposition 13.1.2, i.e. the exact dynamics of a given swap rate $S(t)$ in the annuity measure corresponding to its annuity $A(t)$, see (13.13)–(13.14).

Proposition 13.3.3. *In a multi-factor stochastic volatility quasi-Gaussian model with volatility parameterization (13.75), the dynamics of a swap rate $S(t)$ defined by (13.13) are given by*

$$dS(t) = \sqrt{z(t)} \left((\nabla S) \sigma_x^\top \sigma_x (\nabla S)^\top \right)^{1/2} dU^A(t), \quad (13.78)$$

where all functions are understood to be evaluated at $(t, x(t), y(t))$, and $U^A(t)$ is a one-dimensional Brownian motion in the annuity measure Q^A .

Proof. By standard arguments of using Ito's lemma on $S(t)$ and a martingale property of $S(t)$. \square

Using the Markovian projection method (see Appendix A), the dynamics of (13.78) can be approximated, for the purposes of pricing European options, with

$$dS(t) = \sqrt{z(t)} \varphi(t, S(t)) dU^A(t), \quad (13.79)$$

where

$$\varphi(t, s)^2 = E^A \left((\nabla S) c_x (\nabla S)^\top \middle| S(t) = s \right),$$

and we have denoted

$$c_x = c_x(t, x, y) = \sigma_x(t, x, y)^\top \sigma_x(t, x, y). \quad (13.80)$$

We expect the local volatility term in (13.79) to mostly control the slope of the volatility smile, hence we look for a linear approximation to $\varphi(t, s)$ in s . First, we need to choose the point for expansion. One can use $x(t) = 0$, $y(t) = 0$ as a decent choice; however, using $E^A(x(t))$, $E^A(y(t))$ or approximations thereof is, as always, preferable.

Proposition 13.3.4. *Let*

$$\sigma_x^0(t) = \sigma_x(t, 0, 0). \quad (13.81)$$

Then

$$\mathbb{E}^A(y(t)) \approx \bar{y}(t) \triangleq H(t) \left(\int_0^t H(s)^{-1} \sigma_x^0(s)^\top \sigma_x^0(s) H(s)^{-1} ds \right) H(t).$$

An approximation $\bar{x}(t)$ to $\mathbb{E}^A(x(t))$ is given by

$$\bar{x}(t) = H(t) \left(\int_0^t H(s)^{-1} (\bar{y}(s) \mathbf{1} - \sigma_x^0(s)^\top \sigma_x^0(s) G_A(s)) ds \right), \quad t \in [0, T_1],$$

where (recall that $S(t) = S_{0,N}(t)$, $A(t) = A_{0,N}(t)$)

$$G_A(s) = \frac{1}{A(0)} \sum_{n=0}^{N-1} \tau_n P(0, T_{n+1}) G(s, T_{n+1}).$$

Proof. The result for $\mathbb{E}^A(y(t))$ follows after approximating $\sigma_r(t, \omega)$ in (13.72) by $\sigma_x^0(t)$ and then proceeding as in the proof of Proposition 13.1.4 (see also Section 12.1.1.1).

Using the same approximation and replacing $y(t)$ with $\bar{y}(t)$, we obtain from (13.72) the following SDE for $x_g(t)$, a Gaussian approximation to $x(t)$,

$$dx_g(t) = (\bar{y}(t) \mathbf{1} - \varkappa(t) x_g(t)) dt + \sigma_x^0(t)^\top dW(t).$$

For a given $T > 0$, in the T -forward measure,

$$dx_g(t) = (\bar{y}(t) \mathbf{1} - \sigma_x^0(t)^\top \sigma_x^0(t) G(t, T) - \varkappa(t) x_g(t)) dt + \sigma_x^0(t)^\top dW^T(t),$$

where dW^T is a driftless Brownian motion; hence

$$\mathbb{E}^T(x_g(t)) = H(t) \left(\int_0^t H(s)^{-1} (\bar{y}(s) \mathbf{1} - \sigma_x^0(s)^\top \sigma_x^0(s) G(s, T)) ds \right).$$

Then,

$$\begin{aligned} \mathbb{E}^A(x(t)) &\approx \mathbb{E}^A(x_g(t)) \\ &= \frac{1}{A(0)} \mathbb{E}(\beta(t)^{-1} A(t) x_g(t)) \\ &= \frac{1}{A(0)} \sum_{n=0}^{N-1} \tau_n \mathbb{E}(x_g(t) \beta(t)^{-1} P(t, T_{n+1})) \\ &= \sum_{n=0}^{N-1} \frac{\tau_n P(0, T_{n+1})}{A(0)} \mathbb{E}^{T_{n+1}}(x_g(t)) \\ &= \sum_{n=0}^{N-1} \frac{\tau_n P(0, T_{n+1})}{A(0)} H(t) \\ &\quad \times \int_0^t H(s)^{-1} (\bar{y}(s) \mathbf{1} - \sigma_x^0(s)^\top \sigma_x^0(s) G(s, T_{n+1})) ds, \end{aligned}$$

and the result follows. \square

In preparation for our next result, let us define $S_g(t)$ to be the Gaussian approximation to $S(t)$, i.e. a process with the dynamics given by (13.78) where all functions are evaluated at $(t, 0, 0)$ and there is no stochastic volatility,

$$dS_g(t) = \left. \left((\nabla S) (\sigma_x^0)^\top \sigma_x^0 (\nabla S)^\top \right)^{1/2} \right|_{(t, 0, 0)} dU^A(t).$$

The dynamics of $S_g(t)$ can alternatively be represented using a multi-dimensional Brownian motion,

$$dS_g(t) = \nabla S(t, 0, 0) (\sigma_x^0)^\top dW^A(t). \quad (13.82)$$

Theorem 13.3.5. *Let $\bar{x}(t)$ and $\bar{y}(t)$ be as given in Proposition 13.3.4, and let $c_x(t)$ be defined as in (13.80). For pricing European swaptions, the dynamics of the swap rate $S(t)$ (defined by (13.13)) in the multi-factor quasi-Gaussian model with the volatility structure given by (13.77) can be approximated by the following time-dependent Heston dynamics,*

$$dS(t) \approx \sqrt{z(t)} \lambda_S(t) (b_S(t)S(t) + (1 - b_S(t)) S(0)) dU^A(t),$$

where

$$\begin{aligned} \lambda_S(t) &= \frac{1}{S(0)} \left. \left((\nabla S) c_x (\nabla S)^\top \right)^{1/2} (t, \bar{x}(t), \bar{y}(t)) \right., \\ b_S(t) &= S(0) \left. \left(\frac{1}{2} \frac{(\nabla S) d_x (\nabla S)^\top}{(\nabla S) (c_x (\nabla S)^\top (\nabla S) c_x) (\nabla S)^\top} \right. \right. \\ &\quad \left. \left. + \frac{(\nabla S) (c_x (\nabla^2 S) c_x) (\nabla S)^\top}{(\nabla S) (c_x (\nabla S)^\top (\nabla S) c_x) (\nabla S)^\top} \right) \right|_{(t, \bar{x}(t), \bar{y}(t))}. \end{aligned}$$

Here we have denoted

$$d_x = \sum_{l=1}^d \left(c_x (\nabla S)^\top \right)_l \left(\frac{\partial \sigma_x^\top}{\partial x_l} \sigma_x + \sigma_x^\top \frac{\partial \sigma_x}{\partial x_l} \right),$$

$$\begin{aligned} \frac{\partial \sigma_x^\top}{\partial x_l} &= H(t) H^f(t)^{-1} \text{diag} \left(\lambda_1^f(t) b_1^f(t) \frac{h_l(t + \delta_1)}{h_l(t)}, \dots \right. \\ &\quad \left. \dots, \lambda_d^f(t) b_d^f(t) \frac{h_l(t + \delta_d)}{h_l(t)} \right) D^f(t)^\top, \quad (13.83) \end{aligned}$$

and

$$\nabla^2 S \triangleq \left\{ \frac{\partial^2 S}{\partial x_i \partial x_j} \right\}_{i,j=1}^d.$$

Proof. In (13.79), consider the conditional expected value

$$\mathbb{E}^A \left((\nabla S) c_x (\nabla S)^\top \middle| S(t) = s \right).$$

Expanding the integrand $(\nabla S) c_x (\nabla S)^\top$ around $(t, \bar{x}(t), \bar{y}(t))$ to first order in x , we obtain

$$\begin{aligned} \left((\nabla S) c_x (\nabla S)^\top \right) (t, x(t), y(t)) &\approx \left((\nabla S) c_x (\nabla S)^\top \right) (t, \bar{x}(t), \bar{y}(t)) \\ &\quad + \nabla \left((\nabla S) c_x (\nabla S)^\top \right) \Big|_{(t, \bar{x}(t), \bar{y}(t))} (x(t) - \bar{x}(t)). \end{aligned}$$

Then

$$\begin{aligned} \mathbb{E}^A \left((\nabla S) c_x (\nabla S)^\top \middle| S(t) = s \right) &\approx \left((\nabla S) c_x (\nabla S)^\top \right) (t, \bar{x}(t), \bar{y}(t)) \\ &\quad + \nabla \left((\nabla S) c_x (\nabla S)^\top \right) \Big|_{(t, \bar{x}(t), \bar{y}(t))} \mathbb{E}^A (x(t) - \bar{x}(t) | S(t) = s). \end{aligned}$$

Using a Gaussian approximation for the conditional expected value,

$$\begin{aligned} \mathbb{E}^A (x(t) - \bar{x}(t) | S(t) = s) &\approx \mathbb{E}^A (x_g(t) - \bar{x}(t) | S_g(t) = s) \\ &= \frac{\text{Cov}(S_g(t), x_g(t))}{\text{Var}(S_g(t))} (s - S(0)) \\ &= \frac{(\sigma_x^0)^\top \sigma_x^0 (\nabla S(t, 0, 0))^\top}{(\nabla S(t, 0, 0)) (\sigma_x^0)^\top \sigma_x^0 (\nabla S(t, 0, 0))^\top} (s - S(0)) \\ &\approx \frac{c_x (\nabla S)^\top}{(\nabla S) c_x (\nabla S)^\top} \Big|_{(t, \bar{x}(t), \bar{y}(t))} (s - S(0)). \end{aligned}$$

For any l , $l = 1, \dots, d$, we have

$$\begin{aligned} \frac{\partial}{\partial x_l} \left((\nabla S) c_x (\nabla S)^\top \right) &= ((\nabla^2 S)_l)^\top c_x (\nabla S)^\top + (\nabla S) \left(\frac{\partial c_x}{\partial x_l} \right) (\nabla S)^\top + (\nabla S) c_x (\nabla^2 S)_l \\ &= (\nabla S) \left(\frac{\partial \sigma_x^\top}{\partial x_l} \sigma_x + \sigma_x^\top \frac{\partial \sigma_x}{\partial x_l} \right) (\nabla S)^\top + 2(\nabla S) c_x (\nabla^2 S)_l, \end{aligned}$$

where $(\nabla^2 S)_l$ is the l -th column of the matrix $\nabla^2 S$. Thus

$$\begin{aligned} \nabla \left((\nabla S) c_x (\nabla S)^\top \right) c_x (\nabla S)^\top &= (\nabla S) d_x (\nabla S)^\top + 2(\nabla S) (c_x (\nabla^2 S) c_x) (\nabla S)^\top, \end{aligned}$$

where the matrix d_x is given in the statement of the theorem, and

$$\begin{aligned} & \nabla \left((\nabla S) c_x (\nabla S)^\top \right) \Big|_{(t, \bar{x}(t), \bar{y}(t))} \mathbb{E}^A (x(t) - \bar{x}(t) | S(t) = s) \\ & \approx \left(\frac{(\nabla S) d_x (\nabla S)^\top}{(\nabla S) c_x (\nabla S)^\top} \right. \\ & \quad \left. + 2 \frac{(\nabla S) (c_x (\nabla^2 S) c_x) (\nabla S)^\top}{(\nabla S) c_x (\nabla S)^\top} \right) \Big|_{(t, \bar{x}(t), \bar{y}(t))} (s - S(0)). \end{aligned}$$

Thus, we obtain (all terms on the right-hand side evaluated at $(t, \bar{x}(t), \bar{y}(t))$)

$$\begin{aligned} \varphi(t, s)^2 & \approx (\nabla S) c_x (\nabla S)^\top \\ & + \left(\frac{(\nabla S) d_x (\nabla S)^\top}{(\nabla S) c_x (\nabla S)^\top} + 2 \frac{(\nabla S) (c_x (\nabla^2 S) c_x) (\nabla S)^\top}{(\nabla S) c_x (\nabla S)^\top} \right) (s - S(0)), \end{aligned}$$

so that, linearizing around $s \approx S(0)$,

$$\begin{aligned} \varphi(t, s) & \approx \left((\nabla S) c_x (\nabla S)^\top \right)^{1/2} + \frac{1}{2 \left((\nabla S) c_x (\nabla S)^\top \right)^{1/2}} \\ & \times \left(\frac{(\nabla S) d_x (\nabla S)^\top}{(\nabla S) c_x (\nabla S)^\top} + 2 \frac{(\nabla S) (c_x (\nabla^2 S) c_x) (\nabla S)^\top}{(\nabla S) c_x (\nabla S)^\top} \right) (s - S(0)). \end{aligned}$$

Setting

$$\lambda_S(t) = \frac{\varphi(t, S(0))}{S(0)}, \quad b_S(t) = S(0) \frac{\frac{\partial \varphi}{\partial s}(t, S(0))}{\varphi(t, S(0))},$$

the main statement of the theorem follows. Finally, from the definition of $\sigma_x(t, x, y)$ (see (13.77)),

$$\frac{\partial \sigma_x^\top}{\partial x_l} = H(t) H^f(t)^{-1} \text{diag} \left(\lambda_1^f(t) b_1^f(t) \frac{\partial f_1(t)}{\partial x_l}, \dots, \lambda_d^f(t) b_d^f(t) \frac{\partial f_d(t)}{\partial x_l} \right) D^f(t)^\top$$

and the expression (13.83) follows from

$$\frac{\partial f_i(t)}{\partial x_l} = h_l(t + \delta_i) / h_l(t).$$

□

Remark 13.3.6. Using averaging techniques, the time-dependent Heston dynamics could be easily translated into time-independent ones. The derivation and the result essentially mimic those for the one-dimensional case, see Proposition 13.2.2.

13.3.4 Volatility Calibration

As explained in Section 13.3.2, the model (13.72), (13.74), (13.77) has enough degrees of freedom to calibrate to the smiles of d swaption strips if the vector function $h(t)$, as well as the time-dependent correlation matrix of benchmark rates $X^f(t)$, are specified exogenously. For multi-factor quasi-Gaussian models, the strips are usually taken to be constant-tenor strips with swap tenors matching benchmark rate tenors. With $d = 4$ or 5 factors being a typical choice of dimensionality, a calibration to 4 or 5 swaption strips essentially recovers the whole universe of swaption volatilities, so there is little need to choose calibration targets in a product-specific way.

As with the one-factor stochastic volatility qG model, we favor splitting the calibration into two main steps. First, we calibrate the volatility of volatility curve $\eta(t)$ to the market-implied curvatures of the smiles or, better yet, to average curvatures of volatility smiles across swap tenors, as we only have the flexibility of making $\eta(t)$ time-specific, not tenor-specific. After that, the main calibration is performed, matching the overall levels and slopes of volatility smiles of d swaption strips. We omit the details as they follow closely the algorithm of Section 13.2.3, with the only difference being that on each time step, the calibration problem involves d swaptions and not 1. Since all formulas are closed-form, the calibration is essentially instantaneous.

13.3.5 Mean Reversions, Correlations, and Numerical Schemes

In the multi-factor context, the time-dependent “loadings” vector $h(t)$ essentially defines the interpolation rule, i.e. how the volatilities and correlations of non-benchmark rates are obtained from those of the benchmark rates. We advocate choosing d fixed values of mean reversions and using them for all cases — note that they should all be different, since the inverse of matrix $H_f(t)$ is required to exist. For example, a reasonable choice is to span the interval $[0, 1]$ with mean reversions while always including the point 0, i.e. set

$$\boldsymbol{\varkappa}(t) = \text{diag}((0.015, 0.15, 0.30, 1.20)^\top)$$

for a 4-dimensional model, corresponding to benchmark tenors

$$\{\delta_1, \dots, \delta_4\} = \{6m, 2y, 10y, 30y\}.$$

However, in principle at least, the mean reversions can be calibrated as well, giving us additional d strips to calibrate to. Formulas for mean reversion calibration could be derived in the same way as for the one-dimensional case, see Section 13.1.8.

Additionally, the correlation matrix between benchmark rates, $X^f(t)$, could in principle be used in calibration, particularly when valuing products with strong correlation sensitivity. In this case, to capture market-implied correlation information in the model, one sometimes chooses to best-fit

market-observed prices of CMS spread options by tweaking benchmark rate correlations. We discuss this in more details in the context of LM models in Section 14.5.9. For now we just note that spread option values exhibit a *correlation smile*, i.e. the dependence of implied correlation on the strike of the spread option (see Section 17.4.2), so the choice of the strike to calibrate to should be carefully considered.

Finally, some brief words on numerical implementation. For $d > 1$, PDE methods quickly become impractical — even for the simple case of $d = 2$, there are 3 auxiliary (y) variables to take care of, pushing the dimension of the PDE to 5 which is prohibitively expensive in virtually all applications. However, by using tricks such as “freezing” or projecting some of the auxiliary variables, as in Section 13.1.9.4, a PDE scheme for $d = 2$ or $d = 3$ could possibly be made viable.

With Monte Carlo methods, the usual considerations apply, and no special tricks beyond those for a one-dimensional stochastic volatility model are required.

13.A Appendix: Density Approximation

We prove Theorem 13.1.14 in a number of steps. Denoting for brevity

$$\sigma = \sigma^0(0),$$

and using the notations of Section 13.1.9.1, we can write down the approximate risk-neutral model dynamics as

$$dx(t) = (\bar{y}(t) - \kappa x(t)) dt + \sigma v(x(t)) dW(t), \quad (13.84)$$

$$d\bar{y}(t) = (\sigma^2 - 2\kappa\bar{y}(t)) dt. \quad (13.85)$$

13.A.1 Simplified Forward Measure Dynamics

As a first step we simplify the dynamics of the state process.

Proposition 13.A.1. *For small time T , the distribution of $x(T)$ in the T -forward measure can be approximated by the distribution of $\tilde{x}(T)$, with the dynamics of $\tilde{x}(t)$ given by*

$$d\tilde{x}(t) = -\kappa\tilde{x}(t) dt + \sigma v(\tilde{x}(t)) dW^T(t), \quad (13.86)$$

where $W^T(t)$ is a Brownian motion in the T -forward measure, and $v(x)$ is the same as in (13.84) (and defined by (13.48)).

Remark 13.A.2. Note that the statement is only about approximating the marginal distribution of $x(T)$ with $\tilde{x}(T)$, not the dynamics of $x(\cdot)$ with $\tilde{x}(\cdot)$ in the T -forward measure.

Proof. The process

$$dW^T(t) = dW(t) + \sigma v(x(t)) G(t, T) dt$$

is a driftless Brownian motion in the T -forward measure, hence

$$dx(t) = \left(\bar{y}(t) - \sigma^2 v(x(t))^2 G(t, T) - \kappa x(t) \right) dt + \sigma v(x(t)) dW^T(t).$$

Also,

$$\bar{y}(t) = \sigma^2 \int_0^t e^{-2\kappa(t-s)} ds = \sigma^2 G_2(t),$$

using notation introduced in (13.3). Thus

$$\begin{aligned} x(T) &= \sigma^2 \int_0^T e^{-\kappa(T-t)} \left(G_2(t) - v(x(t))^2 G(t, T) \right) dt \\ &\quad + \sigma \int_0^T e^{-\kappa(T-t)} v(x(t)) dW^T(t). \end{aligned}$$

On the other hand, the instantaneous forward rate $f(t, T)$ is a martingale in the T -forward measure,

$$\mathbb{E}^T(f(T, T)) = f(0, T),$$

which implies that

$$\mathbb{E}^T(x(T)) = 0.$$

Therefore

$$\mathbb{E}^T \left(\int_0^T e^{-\kappa(T-t)} \left(G_2(t) - v(x(t))^2 G(t, T) \right) dt \right) = 0,$$

where the equality is only approximate as we replaced $y(t)$ with $\bar{y}(t)$, but is as accurate as the approximation of $y(t)$ with $\bar{y}(t)$. We replace the equality with a somewhat stronger condition

$$\int_0^T e^{-\kappa(T-t)} \left(G_2(t) - v(x(t))^2 G(t, T) \right) dt = 0,$$

leading to

$$x(T) = \sigma \int_0^T e^{-\kappa(T-t)} v(x(t)) dW^T(t),$$

which is equivalent to (13.86). \square

13.A.2 Effective Volatility

From now on we consider the distribution of $x(T)$ to be given by (13.86), and we drop the tilde to simplify the notations. The (undiscounted) value of a call option on $x(t)$ with strike k is denoted by

$$c(t, k) = \mathbb{E}^T (x(t) - k)^+. \quad (13.87)$$

By the Bachelier formula (Remark 7.2.9), this function is known explicitly for $v(x) \equiv 1$ (since (13.86) is a Gaussian SDE then), and we denote it c_0 :

$$c_0(t, k, \sigma) = (x_0 - k) \Phi\left(\frac{x_0 - k}{\sigma \sqrt{G_2(t)}}\right) + \sigma \sqrt{G_2(t)} \phi\left(\frac{x_0 - k}{\sigma \sqrt{G_2(t)}}\right),$$

where $\phi(z)$ is the standard Gaussian PDF, $\Phi(z)$ is the standard Gaussian CDF, $x_0 = x(0) = 0$, and $G_2(t)$ is defined in (13.49). Using this expression as the base case, we look for the approximate value of $c(t, k)$, for $x(t)$ governed by (13.86), of the form (compare to the methods of Section 7.5)

$$c(t, k) = c_0(t, k, \sigma \varpi(k)). \quad (13.88)$$

Here $\varpi(k)$ has the meaning of the *effective term volatility*. For notational convenience, we also define

$$\zeta(t, k) = \frac{x_0 - k}{\sigma \varpi(k) \sqrt{G_2(t)}}. \quad (13.89)$$

Then

$$c(t, k) = (x_0 - k) \Phi(\zeta(t, k)) + \sigma \varpi(k) \sqrt{G_2(t)} \phi(\zeta(t, k)). \quad (13.90)$$

In the next few sections, we obtain an expression for $\zeta(t, k)$ in the small- σ limit. To do so, we, firstly, derive a PDE for $c(t, k)$ defined by (13.87). Then, we substitute the expression (13.90) into the PDE to derive an equation on $\zeta(t, k)$. We drop the terms of order $O(\sigma^2)$ and smaller, and solve the simplified equation for $\zeta(t, k)$. Finally, we obtain a CDF and a PDF of $x(t)$ by differentiating $c(t, k)$ in strike.

13.A.3 The Forward Equation for Call Options

In this section we derive a PDE for $c(t, k)$ in the variables t (time to expiry) and k (strike), just as we did in Proposition 7.4.2 for vanilla local volatility models. Let $\psi(t, k)$ be the density of $x(t)$ and $\Psi(t, k)$ its CDF. We use subscripts to denote partial derivatives, and primes to denote derivatives of functions of a single variable.

Proposition 13.A.3. *The function $c(t, k)$ defined by (13.87) for $x(t)$ following (13.86) satisfies*

$$c_t(t, k) = -\varkappa [c(t, k) - kc_k(t, k)] + \frac{\sigma^2}{2} v(k)^2 c_{kk}(t, k) \quad (13.91)$$

with the initial condition

$$c(0, k) = \delta(k),$$

where $\delta(k)$ is the Dirac delta function at 0.

Proof. Follows that of Proposition 7.4.2 closely, with the use of the following identity

$$\begin{aligned} \mathbb{E}(x(t)1_{\{x(t)>k\}}) &= \mathbb{E}((x(t) - k)1_{\{x(t)>k\}}) + k\mathbb{E}(1_{\{x(t)>k\}}) \\ &= c(t, k) - kc_k(t, k). \end{aligned}$$

□

13.A.4 Asymptotic Expansion

Lemma 13.A.4. *The following holds for $c(t, k)$ as defined by (13.90),*

$$\frac{c_t}{\phi(\zeta)} = \sigma \varpi(k) \left(\sqrt{G_2(t)} \right)', \quad (13.92)$$

$$\frac{c - kc_k}{\phi(\zeta)} = (\varpi(k) - k\varpi'(k)) \sigma \sqrt{G_2(t)}, \quad (13.93)$$

$$\frac{c_{kk}}{\phi(\zeta)} = -\zeta_k + \sigma \sqrt{G_2(t)} (\varpi''(k) - \varpi'(k)\zeta_k \zeta), \quad (13.94)$$

where ϖ, ζ are defined by (13.88), (13.89).

Proof. Follows by straightforward differentiation of $c(t, k)$ defined by (13.90) with the help of the identities

$$\phi'(z) = -z\phi(z), \quad \phi''(z) = (z^2 - 1)\phi(z). \quad (13.95)$$

□

Proposition 13.A.5. *If the function $\varpi(k)$ is such that*

$$\pi(k) \triangleq \frac{k}{\varpi(k)}$$

satisfies the ODE

$$v(k)^2 (\pi'(k))^2 + 2\varkappa G_2(t) \pi(k) \left(\frac{k}{\pi(k)} \right)' - 1 = 0, \quad k \in \mathbb{R},$$

with the boundary condition

$$\pi(0) = 0,$$

then

$$c(t, k) = c_0(t, k, \sigma\varpi(k)) + O(\sigma^2),$$

i.e. $c_0(t, k, \sigma\varpi(k))$ is an approximation to $c(t, k)$ (given in (13.90)) to the first order in σ , $\sigma \rightarrow 0$.

Proof. Substituting (13.92)–(13.94) into the PDE (13.91) and keeping only the terms of order $O(\sigma^2)$ we obtain,

$$\begin{aligned} \sigma\varpi(k) \frac{d}{dt} \sqrt{G_2(t)} &= -\varkappa(\varpi(k) - k\varpi'(k)) \sigma \sqrt{G_2(t)} \\ &\quad + \frac{\sigma^2}{2} v(k)^2 \left(-\zeta'(k) - \sigma \sqrt{G_2(t)} \varpi'(k) \zeta'(k) \zeta(k) \right). \end{aligned}$$

Dividing by σ and using the fact that

$$\sigma\varpi(k) \sqrt{G_2(t)} \zeta(k) = -k, \quad (13.96)$$

we obtain

$$\begin{aligned} \varpi(k) \frac{d}{dt} \sqrt{G_2(t)} &= -\varkappa(\varpi(k) - k\varpi'(k)) \sqrt{G_2(t)} \\ &\quad + \frac{\sigma}{2} v(k)^2 \left(\left(\frac{k\varpi'(k)}{\varpi(k)} - 1 \right) \zeta'(k) \right). \quad (13.97) \end{aligned}$$

By definition of $\pi(k)$, with the help of (13.96), and using the definition of $G_2(t)$,

$$\begin{aligned} \zeta'(k) &= -\frac{1}{\sigma\sqrt{G_2(t)}} \pi'(k), \quad \frac{k\varpi'(k)}{\varpi(k)} - 1 = -\pi'(k)\varpi(k), \\ \frac{d}{dt} \sqrt{G_2(t)} &= \frac{dG_2(t)/dt}{2\sqrt{G_2(t)}}, \quad \frac{d}{dt} G_2(t) = e^{-2\varkappa t} = 1 - 2\varkappa G_2(t), \end{aligned}$$

which, substituted into (13.97), gives us, after some simplifications,

$$2\varkappa G_2(t) \frac{k\varpi'(k)}{\varpi(k)} + v(k)^2 (\pi'(k))^2 = 1.$$

In addition,

$$\frac{k\varpi'(k)}{\varpi(k)} = \pi(k) \left(\frac{k}{\pi(k)} \right)'$$

so, finally,

$$v(k)^2 (\pi'(k))^2 + 2\varkappa G_2(t) \pi(k) \left(\frac{k}{\pi(k)} \right)' - 1 = 0.$$

To obtain the boundary conditions on $\pi(k)$, we recall that $\varpi(k) = k/\pi(k)$ and, as $\varpi(k)$ has to be bounded at $k = 0$, we must have $\pi(k) = 0$. \square

13.A.5 Proof of Theorem 13.1.14

The statement of the theorem follows by using Proposition 13.A.3 to simplify the model dynamics to (13.86), and then differentiating $c_0(T, x, \sigma\varpi(x))$ with $\varpi(x)$ from Proposition 13.A.5 once with respect to x to obtain the approximate CDF of $x(T)$,

$$\Psi(T, x) = \frac{\partial c_0(T, x, \sigma\varpi(x))}{\partial x} + 1.$$

We omit tedious but straightforward details.

The Libor Market Model I

Many of the models considered so far describe the evolution of the yield curve in terms of a small set of Markov state variables. While proper calibration procedures allow for successful application of such models to the pricing and hedging of a surprising variety of securities, many exotic derivatives require richer dynamics than what is possible with low-dimensional Markov models. For instance, exotic derivatives may be strongly sensitive to the joint evolution of multiple points of the yield curve, necessitating the usage of several driving Brownian motions. Also, most exotic derivatives may not be related in any obvious way to vanilla European options, making it hard to confidently identify a small, representative set of vanilla securities to which a low-dimensional Markovian model can feasibly be calibrated. What is required in such situations is a model sufficiently rich to capture the full correlation structure across the entire yield curve, and to allow for volatility calibration to a large enough set of European options that the volatility characteristics of most exotic securities can be considered “spanned” by the calibration. Candidates for such a model include the multi-factor short rate models in Chapter 12 and the multi-factor quasi-Gaussian models in Section 13.3. In this chapter, we shall cover an alternative approach to the construction of multi-factor interest rate models, the so-called *Libor market (LM)* model framework. Originally developed in Brace et al. [1997], Jamshidian [1997], and Miltersen et al. [1997], the LM model class enjoys significant popularity with practitioners and is in many ways easier to grasp than, say, the multi-factor quasi-Gaussian models in Chapter 13.

This chapter develops the basic LM model and provides a series of extensions to the original log-normal framework in Brace et al. [1997] and Miltersen et al. [1997] in order to better capture observed volatility smiles. To facilitate calibration of the model, efficient techniques for the pricing of European securities are developed. We provide a detailed discussion of the modeling of forward rate correlations which, along with the pricing formulas for caps and swaptions, serves as the basis for most of the calibration

strategies that we proceed to examine. Many of these strategies are generic in nature and apply to multi-factor models other than the LM class, including the models discussed in Chapters 12 and 13. We wrap up the chapter with a careful discussion of schemes for Monte Carlo simulation of LM models. A number of advanced topics in LM modeling is postponed to Chapter 15.

14.1 Introduction and Setup

14.1.1 Motivation and Historical Notes

Chapter 4 introduced the HJM framework which, in its most general form, involves multiple driving Brownian motions and an infinite set of state variables (namely the set of instantaneous forward rates). As argued earlier, the HJM framework contains any arbitrage-free interest rate model adapted to a finite set of Brownian motions. Working directly with instantaneous forward rates is, however, not particularly attractive in applications, for a variety of reasons. First, instantaneous forward rates are never quoted in the market, nor do they figure directly in the payoff definition of any traded derivative contract. As discussed in Chapter 5, realistic securities (swaps, caps, futures, etc.) instead involve simply compounded (Libor) rates, effectively representing *integrals* of instantaneous forward rates. The disconnect between market observables and model primitives often makes development of market-consistent pricing expression for simple derivatives cumbersome. Second, an infinite set of instantaneous forward rates can generally¹ not be represented exactly on a computer, but will require discretization into a finite set. Third, prescribing the form of the volatility function of instantaneous forward rates is subject to a number of technical complications, requiring sub-linear growth to prevent explosions in the forward rate dynamics, which precludes the formulation of a log-normal forward rate model (see Sandmann and Sondermann [1997] and the discussion in Sections 4.5.3 and 11.1.3).

As discovered in Brace et al. [1997], Jamshidian [1997], and Miltersen et al. [1997], the three complications above can all be addressed simultaneously by simply formulating the model in terms of a non-overlapping set of simply compounded Libor rates. Not only do we then conveniently work with a finite set of directly observable rates that can be represented on a computer but, as we shall show, an explosion-free log-normal forward rate model also becomes possible. Despite the change to simply compounded rates, we should emphasize that the Libor market model will still be a special case of an HJM model, albeit one where we only indirectly specify the volatility function of the instantaneous forward rates.

¹As we have seen in earlier chapters, for special choices of the forward rate volatility we can sometimes identify a finite-dimensional Markovian representation of the forward curve that eliminates the need to store the entire curve. This is not possible in general, however.

14.1.2 Tenor Structure

The starting point for our development of the LM model is a fixed tenor structure

$$0 = T_0 < T_1 < \dots < T_N. \quad (14.1)$$

The intervals $\tau_n = T_{n+1} - T_n$, $n = 0, \dots, N - 1$, would typically be set to be either 3 or 6 months, corresponding to the accrual period associated with observable Libor rates. Rather than keeping track of an entire yield curve, at any point in time t we are (for now; but see Section 15.1) focused only on a finite set of zero-coupon bonds $P(t, T_n)$ for the set of n 's for which $T_N \geq T_n > t$; notice that this set shrinks as t moves forward, becoming empty when $t > T_N$. To formalize this “roll-off” of zero-coupon bonds in the tenor structure as time progresses, it is often useful to work with an *index function* $q(t)$, defined by the relation

$$T_{q(t)-1} \leq t < T_{q(t)}. \quad (14.2)$$

We think of $q(t)$ as representing the tenor structure index of the shortest-dated discount bond still alive.

On the fixed tenor structure, we proceed to define Libor forward rates according to the relation (see (4.2))

$$L(t, T_n, T_{n+1}) = L_n(t) = \frac{1}{\tau_n} \left(\frac{P(t, T_n)}{P(t, T_{n+1})} - 1 \right), \quad N - 1 \geq n \geq q(t).$$

We note that when considering a given forward Libor rate $L_n(t)$, we always assume $n \geq q(t)$ unless stated otherwise. For any $T_n > t$,

$$P(t, T_n) = P(t, T_{q(t)}) \prod_{i=q(t)}^{n-1} (1 + L_i(t)\tau_i)^{-1}. \quad (14.3)$$

Notice that knowledge of $L_n(t)$ for all $n \geq q(t)$ is generally *not* sufficient to reconstruct discount bond prices on the entire (remaining) tenor structure; the front “stub” discount bond price $P(t, T_{q(t)})$ must also be known.

14.2 LM Dynamics and Measures

14.2.1 Setting

In the Libor market model, the set of Libor forward rates $L_{q(t)}(t), L_{q(t)+1}(t), \dots, L_{N-1}(t)$ constitutes the set of state variables for which we wish to specify dynamics. As a first step, we pick a probability measure P and assume that those dynamics originate from an m -dimensional Brownian motion $W(t)$, in the sense that all Libor rates are measurable with

respect to the filtration generated by $W(t)$. Further assuming that the Libor rates are square integrable, it follows from the martingale representation theorem that, for all $n \geq q(t)$,

$$dL_n(t) = \sigma_n(t)^\top (\mu_n(t) dt + dW(t)), \quad (14.4)$$

where μ_n and σ_n are m -dimensional processes, respectively, both adapted to the filtration generated by $W(t)$. From the diffusion invariance principle (see Section 1.5) we know that $\sigma_n(t)$ is measure invariant, whereas $\mu_n(t)$ is not.

As it turns out, for a given choice of $\sigma_n(t)$ in the specification (14.4), it is quite straightforward to work out explicitly the form of $\mu_n(t)$ in various martingale measures of practical interest. We turn to this shortly, but let us first stress that (14.4) allows us to use a *different* volatility function σ_n for each of the forward rates $L_n(t)$, $n = q(t), \dots, N - 1$, in the tenor structure. This obviously gives us tremendous flexibility in specifying the volatility structure of the forward curve evolution, but in practice will require us to impose quite a bit of additional structure on the model to ensure realism and to avoid an excess of parameters. We shall return to this topic later in this chapter.

14.2.2 Probability Measures

As shown in Lemma 4.2.3, $L_n(t)$ must be a martingale in the T_{n+1} -forward measure $Q^{T_{n+1}}$, such that, from (14.4),

$$dL_n(t) = \sigma_n(t)^\top dW^{n+1}(t), \quad (14.5)$$

where $W^{n+1}(t) \triangleq W^{T_{n+1}}(t)$ is an m -dimensional Brownian motion in $Q^{T_{n+1}}$. It is to be emphasized that only *one* specific Libor forward rate — namely L_n — is a martingale in the T_{n+1} -forward measure. To establish dynamics in other probability measures, the following proposition is useful.

Proposition 14.2.1. *Let $L_n(t)$ satisfy (14.5). In measure Q^{T_n} the process for $L_n(t)$ is*

$$dL_n(t) = \sigma_n(t)^\top \left(\frac{\tau_n \sigma_n(t)}{1 + \tau_n L_n(t)} dt + dW^n(t) \right),$$

where $W^n(t)$ is an m -dimensional Brownian motion in measure Q^{T_n} .

Proof. From Theorem 1.4.2 we know that the density $\varsigma(t)$ relating the measures $Q^{T_{n+1}}$ and Q^{T_n} is given by

$$\begin{aligned} \varsigma(t) &= E_t^{T_{n+1}} \left(\frac{dQ^{T_n}}{dQ^{T_{n+1}}} \right) \\ &= \frac{P(t, T_n)/P(0, T_n)}{P(t, T_{n+1})/P(0, T_{n+1})} = (1 + \tau_n L_n(t)) \frac{P(0, T_{n+1})}{P(0, T_n)}. \end{aligned}$$

Clearly, then,

$$d\varsigma(t) = \frac{P(0, T_{n+1})}{P(0, T_n)} \tau_n dL_n(t) = \frac{P(0, T_{n+1})}{P(0, T_n)} \tau_n \sigma_n(t)^\top dW^{n+1}(t),$$

or

$$d\varsigma(t)/\varsigma(t) = \frac{\tau_n \sigma_n(t)^\top dW^{n+1}(t)}{1 + \tau_n L_n(t)}.$$

From the Girsanov theorem (Theorem 1.5.1), it follows that the process

$$dW^n(t) = dW^{n+1}(t) - \frac{\tau_n \sigma_n(t)}{1 + \tau_n L_n(t)} dt \quad (14.6)$$

is a Brownian motion in Q^{T_n} . The proposition then follows directly from (14.5). \square

To gain some further intuition for the important result in Proposition 14.2.1, let us derive it in less formal fashion. For this, consider the forward discount bond $P(t, T_n, T_{n+1}) = P(t, T_{n+1})/P(t, T_n) = (1 + \tau_n L_n(t))^{-1}$. An application of Ito's lemma to $P(t, T_n, T_{n+1})$, with the help of (14.5), shows that

$$\begin{aligned} dP(t, T_n, T_{n+1}) &= \tau_n^2 (1 + \tau_n L_n(t))^{-3} \sigma_n(t)^\top \sigma_n(t) dt \\ &\quad - \tau_n (1 + \tau_n L_n(t))^{-2} \sigma_n(t)^\top dW^{n+1}(t) \\ &= \tau_n (1 + \tau_n L_n(t))^{-2} \sigma_n(t)^\top \left\{ \tau_n (1 + \tau_n L_n(t))^{-1} \sigma_n(t) dt - dW^{n+1}(t) \right\}. \end{aligned}$$

As $P(t, T_n, T_{n+1})$ must be a martingale in the Q^{T_n} -measure, it follows that

$$-dW^n(t) = \tau_n (1 + \tau_n L_n(t))^{-1} \sigma_n(t) dt - dW^{n+1}(t)$$

is a Brownian motion in Q^{T_n} , consistent with the result in Proposition 14.2.1.

While Proposition 14.2.1 only relates “neighboring” measures $Q^{T_{n+1}}$ and Q^{T_n} , it is straightforward to use the proposition iteratively to find the drift of L_n in any of the probability measures discussed in Section 4.2. Let us give some examples.

Lemma 14.2.2. *Let $L_n(t)$ satisfy (14.5). Under the terminal measure Q^{T_N} the process for $L_n(t)$ is*

$$dL_n(t) = \sigma_n(t)^\top \left(- \sum_{j=n+1}^{N-1} \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)} dt + dW^N(t) \right),$$

where $W^N(t)$ is an m -dimensional Brownian motion in measure Q^{T_N} .

Proof. From (14.6) we know that

$$\begin{aligned} dW^N(t) &= dW^{N-1}(t) + \frac{\tau_{N-1}\sigma_{N-1}(t)}{1+\tau_{N-1}L_{N-1}(t)} dt \\ &= dW^{N-2}(t) + \frac{\tau_{N-2}\sigma_{N-2}(t)}{1+\tau_{N-2}L_{N-2}(t)} dt + \frac{\tau_{N-1}\sigma_{N-1}(t)}{1+\tau_{N-1}L_{N-1}(t)} dt. \end{aligned}$$

Continuing this iteration down to W^{n+1} , we get

$$dW^N(t) = dW^{n+1}(t) + \sum_{j=n+1}^{N-1} \frac{\tau_j\sigma_j(t)}{1+\tau_jL_j(t)} dt.$$

The lemma now follows from (14.5). \square

Lemma 14.2.3. *Let $L_n(t)$ satisfy (14.5). Under the spot measure Q^B (see Section 4.2.3) the process for $L_n(t)$ is*

$$dL_n(t) = \sigma_n(t)^\top \left(\sum_{j=q(t)}^n \frac{\tau_j\sigma_j(t)}{1+\tau_jL_j(t)} dt + dW^B(t) \right), \quad (14.7)$$

where $W^B(t)$ is an m -dimensional Brownian motion in measure Q^B .

Proof. Recall from Section 4.2.3 that the spot measure is characterized by a rolling or “jumping” numeraire

$$B(t) = P(t, T_{q(t)}) \prod_{n=0}^{q(t)-1} (1 + \tau_n L_n(T_n)). \quad (14.8)$$

At any time t , the random part of the numeraire is the discount bond $P(t, T_{q(t)})$, so effectively we need to establish dynamics in the measure $Q^{T_{q(t)}}$. Applying the iteration idea shown in the proof of Lemma 14.2.2, we get

$$dW^{n+1}(t) = dW^{q(t)}(t) + \sum_{j=q(t)}^n \frac{\tau_j\sigma_j(t)}{1+\tau_jL_j(t)} dt,$$

as stated. \square

The spot and terminal measures are, by far, the most commonly used probability measures in practice. Occasionally, however, it may be beneficial to use one of the hybrid measures discussed earlier in this book, for instance if one wishes to enforce that a particular Libor rate $L_n(t)$ be a martingale. As shown in Section 4.2.4, we could pick as a numeraire the asset price process

$$\tilde{P}_{n+1}(t) = \begin{cases} P(t, T_{n+1}), & t \leq T_{n+1}, \\ B(t)/B(T_{n+1}), & t > T_{n+1}, \end{cases} \quad (14.9)$$

where $B(t)$ is the spot numeraire (14.8). Using the same technique as in the proofs of Lemmas 14.2.2 and 14.2.3, it is easily seen that when $i \geq n$, then

$$dL_i(t) = \begin{cases} \sigma_i(t)^\top \left(\sum_{j=n+1}^i \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)} dt + d\tilde{W}^{n+1}(t) \right), & t \leq T_{n+1}, \\ \sigma_i(t)^\top \left(\sum_{j=q(t)}^i \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)} dt + d\tilde{W}^{n+1}(t) \right), & t > T_{n+1}, \end{cases}$$

where $\tilde{W}^{n+1}(t)$ is a Brownian motion in the measure induced by the numeraire $\tilde{P}_{n+1}(t)$. Note in particular that $L_n(t)$ is a martingale as desired, and that we have defined a numeraire which — unlike $P(t, T_{n+1})$ — will be alive at any time t .

We should note that an equally valid definition of a hybrid measure will replace (14.9) with the asset process

$$\bar{P}_{n+1}(t) = \begin{cases} B(t), & t \leq T_{n+1}, \\ B(T_{n+1}) P(t, T_N) / P(T_{n+1}, T_N), & t > T_{n+1}. \end{cases} \quad (14.10)$$

This type of numeraire process is often useful in discretization of the LM model for simulation purposes; see Section 14.6.1.2 for details.

14.2.3 Link to HJM Analysis

As discussed earlier, the LM model is a special case of the general HJM class of diffusive interest rate models. To explore this relationship a bit further, we recall that HJM models generally have risk-neutral dynamics of the form

$$df(t, T) = \sigma_f(t, T)^\top \int_t^T \sigma_f(t, u) du dt + \sigma_f(t, T)^\top dW(t),$$

where $f(t, T)$ is the time t instantaneous forward rate to time T and $\sigma_f(t, T)$ is the instantaneous forward rate volatility function. From the results in Chapter 4, it follows that dynamics for the forward bond $P(t, T_n, T_{n+1})$ are of the form

$$\frac{dP(t, T_n, T_{n+1})}{P(t, T_n, T_{n+1})} = O(dt) - (\sigma_P(t, T_{n+1})^\top - \sigma_P(t, T_n)^\top) dW(t),$$

where $O(dt)$ is a drift term and

$$\sigma_P(t, T) = \int_t^T \sigma_f(t, u) du.$$

By definition $L_n(t) = \tau_n^{-1}(P(t, T_n, T_{n+1})^{-1} - 1)$, so that

$$dL_n(t) = O(dt) + \tau_n^{-1}(1 + \tau_n L_n(t)) \int_{T_n}^{T_{n+1}} \sigma_f(t, u)^\top du dW(t).$$

By the diffusion invariance principle, it follows from (14.5) that the LM model volatility $\sigma_n(t)$ is related to the HJM instantaneous forward volatility function $\sigma_f(t, T)$ by

$$\sigma_n(t) = \tau_n^{-1}(1 + \tau_n L_n(t)) \int_{T_n}^{T_{n+1}} \sigma_f(t, u) du. \quad (14.11)$$

Note that, as expected, $\sigma_n(t) \rightarrow \sigma_f(t, T_n)$ as $\tau_n \rightarrow 0$.

It should be obvious from (14.11) that a complete specification of $\sigma_f(t, T)$ uniquely determines the LM volatility $\sigma_n(t)$ for all t and all n . On the other hand, specification of $\sigma_n(t)$ for all t and all n does *not* allow us to imply a unique HJM forward volatility function $\sigma_f(t, T)$ — all we are specifying is essentially a strip of contiguous integrals of this function in the T -direction. This is hardly surprising, inasmuch as the LM model only concerns itself with a finite set of discretely compounded forward rates and cannot be expected to uniquely characterize the behaviors of instantaneous forward rates and their volatilities. Along the same lines, we note that the LM model does not uniquely specify the behavior of the short rate $r(t) = f(t, t)$; as a consequence, the rolling money market account $\beta(t)$ and the risk-neutral measure are not natural constructions in the LM model². Section 15.3 discusses these issues in more detail.

14.2.4 Separable Deterministic Volatility Function

So far, our discussion of the LM model has been generic, with little structure imposed on the $N - 1$ volatility functions $\sigma_n(t)$, $n = 1, 2, \dots, N - 1$. To build a workable model, however, we need to be more specific about our choice of $\sigma_n(t)$. A common prescription of $\sigma_n(t)$ takes the form

$$\sigma_n(t) = \lambda_n(t)\varphi(L_n(t)), \quad (14.12)$$

where $\lambda_n(t)$ is a bounded vector-valued deterministic function and $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is a time-homogeneous local volatility function. This specification is conceptually very similar to the local volatility models in Chapter 7, although here $\sigma_n(\cdot)$ is vector-valued and the model involves joint dynamics of multiple state variables (the $N - 1$ Libor forward rates).

At this point, the reader may reasonably ask whether the choice (14.12) in fact leads to a system of SDEs for the various Libor forward rates that is “reasonable”, in the sense of existence and uniqueness of solutions. While we here shall not pay much attention to such technical regularity issues, it should be obvious that not all functions φ can be allowed. One relevant result is given below.

²In fact, as discussed in Jamshidian [1997], one does not need to assume that a short rate process exists when constructing an LM model.

Proposition 14.2.4. Assume that (14.12) holds with $\varphi(0) = 0$ and that $L_n(0) \geq 0$ for all n . Also assume that φ is locally Lipschitz continuous and satisfies the growth condition

$$\varphi(x)^2 \leq C(1 + x^2), \quad x > 0,$$

where C is some positive constant. Then non-explosive, pathwise unique solutions of the no-arbitrage SDEs for $L_n(t)$, $q(t) \leq n \leq N - 1$, exist under all measures Q^{T_i} , $q(t) \leq i \leq N$. If $L_n(0) > 0$, then $L_n(t)$ stays positive at all t .

Proof. (Sketch) Due to the recursive relationship between measures, it suffices to consider the system of SDEs (14.7) under the spot measure Q^B :

$$dL_n(t) = \varphi(L_n(t)) \lambda_n(t)^\top (\mu_n(t) dt + dW^B(t)), \quad (14.13)$$

$$\mu_n(t) = \sum_{j=q(t)}^n \frac{\tau_j \varphi(L_j(t)) \lambda_j(t)}{1 + \tau_j L_j(t)}. \quad (14.14)$$

Under our assumptions, it is easy to see that each term in the sum for μ_n is locally Lipschitz continuous and bounded. The growth condition on φ in turn ensures that the product $\varphi(L_n(t)) \lambda_n(t)^\top \mu_n(t)$ is also locally Lipschitz continuous and, due to the boundedness of μ_n , satisfies a linear growth condition. Existence and uniqueness now follow from Theorem 1.6.1. The result that 0 is a non-accessible boundary for the forward rates if started above 0 follows from standard speed-scale boundary classification results; see Andersen and Andreasen [2000b] for the details. \square

Some standard parameterizations of φ are shown in Table 14.1. Of those, only the log-normal specification and the LCEV specification directly satisfy the criteria in Proposition 14.2.4. The CEV specification violates Lipschitz continuity at $x = 0$, and as a result uniqueness of the SDE fails. As shown in Andersen and Andreasen [2000b], we restore uniqueness by specifying that forward rates are *absorbed* at the origin (see also Section 7.2.3). As for the displaced log-normal specification $\varphi(x) = ax + b$, we here violate the assumption that $\varphi(0) = 0$, and as a result we cannot always guarantee that forward rates stay positive. Also, to prevent explosion of the forward rate drifts, we need to impose additional restrictions to prevent terms of the form $1 + \tau_n L_n(t)$ (in the denominator) from becoming zero. As displaced log-normal models are of considerable practical importance, we list the relevant restrictions in Lemma 14.2.5 below.

Lemma 14.2.5. Consider a local volatility Libor market model with local volatility function $\varphi(x) = bx + a$, where $b > 0$ and $a \neq 0$. Assume that $bL_n(0) + a > 0$ and $a/b < \tau_n^{-1}$ for all $n = 1, 2, \dots, N - 1$. Then non-explosive, pathwise unique solutions of the no-arbitrage SDEs for $L_n(t)$, $q(t) \leq n \leq N - 1$, exist under all measures Q^{T_i} , $q(t) \leq i \leq N$. All $L_n(t)$ are bounded from below by $-a/b$.

Name	$\varphi(x)$
Log-normal	x
CEV	$x^p, \quad 0 < p < 1$
LCEV	$x \min(\varepsilon^{p-1}, x^{p-1}), \quad 0 < p < 1, \varepsilon > 0$
Displaced log-normal	$bx + a, \quad b > 0, a \neq 0$

Table 14.1. Common DVF Specifications

Proof. Define $H_n(t) = bL_n(t) + a$. By Ito's lemma, we have

$$dH_n(t) = b dL_n(t) = bH_n(t)\lambda_n(t)^\top (\mu_n(t) dt + dW^B(t)),$$

$$\mu_n(t) = \sum_{j=q(t)}^n \frac{\tau_j H_j(t)\lambda_j(t)}{1 + \tau_j (H_j(t) - a)/b}.$$

From the assumptions of the lemma, we have $H_n(0) > 0$ for all n , allowing us to apply the result of Proposition 14.2.4 to $H_n(t)$, provided that we can guarantee that $\mu_j(t)$ is bounded for all positive H_j , $j = q(t), \dots, n$. This follows from $1 - \tau_j a/b > 0$ or $a/b < \tau_j^{-1}$. \square

We emphasize that the requirement $a/b < \tau_n^{-1}$ implies that only in the limit of $\tau_j \rightarrow 0$ — where the discrete forward Libor rates become instantaneous forward rates — will a pure Gaussian LM model specification ($b = 0$) be meaningful; such a model was outlined in Section 4.5.1. On the flip-side, according to Proposition 14.2.4, a finite-sized value of τ_j ensures that a well-behaved log-normal forward rate model exists, something that we saw earlier (Section 11.1.3) was *not* the case for models based on instantaneous forward rates. The existence of log-normal forward rate dynamics in the LM setting was, in fact, a major driving force behind the development and popularization of the LM framework, and all early examples of LM models (see Brace et al. [1997], Jamshidian [1997], and Miltersen et al. [1997]) were exclusively log-normal.

We recall from earlier chapters that it is often convenient to specify displaced log-normal models as $\varphi(L_n(t)) = (1 - b)L_n(0) + bL_n(t)$, in which case the constant a in Lemma 14.2.5 is different from one Libor rate to the next. In this case, we must require

$$(1 - b)/b < (L_n(0)\tau_n)^{-1}, \quad n = 1, \dots, N - 1.$$

As $L_n(0)\tau_n$ is typically in the magnitude of a few percent, the regularity requirement on b in (14.2.4) is not particularly restrictive.

14.2.5 Stochastic Volatility

As discussed earlier in this book, to ensure that the evolution of the volatility smile is reasonably stationary, it is best if the skew function φ in (14.14)

is (close to) monotonic in its argument. Typically we are interested in specifications where $\varphi(x)/x$ is downward-sloping, to establish the standard behavior of interest rate implied volatilities tending to increase as interest rates decline. In reality, however, markets often exhibit non-monotonic volatility smiles or “smirks” with high-struck options trading at implied volatilities above the at-the-money levels. An increasingly popular mechanism to capture such behavior in LM models is through the introduction of stochastic volatility. We have already encountered stochastic volatility models in Chapters 8, 9 and, in the context of term structure models, in Sections 13.2 and 13.3; we now discuss how to extend the notion of stochastic volatility models to the simultaneous modeling of multiple Libor forward rates.

As our starting point, we take the process (14.14), preferably equipped with a φ that generates either a flat or monotonically downward-sloping volatility skew, but allow the term on the Brownian motion to be scaled by a stochastic process. Specifically, we introduce a mean-reverting scalar process $z(t)$, with dynamics of the form

$$dz(t) = \theta(z_0 - z(t)) dt + \eta\psi(z(t)) dZ(t), \quad z(0) = z_0, \quad (14.15)$$

where θ , z_0 , and η are positive constants, Z is a Brownian motion under the spot measure Q^B , and $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a well-behaved function. We impose that (14.15) will not generate negative values of $z(t)$, which requires $\psi(0) = 0$. We will interpret the process in (14.15) as the (scaled) variance process for our forward rate diffusions, in the sense that the square root of $z(t)$ will be used as a stochastic, multiplicative scaling of the diffusion term in (14.14). That is, our forward rate processes in Q^B are, for all $n \geq q(t)$,

$$dL_n(t) = \sqrt{z(t)}\varphi(L_n(t))\lambda_n(t)^\top \left(\sqrt{z(t)}\mu_n(t) dt + dW^B(t) \right), \quad (14.16)$$

$$\mu_n(t) = \sum_{j=q(t)}^n \frac{\tau_j \varphi(L_j(t)) \lambda_j(t)}{1 + \tau_j L_j(t)},$$

where $z(t)$ satisfies (14.15). This construction naturally follows the specification of vanilla stochastic volatility models in Chapter 8, and the specification of stochastic volatility quasi-Gaussian models in Chapter 13. As we discussed previously, it is often natural to scale the process for $z(t)$ such that $z(0) = z_0 = 1$.

Let us make two important comments about (14.16). First, we emphasize that a single common factor $\sqrt{z(t)}$ simultaneously scales all forward rate volatilities; movements in volatilities are therefore perfectly correlated across the various forward rates. In effect, our model corresponds only to the first principal component of the movements of the instantaneous forward rate volatilities. This is a common assumption that provides good balance between realism and parsimony, and we concentrate mostly on this case — although we do relax it later in the book, in Chapter 15. Second, we note

that the clean form of the z -process (14.15) in the measure Q^B generally does not carry over to other probability measures, as we would expect from Proposition 8.3.9. To state the relevant result, let $\langle Z(t), W(t) \rangle$ denote the vector of quadratic covariations between $Z(t)$ and the m components of $W(t)$ (recall the definition of covariation in Remark 1.1.7). We then have

Lemma 14.2.6. *Let dynamics for $z(t)$ in the measure Q^B be as in (14.15). Then the SDE for $z(t)$ in measure $Q^{T_{n+1}}$, $n \geq q(t) - 1$, is*

$$\begin{aligned} dz(t) &= \theta(z_0 - z(t)) dt + \eta \psi(z(t)) \\ &\quad \times \left(-\sqrt{z(t)} \mu_n(t)^\top \langle dZ(t), dW^B(t) \rangle + dZ^{n+1}(t) \right), \end{aligned} \quad (14.17)$$

where $\mu_n(t)$ is given in (14.16) and $Z^{n+1}(t)$ is a Brownian motion in measure $Q^{T_{n+1}}$.

Proof. From earlier results, we have

$$dW^{n+1}(t) = \sqrt{z(t)} \mu_n(t) dt + dW^B(t).$$

Let us introduce the m -dimensional vector

$$a(t) = \langle dZ(t), dW^B(t) \rangle / dt,$$

so that we can write

$$dZ(t) = a(t)^\top dW^B(t) + \sqrt{1 - \|a(t)\|^2} d\tilde{W}(t),$$

where $\tilde{W}(t)$ is a scalar Brownian motion independent of $W^B(t)$. In the measure $Q^{T_{n+1}}$, we then have

$$\begin{aligned} dZ(t) &= a(t)^\top \left(dW^{n+1}(t) - \sqrt{z(t)} \mu_n(t) dt \right) + \sqrt{1 - \|a(t)\|^2} d\tilde{W}(t) \\ &= dZ^{n+1}(t) - a(t)^\top \sqrt{z(t)} \mu_n(t) dt, \end{aligned}$$

and the result follows. \square

The process (14.17) is awkward to deal with, due to presence of the drift term $\mu_n(t)^\top \langle dZ(t), dW^B(t) \rangle$ which will, in general, depend on the state of the Libor forward rates at time t . For tractability, on the other hand, we would like for the z -process to only depend on $z(t)$ itself. To achieve this, and to generally simplify measure shifts in the model, we make the following assumption³ about (14.15)–(14.16):

Assumption 14.2.7. *The Brownian motion $Z(t)$ of the variance process $z(t)$ is independent of the vector-valued Brownian motion $W^B(t)$.*

³We briefly return to the general case in Section 15.6.

We have already encountered the same assumption in the context of stochastic volatility quasi-Gaussian models, see Section 13.2.1, where we also discussed the implications of such a restriction.

The diffusion coefficient of the variance process, the function ψ , is traditionally chosen to be of power form, $\psi(x) = x^\alpha, \alpha > 0$. While it probably makes sense to keep the function monotonic, the power specification is likely a nod to tradition rather than anything else. Nevertheless, some particular choices lead to analytically tractable specifications, as we saw in Chapter 8; for that reason, $\alpha = 1/2$ (the Heston model) is popular.

Remark 14.2.8. Going forward we shall often use the stochastic volatility model in this section as a benchmark for theoretical and numerical work. As the stochastic volatility model reduces to the local volatility model in Section 14.2.4 when $z(t)$ is constant, all results for the stochastic volatility model will carry over to the DVF setting.

14.2.6 Time-Dependence in Model Parameters

In the models we outlined in Sections 14.2.4 and 14.2.5, the main role of the vector-valued function of time $\lambda_n(t)$ was to establish a term structure “spine” of at-the-money option volatilities. To build volatility smiles around this spine, we further introduced a universal skew-function φ , possibly combined with a stochastic volatility scale $z(t)$ with time-independent process parameters. In practice, this typically gives us a handful of free parameters with which we can attempt to match the market-observed volatility smiles for various cap and swaption tenors. As it turns out, a surprisingly good fit to market skew data can, in fact, often be achieved with the models of Sections 14.2.4 and 14.2.5. For a truly precise fit to volatility skews across all maturities and swaption tenors it may, however, be necessary to allow for time-dependence in both the process parameters for $z(t)$ and, more importantly, the skew function φ . The resulting model is conceptually similar to the model in Section 14.2.5, but involves a number of technical intricacies that draw heavily on the material presented in Chapter 9. To avoid cluttering this first chapter on LM models with technical detail, we postpone the treatment of time-inhomogeneous φ and z -process parameters to Chapter 15.

14.3 Correlation

In one-factor models for interest rates — such as the ones presented in Chapters 10 and 11 — all points on the forward curve always move in the same direction. While this type of forward curve move indeed is the most commonly observed type of shift to the curve, “rotational steepenings” and

the formation of “humps” may also take place, as may other more complex types of curve changes. The empirical presence of such non-trivial curve movements is an indication of the fact that various points on the forward curve do not move co-monotonically with each other, i.e. they are imperfectly correlated. A key characteristic of the LM model is the consistent use of vector-valued Brownian motion drivers, of dimension m , which gives us control over the instantaneous correlation between various points on the forward curve.

Proposition 14.3.1. *The correlation between forward rate increments $dL_k(t)$ and $dL_j(t)$ in the SV model (14.16) is*

$$\text{Corr}(dL_k(t), dL_j(t)) = \frac{\lambda_k(t)^\top \lambda_j(t)}{\|\lambda_k(t)\| \|\lambda_j(t)\|}.$$

Proof. Using the covariance notation of Remark 1.1.7, we have, for any j and k ,

$$d\langle L_k(t), L_j(t) \rangle = z(t) \varphi(L_k(t)) \varphi(L_j(t)) \lambda_k(t)^\top \lambda_j(t) dt.$$

Using this in the definition of the correlation,

$$\text{Corr}(dL_k(t), dL_j(t)) = \frac{\langle dL_k(t), dL_j(t) \rangle}{\sqrt{\langle dL_k(t) \rangle \langle dL_j(t) \rangle}},$$

which gives the result of the proposition. \square

A trivial corollary of Proposition 14.3.1 is the fact that $\text{Corr}(dL_k(t), dL_j(t)) = 1$ always when $m = 1$, i.e. when we only have one Brownian motion. As we add more Brownian motions, our ability to capture increasingly complicated correlation structures progressively improves (in a sense that we shall examine further shortly), but at a cost of increasing the model complexity and, ultimately, computational effort. To make rational decisions about the choice of model dimension m , let us turn to the empirical data.

14.3.1 Empirical Principal Components Analysis

For some fixed value of τ (e.g. 0.25 or 0.5), let us define “sliding” forward rates⁴ $l(t, x)$ with tenor x as

$$l(t, x) = L(t, t + x, t + x + \tau).$$

⁴The use of sliding forward rates, i.e. forward rates with a fixed time **to** maturity rather than a fixed time **of** maturity, is often known as the *Musiela parameterization*.

For a given set of tenors x_1, \dots, x_{N_x} and a given set of calendar times t_0, t_1, \dots, t_{N_t} , we can use market observations⁵ to set up the $N_x \times N_t$ observation matrix O with elements

$$O_{i,j} = \frac{l(t_j, x_i) - l(t_{j-1}, x_i)}{\sqrt{t_j - t_{j-1}}}, \quad i = 1, \dots, N_x, \quad j = 1, \dots, N_t.$$

Notice the normalization with $\sqrt{t_j - t_{j-1}}$ which annualizes the variance of the observed forward rate increments. Also note that we use absolute increments in forward rates here. This is arbitrary — we could have used, say, relative increases as well, if we felt that rates were more log-normal than Gaussian. For small sampling periods, the precise choice is of little importance.

Assuming time-homogeneity and ignoring small drift terms, the data collected above will imply a sample $N_x \times N_x$ variance-covariance matrix equal to

$$C = \frac{OO^\top}{N_t}. \quad (14.18)$$

For our LM model to conform to empirical data, we need to use a sufficiently high number m of Brownian motions to closely replicate this variance-covariance matrix. A formal analysis of what value of m will suffice can proceed with the tools of principal components analysis (PCA), as established in Section 3.1.3.

14.3.1.1 Example: USD Forward Rates

To give a concrete example of a PCA run, we set $N_x = 9$ and use tenors of $\{x_1, \dots, x_9\} = \{0.5, 1, 2, 3, 5, 7, 10, 15, 20\}$ years. We fix $\tau = 0.5$ (i.e., all forward rates are 6 months discrete rates) and use 4 years of weekly data from the USD market, spanning January 2003 to January 2007, for a total of $N_t = 203$ curve observations. The eigenvalues of the matrix C in (14.18) are listed in Table 14.2, along with the percentage of variance that is explained by using only the first m principal components.

m	1	2	3	4	5	6	7	8	9
Eigenvalue	7.0	0.94	0.29	0.064	0.053	0.029	0.016	0.0091	0.0070
% Variance	83.3	94.5	97.9	98.7	99.3	99.6	99.8	99.9	100

Table 14.2. PCA for USD Rates. All eigenvalues are scaled up by 10^4 .

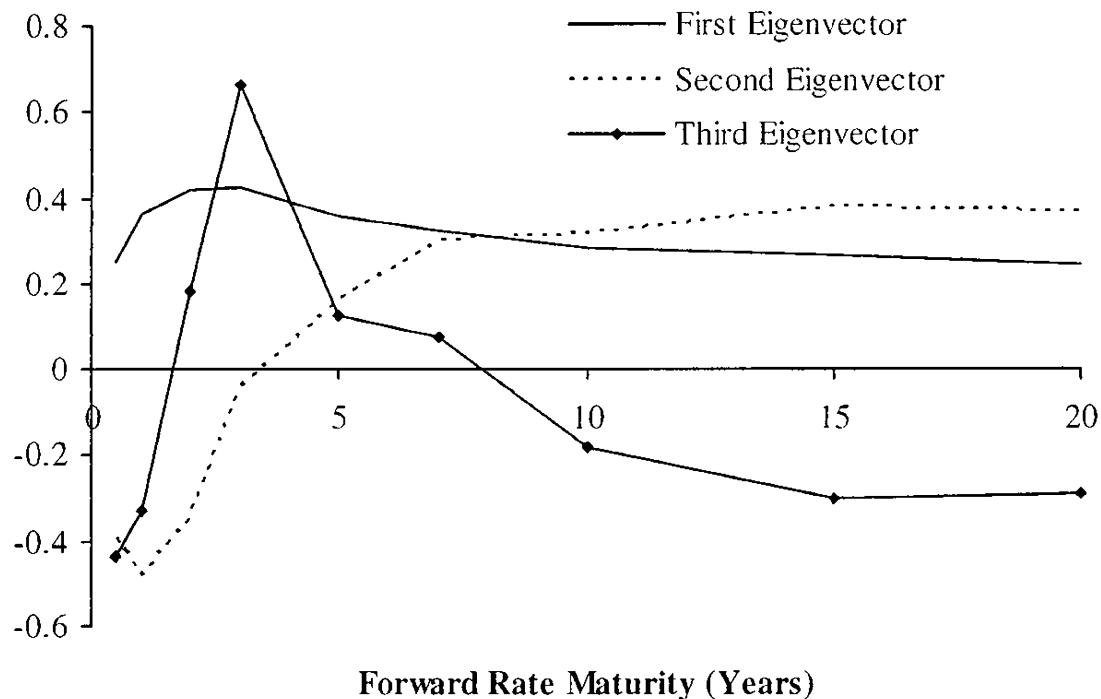
As we see from the table, the first principal component explains about 83% of the observed variance, and the first three principal components together

⁵For each date in the time grid t_j we construct the forward curve from market observable swaps, futures, and deposits, using the techniques from Chapter 6.

explain nearly 98%. This pattern carries over to most major currencies, and in many applications we would consequently expect that using $m = 3$ or $m = 4$ Brownian motions in a LM model would adequately capture the empirical covariation of the points on the forward curve. An exception to this rule-of-thumb occurs when a particular derivative security depends strongly on the correlation between forward rates with tenors that are close to each other; in this case, as we shall see in Section 14.3.4, a high number of principal components is required to provide for sufficient decoupling of nearby forward rates.

The eigenvectors corresponding to the largest three eigenvectors in Table 14.2 are shown in the Figure 14.1; the figure gives us a reasonable idea about what the (suitably scaled) first three elements of the $\lambda_k(t)$ vectors should look like as functions of $T_k - t$. Loosely speaking, the first principal component can be interpreted as a near-parallel shift of the forward curve, whereas the second and third principal components correspond to forward curve twists and bends, respectively.

Fig. 14.1. Eigenvectors



Notes: Eigenvectors for the largest three eigenvalues in Table 14.2.

14.3.2 Correlation Estimation and Smoothing

Empirical estimates for forward rate correlations can proceed along the lines of Section 14.3.1. Specifically, if we introduce the diagonal matrix

$$c \triangleq \begin{pmatrix} \sqrt{C_{1,1}} & 0 & \ddots & 0 \\ 0 & \sqrt{C_{2,2}} & \ddots & \ddots \\ \ddots & \ddots & \ddots & 0 \\ 0 & \ddots & 0 & \sqrt{C_{N_x,N_x}} \end{pmatrix},$$

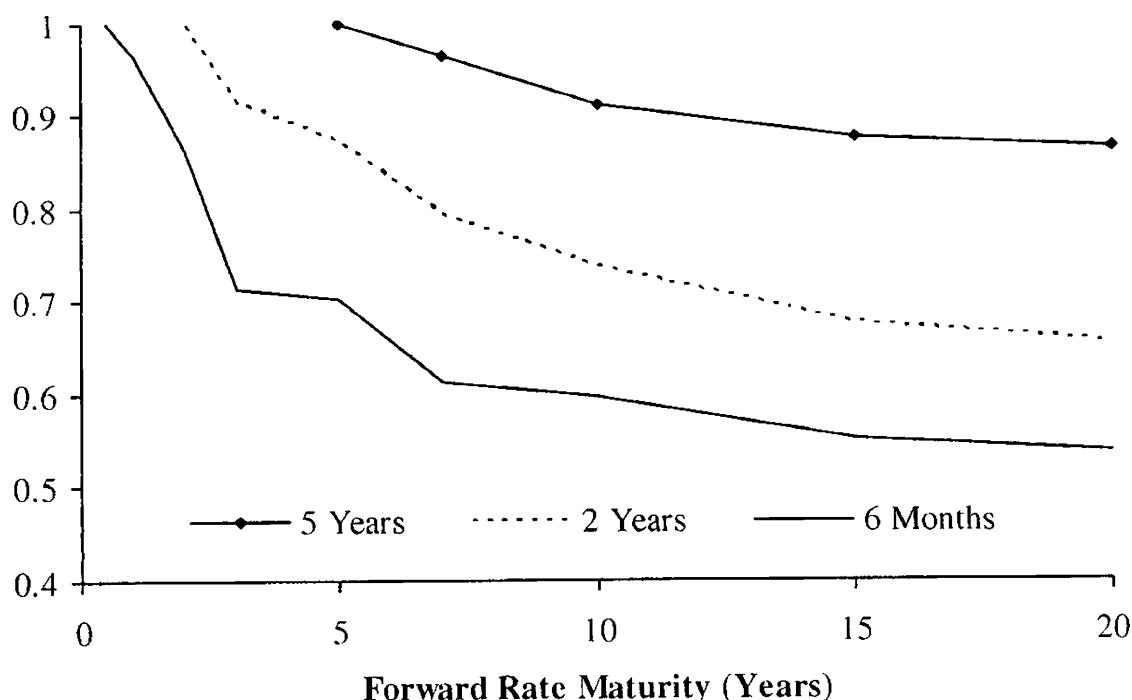
then the empirical $N_x \times N_x$ forward rate correlation matrix R becomes

$$R = c^{-1}Cc^{-1}.$$

Element $R_{i,j}$ of R provides a sample estimate of the instantaneous correlation between increments in $l(t, x_i)$ and $l(t, x_j)$, under the assumption that this correlation is time-homogeneous.

The matrix R is often relatively noisy, partially as a reflection of the fact that correlations are well-known to be quite variable over time, and partially as a reflection of the fact that the empirical correlation estimator has rather poor sample properties with large confidence bounds (see Johnson et al. [1995] for details). Nevertheless, several stylistic facts can be gleaned from the data, as demonstrated in Figure 14.2 where we have graphed a few slices of the correlation matrix for the USD data in Section 14.3.1.1.

Fig. 14.2. Forward Rate Correlations



Notes: For each of three fixed forward rate maturities (6 months, 2 years, and 5 years), the figure shows the correlation between the fixed forward rate and forward rates with other maturities (as indicated on the x -axis of the graph).

To make a few qualitative observations about Figure 14.2, we notice that correlations between forward rates $l(\cdot, x_k)$ and $l(\cdot, x_j)$ generally decline in $|x_k - x_j|$; this decline appears near-exponential for x_k and x_j close to each other, but with a near-flat asymptote for large $|x_k - x_j|$. It appears that the rate of the correlation decay and the level of the asymptote depend not only on $|x_k - x_j|$, but also on $\min(x_k, x_j)$. Specifically, the decay rate decreases with $\min(x_k, x_j)$, and the asymptote level increases with $\min(x_k, x_j)$.

In practice, unaltered empirical correlation matrices are typically too noisy for comfort, and might contain non-intuitive entries (e.g., correlation between a 10 year forward and a 2 year forward might come out higher than between a 10 year forward and a 5 year forward). As such, it is common practice in multi-factor yield curve modeling to work with simple parametric forms; this not only smoothes the correlation matrix, but also conveniently reduces the effective parameter dimension of the correlation matrix object, from $N_x(N_x - 1)/2$ distinct matrix elements to the number of parameters in the parametric form.

Several candidate parametric forms for the correlation have been proposed in the literature, see Schoenmakers and Coffey [2000], Jong et al. [2001], and Rebonato [2002], among many others. Rather than list all of these, we instead focus on a few reasonable forms that we have designed to encompass most or all of the empirical facts listed above. Our first parametric form is as follows:

$$\text{Corr}(dL_k(t), dL_j(t)) = q_1(T_k - t, T_j - t),$$

where

$$\begin{aligned} q_1(x, y) &= \rho_\infty + (1 - \rho_\infty) \exp(-a(\min(x, y)) |y - x|), \\ a(z) &= a_\infty + (a_0 - a_\infty)e^{-\kappa z}, \end{aligned} \tag{14.19}$$

subject to $0 \leq \rho_\infty \leq 1$, $a_0, a_\infty, \kappa \geq 0$. Fundamentally, $q_1(x, y)$ exhibits correlation decay at a rate of a as $|y - x|$ is increased, with the decay rate a itself being an exponential function of $\min(x, y)$. We would always expect to have $a_0 \geq a_\infty$, in which case

$$\frac{\partial q_1(x, y)}{\partial x} = (1 - \rho_\infty)e^{-a(x)(y-x)} [a(x) + (y - x)\kappa(a_0 - a_\infty)e^{-\kappa x}], \quad x < y,$$

is non-negative, as one would expect.

Variations on (14.19) are abundant in the literature — the case $a_0 = a_\infty$ is particularly popular — and q_1 generally has sufficient degrees of freedom to provide a reasonable fit to empirical data. One immediate issue, however, is a lack of control of the asymptotic correlation level at $|x - y| \rightarrow \infty$ which, as we argued above, is typically not independent of x and y . As the empirical data suggests that ρ_∞ tends to increase with $\min(x, y)$, we could introduce yet another decaying function

$$\rho_\infty(z) = b_\infty + (b_0 - b_\infty)e^{-\alpha z}, \tag{14.20}$$



and extend q_1 to the “triple-decaying” form

$$q_2(x, y) = \rho_\infty(\min(x, y)) + (1 - \rho_\infty(\min(x, y))) \exp(-a(\min(x, y))) e^{-|x-y|}$$

with $a(z)$ given in (14.19), and where $0 \leq b_0, b_\infty \leq 1, \alpha \geq 0$. Empirical data suggests that normally $b_0 \leq b_\infty$, in which case we have

$$\begin{aligned} \frac{\partial q_2(x, y)}{\partial x} &= -\alpha(b_0 - b_\infty)e^{-\alpha x} \left(1 - e^{-a(x)(y-x)}\right) \\ &\quad + (1 - \rho_\infty(x))e^{-a(x)(y-x)} [a(x) + (y-x)\kappa(a_0 - a_\infty)e^{-\kappa x}], \quad x < y \end{aligned}$$

which remains non-negative if $b_0 \leq b_\infty$ and $a_0 \geq a_\infty$.

In a typical application, the four parameters of q_1 and the six parameters of q_2 are found by least-squares optimization against an empirical correlation matrix. Any standard optimization algorithm, such as the Levenberg-Marquardt algorithm in Press et al. [1992], can be used for this purpose. Some parameters are here subject to simple box-style constraints (e.g. $\rho_\infty \in [0, 1]$), which poses no particular problems for most commercial optimizers. In any case, we can always use functional mappings to rewrite our optimization problem in terms of variables with unbounded domains. For instance, for the form q_1 , we can set

$$\rho_\infty = \frac{1}{2} + \frac{\arctan(u)}{\pi}, \quad u \in (-\infty, \infty),$$

and optimize on the variable u instead of ρ_∞ . Note that we sometimes may wish to optimize correlation parameters against more market-driven targets than empirical correlation matrices, an idea that we shall investigate further in Section 14.5.9.

14.3.2.1 Example: Fit to USD Data

Let R be the 9×9 empirical correlation matrix generated from the data in Section 14.3.1.1, and let $R_2(\xi)$, $\xi \triangleq (a_0, a_\infty, \kappa, b_0, b_\infty, \alpha)^\top$, be the 9×9 correlation matrix generated from the form q_2 , when using the 9 specific forward tenors in 14.3.1.1. To determine the optimal parameter vector ξ^* , we minimize an unweighted Frobenius (least-squares) matrix norm, subject to a non-negativity constraint

$$\xi^* = \underset{\xi}{\operatorname{argmin}} \left(\operatorname{tr} \left((R - R_2(\xi))(R - R_2(\xi))^\top \right) \right), \text{ subject to } \xi \geq 0.$$

The resulting fit is summarized in Table 14.3; Figure 14.3 in Section 14.3.4.1 contains a 3D plot of the correlation matrix $R_2(\xi^*)$.

The value of the Frobenius norm at ξ^* is 0.070, which translates into an average absolute correlation error (excluding diagonal elements) of around

a_0	a_∞	κ	b_0	b_∞	α
0.312	0.157	0.264	0.490	0.946	0.325

Table 14.3. Best-Fit Parameters for q_2 in USD Market

2%. If we use the four parameter form q_1 instead of q_2 in the optimization exercise, the Frobenius norm at the optimum increases to 0.164. As we would expect from Figure 14.2, allowing correlation asymptotes to increase in tenors thus adds significant explanatory power to the parametric form.

14.3.3 Negative Eigenvalues

While some functional forms are designed to always return valid correlation matrices (the function in Schoenmakers and Coffey [2000] being one such example), many popular forms — including our q_1 and q_2 above — can, when stressed, generate matrices R that fail to be positive definite. While this rarely happens in real applications, it is not inconceivable that on occasion one or more eigenvalues of R may turn out to be negative, requiring us to somehow “repair” the matrix. A similar problem can also arise due to rounding errors when working with large empirical correlation matrices.

Formally, when faced with an R matrix that is not positive definite, we would ideally like to replace it with a modified matrix R^* which i) is a valid correlation matrix; and ii) is as close as possible to R , in the sense of some matrix norm. The problem of locating R^* then involves computing the norm

$$\{\|R - X\| : X \text{ is a correlation matrix}\}$$

and setting R^* equal to the matrix X that minimizes this distance. If $\|\cdot\|$ is a weighted Frobenius norm, good numerical algorithms for the computation of R^* have recently emerged, see Higham [2002] for a review and a clean approach.

If the negative eigenvalues are small in absolute magnitude (which is often the case in practice), it is often reasonable to abandon a full-blown optimization algorithm in favor of a more heuristic approach where we simply raise all offending negative eigenvalues to some positive cut-off value. To present one obvious algorithm, let us start by writing

$$R = E\Lambda E^\top,$$

where Λ is a diagonal matrix of eigenvalues, and E is a matrix with the eigenvectors of R in its columns. Let Λ^* be the diagonal matrix with all-positive entries

$$\Lambda_{i,i}^* = \max(\epsilon, \Lambda_{i,i}), \quad i = 1, \dots, N_x,$$

for some small cut-off value $\epsilon > 0$. Then set

$$C^* = E\Lambda^*E^\top,$$

which we interpret as a *covariance* matrix, i.e. of the form

$$C^* = c^*R^*c^*,$$

where c^* is a diagonal matrix with elements $c_{i,i}^* = \sqrt{C_{i,i}^*}$ and R^* is the valid, positive definite correlation matrix we seek. R^* is then computed as

$$R^* = (c^*)^{-1}C^*(c^*)^{-1}. \quad 14.22)$$

We emphasize that R^* as defined in (14.22) will have 1's in its diagonal, whereas C^* will not. Both C^* and R^* are, by construction, positive definite.

14.3.4 Correlation PCA

We now turn to a problem that arises in certain important applications, such as the calibration procedure we shall discuss in Section 14.5. Consider a p -dimensional Gaussian variable Y , where all elements of Y have zero mean and unit variance. Let Y have a positive definite correlation matrix R , given by

$$R = E(YY^\top).$$

Consider now writing, as an approximation,

$$Y \approx DX, \quad 14.23$$

where X is an m -dimensional vector of independent standard Gaussian variables, $m < p$, and D is a $(p \times m)$ -dimensional matrix. We wish to strictly enforce that DX remains a vector of variables with zero means and unit variances, thereby ensuring that the matrix DD^\top has the interpretation of a valid correlation matrix. In particular, we require that DD^\top has ones on its diagonal.

Let $v(D)$ be the p -dimensional vector of the diagonal elements of DD^\top , i.e. $v_i = (DD^\top)_{i,i}$, $i = 1, \dots, p$. Working as before with an unweighted⁶ Frobenius norm, we set

$$h(D; R) = \text{tr} \left((R - DD^\top) (R - DD^\top)^\top \right), \quad 14.24$$

and define the optimal choice of D , denoted D^* , as

$$D^* = \underset{D}{\operatorname{argmin}} h(D; R), \quad \text{subject to } v(D) = \mathbf{1}, \quad 14.25$$

where $\mathbf{1}$ is a p -dimensional vector of 1's.

⁶The introduction of user-specified weights into this norm is a straightforward extension.

Proposition 14.3.2. Let μ be a p -dimensional vector, and let D_μ be given as the unconstrained optimum

$$D_\mu = \underset{D}{\operatorname{argmin}} h(D; R + \operatorname{diag}(\mu)),$$

with h given in (14.24). Define D^* as in (14.25) and let μ^* be the solution to

$$v(D_\mu) - 1 = 0.$$

Then $D^* = D_{\mu^*}$.

Proof. We only provide a sketch of the proof; for more details, see Zhang and Wu [2003]. First, we introduce the Lagrangian

$$\mathfrak{L}(D, \mu) = h(D; R) - 2\mu^\top (v(D) - 1).$$

(The factor 2 on μ^\top simplifies results.) Standard matrix calculus shows that

$$\frac{dh(D; R)}{dD} = \left\{ \frac{dh(D; R)}{dD_{i,j}} \right\} = -4RD + 4DD^\top D.$$

We can use this result to compute the derivative of the Lagrangian with respect to D , which in turn yields the following condition for an optimum

$$-(R + \operatorname{diag}(\mu))D + DD^\top D = 0, \quad (14.26)$$

where we still must enforce the condition $v(D) = 1$. Equation (14.26) identifies the optimum as minimizing the (unconstrained) optimization norm $h(D; R + \operatorname{diag}(\mu))$. \square

Remark 14.3.3. For any fixed value of μ , D_μ can be computed easily by standard PCA methods provided we interpret $R + \operatorname{diag}(\mu)$ as the target covariance matrix.

With Proposition 14.3.2, determination of D^* is reduced to solving the p -dimensional root-search problem $v(D_\mu) - 1 = 0$ for μ . Many standard methods will suffice; for instance, one can use straightforward secant search methods such as the Broyden algorithm on p. 389 of Press et al. [1992].

As is the case for ordinary PCA approximations of covariance matrices, the “correlation PCA” algorithm outlined so far will return a correlation matrix approximation $D^*(D^*)^\top$ that has reduced rank (from p down to m), a consequence of the PCA steps taken in estimating D_μ .

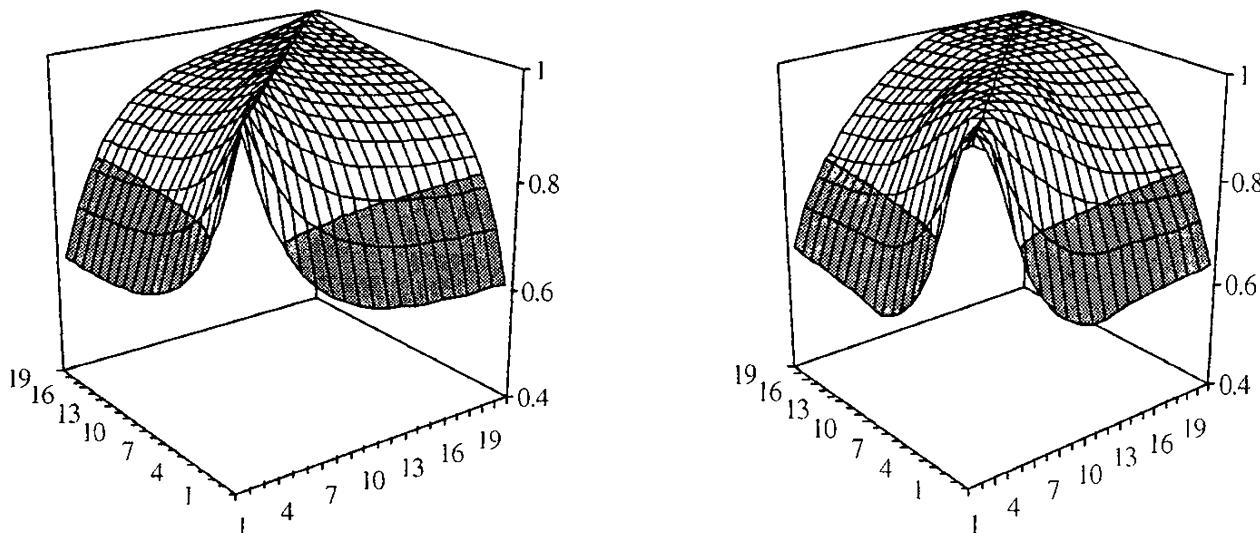
Computation of optimal rank-reduced correlation approximations is a relatively well-understood problem, and we should note the existence of several recent alternatives to the basic algorithm we outlined above. A survey can be found in Pietersz and Groenen [2004] where an algorithm

based on *majorization* is also developed⁷. A recent paper by Li and Qi [2009] not covered in the survey of Pietersz and Groenen [2004] develops an approach based on representing the problem as a non-convex semi-definite programming problem. The authors show how a numerical algorithm that is highly efficient for large scale problems can be constructed. We should also note that certain heuristic (and non-optimal) methods have appeared in the literature, some of which are closely related to the simple algorithm we outlined in Section 14.3.3 for repair of correlation matrices. We briefly outline one such approach below (in Section 14.3.4.2), but first we consider a numerical example.

14.3.4.1 Example: USD Data

We here consider performing a correlation PC analysis on the correlation matrix R generated from our best-fit form q_2 in Section 14.3.2.1. The 3D plots in Figure 14.3 below show the correlation fit we get with a rank-3 correlation matrix.

Fig. 14.3. Forward Rate Correlation Matrix in USD



Notes: The left-hand panel shows the correlation matrix R for form q_2 calibrated to USD data. The right-hand panel shows the best-fitting rank-3 correlation matrix, computed by the algorithm in Proposition 14.3.2. In both graphs, the x - and y -axes represent the Libor forward rate maturities in years.

Looking at Figure 14.3, the effect of rank reduction is, loosely, that the exponential decay of our original matrix R away from the diagonal has been

⁷In our experience, the majorization method in Pietersz and Groenen [2004] is faster than the method in Proposition 14.3.2 but less robust, at least for large and irregular correlation matrices.

replaced with a “sigmoid” shape (to paraphrase Riccardo Rebonato) that is substantially too high close to the matrix diagonal. As the rank of the approximating correlation matrix is increased, the sigmoid shape is — often rather slowly — pulled towards the exponential shape of the full-rank data. Intuitively, we should not be surprised at this result: with the rank m being a low number, we effectively only incorporate smooth, large-scale curve movements (e.g. parallel shifts and twists) into our statistical model, and there is no mechanism to “pull apart” Libor forward rates with maturities close to each other.

Analysis of this difference — rather than the simple PCA considerations of Section 14.3.1 — often forms the basis for deciding how many factors m to use in the model, especially for pricing derivatives with strong correlation dependence. For the reader’s guidance, we find that $m = 5$ to 10 suffices to recover the full-rank correlation shape in most cases.

14.3.4.2 Poor Man’s Correlation PCA

For the case where the $p \times p$ correlation matrix R is well-represented by a rank- m representation of the form (14.23), it may sometimes be sufficiently accurate to compute the loading matrix D by a simpler algorithm based on *standard* PCA applied directly to the correlation matrix. Specifically, suppose that we as a first step compute

$$R_m = E_m \Lambda_m E_m^\top,$$

where Λ_m is an $m \times m$ diagonal matrix of the m largest eigenvalues of R , and E_m is a $p \times m$ matrix of eigenvectors corresponding to these eigenvalues. While the error $R_m - R$ minimizes a least-squares norm, R_m itself is obviously not a valid approximation to the correlation matrix R as no steps were taken to ensure that R_m has a unit diagonal. A simple way to accomplish this borrows the ideas of Section 14.3.3 and writes

$$R \approx r_m^{-1} R_m r_m^{-1}, \quad (14.27)$$

where r_m is a diagonal matrix with elements $(r_m)_{i,i} = \sqrt{(R_m)_{i,i}}$, $i = 1, \dots, p$. We note that this approximation sets the matrix D in (14.23) to

$$D = r_m^{-1} E_m \sqrt{\Lambda_m}.$$

It is clear that the difference between the “poor man’s” PCA result (14.27) and the optimal result in Proposition 14.3.2 will generally be small if R_m is close to having a unit diagonal, as the heuristic step taken in (14.27) will then have little effect. For large, complex correlation matrices, however, the optimal approximation in Proposition 14.3.2 will often be quite different from (14.27) unless m is quite large.

14.4 Pricing of European Options

The previous section laid the foundation for calibrating an LM model to empirical forward curve correlation data, a topic that we shall return to in more detail in Section 14.5. Besides correlation calibration, however, we need to ensure that the forward rate *variances* implied by the LM model are in line with market data. In most applications — and certainly in all those that involve pricing and hedging of traded derivatives — this translates into a requirement that the vectors $\lambda_k(t)$ of the model are such that it will successfully reproduce the prices of liquid plain-vanilla derivatives, i.e. swaptions and caps. A condition for practical uses of the LM model is thus that we can find pricing formulas for vanilla options that are fast enough to be embedded into an iterative calibration algorithm.

14.4.1 Caplets

Deriving formulas for caplets is generally straightforward in the LM model, a consequence of the fact that Libor rates — which figure directly in the payout formulas for caps — are the main primitives of the LM model itself. Indeed, the word “market” in the term “Libor market model” originates from the ease with which the model can accommodate market-pricing of caplets by the Black formula.

As our starting point here, we use the generalized version of the LM model with skews and stochastic volatility; see (14.15) and (14.16). Other, simpler models, are special cases of this framework, and the fundamental caplet pricing methodology will carry over to these cases in a transparent manner. We consider the price of a c -strike caplet $V_{\text{caplet}}(\cdot)$ maturing at time T_n and settling at time T_{n+1} . That is,

$$V_{\text{caplet}}(T_{n+1}) = \tau_n (L_n(T_n) - c)^+.$$

For the purpose of pricing the caplet, the m -dimensional Brownian motion $W^{n+1}(t)$ can here be reduced to one dimension, as shown in the following result.

Proposition 14.4.1. *Assume that the forward rate dynamics in the spot measure are as in (14.15)–(14.16), and that Assumption 14.2.7 holds. Then*

$$V_{\text{caplet}}(0) = P(0, T_{n+1}) \tau_n E^{T_{n+1}} \left((L_n(T_n) - c)^+ \right),$$

where

$$\begin{aligned} dL_n(t) &= \sqrt{z(t)} \varphi(L_n(t)) \|\lambda_n(t)\| dY^{n+1}(t), \\ dz(t) &= \theta(z_0 - z(t)) dt + \eta \psi(z(t)) dZ(t), \end{aligned} \tag{14.28}$$

and $Y^{n+1}(t)$ and $Z(t)$ are independent scalar Brownian motions in measure $Q^{T_{n+1}}$. Specifically, $Y^{n+1}(t)$ is given by

$$Y^{n+1}(t) = \int_0^t \frac{\lambda_n(s)^\top}{\|\lambda_n(s)\|} dW^{n+1}(s).$$

Proof. $Y^{n+1}(t)$ is clearly Gaussian, with mean 0 and variance \sqrt{t} , identifying $Y^{n+1}(t)$ as a Brownian motion such that $\|\lambda_n(t)\|dY^{n+1}(t) = \lambda_n(t)^\top dW^{n+1}(t)$. The remainder of the proposition follows from the martingale property of $L_n(t)$ in $Q^{T_{n+1}}$, combined with the assumed independence of the forward rates and the process for $z(t)$. \square

While rather obvious, Proposition 14.4.1 is useful as it demonstrates that caplet pricing reduces to evaluation of an expectation of $(L_n(T_n) - c)^+$, where the process for $L_n(t)$ is now identical to the types of scalar stochastic volatility diffusions covered in detail in Chapters 8 and 9; the pricing of caplets can therefore be accomplished with the formulas listed in these chapters. In the same way, when dealing with LM models of the simpler local volatility type, we compute caplet prices directly from formulas in Chapter 7.

14.4.2 Swaptions

Whereas pricing of caplets is, by design, convenient in LM models, swaption pricing requires a bit more work and generally will involve some amount of approximation if a quick algorithm is required. In this section, we will outline one such approximation which normally has sufficient accuracy for calibration applications. A more accurate (but also more complicated) approach can be found in Section 15.2.

First, let us recall some notations. Let $V_{\text{swaption}}(t)$ denote the time t value of a payer swaption that matures at time $T_j \geq t$, with the underlying security being a fixed-for-floating swap making payments at times T_{j+1}, \dots, T_k , where $j < k \leq N$. We define an annuity factor for this swap as (see (4.8))

$$A(t) \triangleq A_{j,k-j}(t) = \sum_{n=j}^{k-1} P(t, T_{n+1}) \tau_n, \quad \tau_n = T_{n+1} - T_n. \quad (14.29)$$

Assuming that the swap underlying the swaption pays a fixed coupon of c against Libor, the payout of V_{swaption} at time T_j is (see Section 4.1.3)

$$V_{\text{swaption}}(T_j) = A(T_j) (S(T_j) - c)^+,$$

where we have defined a par forward swap rate (see (4.10))

$$S(t) \triangleq S_{j,k-j}(t) = \frac{P(t, T_j) - P(t, T_k)}{A(t)}.$$

Assume, as in Section 14.4.1, that we are working in the setting of a stochastic volatility LM model, of the type defined in Section 14.2.5; the procedure we shall now outline will carry over to simpler models unchanged.

Proposition 14.4.2. *Assume that the forward rate dynamics in the spot measure are as in (14.15)–(14.16). Let Q^A be the measure induced by using $A(t)$ in (14.29) as a numeraire, and let $W^A(t)$ be an m -dimensional Brownian motion in Q^A . Then, in measure Q^A ,*

$$dS(t) = \sqrt{z(t)}\varphi(S(t)) \sum_{n=j}^{k-1} w_n(t)\lambda_n(t)^\top dW^A(t), \quad (14.30)$$

where the stochastic weights are

$$\begin{aligned} w_n(t) &= \frac{\varphi(L_n(t))}{\varphi(S(t))} \times \frac{\partial S(t)}{\partial L_n(t)} = \frac{\varphi(L_n(t))}{\varphi(S(t))} \times \frac{S(t)\tau_n}{1 + \tau_n L_n(t)} \\ &\quad \times \left[\frac{P(t, T_k)}{P(t, T_j) - P(t, T_k)} + \frac{\sum_{i=n}^{k-1} \tau_i P(t, T_{i+1})}{A(t)} \right]. \end{aligned} \quad (14.31)$$

Proof. It follows from Lemma 4.2.4 that $S(t)$ is a martingale in measure Q^A , hence we know that the drift of the process for $S(t)$ must be zero in this measure. From its definition, $S(t)$ is a function of $L_j(t), L_{j+1}(t), \dots, L_{k-1}(t)$, and an application of Ito's lemma shows that

$$dS(t) = \sum_{n=j}^{k-1} \sqrt{z(t)}\varphi(L_n(t)) \frac{\partial S(t)}{\partial L_n(t)} \lambda_n(t)^\top dW^A(t).$$

Evaluating the partial derivative proves the proposition. \square

It should be immediately obvious that the dynamics of the par rate in (14.30) are too complicated to allow for analytical treatment. The main culprit are the random weights $w_n(t)$ in (14.31) which depend on the entire forward curve in a complex manner. All is not lost, however, as one would intuitively expect that $S(t)$ is well-approximated by a weighted sum of its “component” forward rates $L_j(t), L_{j+1}(t), \dots, L_{k-1}(t)$, with weights varying little over time. In other words, we expect that, for each n , $\partial S(t)/\partial L_n(t)$ is a near-constant quantity.

Consider now the ratio $\varphi(L_n(t))/\varphi(S(t))$ which multiplies $\partial S(t)/\partial L_n(t)$ in (14.31). For forward curves that are reasonably flat and forward curve movements that are predominantly parallel (which is consistent with our earlier discussion in Section 14.3.1.1), it is often reasonable to assume that the ratio is close to constant. This assumption obviously hinges on the precise form of φ , but tends to hold well for many of the functions that we would consider using in practice. To provide some loose motivation for this statement, consider first the extreme case where $\varphi(x) = \text{const}$ (i.e. the

model is Gaussian) in which case the ratio $\varphi(L_n(t))/\varphi(S(t))$ is constant, by definition. Second, let us consider the log-normal case where $\varphi(x) = x$. In this case, a parallel shift h of the forward curve at time t would move the ratio to

$$\frac{L_n(t) + h}{S(t) + h} = \frac{L_n(t)}{S(t)} + h \frac{S(t) - L_n(t)}{S(t)^2} + O(h^2),$$

which is small if the forward curve slope (and thereby $S(t) - L_n(t)$) is small. As the φ 's that we use in practical applications are mostly meant to produce skews that lie somewhere between log-normal and Gaussian ones, assuming that $\varphi(L_n(t))/\varphi(S(t))$ is constant thus appears reasonable.

The discussion above leads to the following approximation, where we “freeze” the weights $w_n(t)$ at their time 0 values.

Proposition 14.4.3. *The time 0 price of the swaption is given by*

$$V_{\text{swaption}}(0) = A(0)E^A((S(T_j) - c)^+). \quad (14.32)$$

Let $w_n(t)$ be as in Proposition 14.4.2 and set

$$\lambda_S(t) = \sum_{n=j}^{k-1} w_n(0) \lambda_n(t),$$

The swap rate dynamics in Proposition 14.4.2 can be approximated as

$$\begin{aligned} dS(t) &\approx \sqrt{z(t)} \varphi(S(t)) \|\lambda_S(t)\| dY^A(t), \\ dz(t) &= \theta(z_0 - z(t)) dt + \eta \psi(z(t)) dZ(t), \end{aligned} \quad (14.33)$$

where $Y^A(t)$ and $Z(t)$ are independent scalar Brownian motions in measure Q^A , and

$$\|\lambda_S(t)\| dY^A(t) = \sum_{n=j}^{k-1} w_n(0) \lambda_n(t)^\top dW^A(t).$$

Proof. Equation (14.32) follows from standard properties of Q^A . The remainder of the proposition is proven the same way as Proposition 14.4.1, after approximating $w_n(t) \approx w_n(0)$. \square

We emphasize that the scalar term $\|\lambda_S(t)\|$ is purely deterministic, whereby the dynamics of $S(t)$ in the annuity measure have precisely the same form as the Libor rate SDE in Proposition 14.4.1. Therefore, computation of the Q^A -expectation in (14.32) can lean directly on the analytical results we established for scalar stochastic volatility processes in Chapter 8 and, for simpler DVF-type LM models, in Chapter 7. We review relevant results and apply them to LM models in Chapter 15; here, to give an example, we merely list a representative result for a displaced log-normal local volatility LM model.

Proposition 14.4.4. *Let each rate $L_n(t)$ follow a displaced log-normal process in its own forward measure,*

$$dL_n(t) = (bL_n(t) + (1 - b)L_n(0))\lambda_n(t)^\top dW^{n+1}(t), \quad n = 1, \dots, N - 1.$$

Then the time 0 price of the swaption is approximated by

$$V_{\text{swaption}}(0) \approx A(0)c_B(0, S(0)/b; T_j, c - S(0) + S(0)/b; b\bar{\lambda}_S),$$

where $c_B(t, S; T, K; \sigma)$ is the Black call option formula with volatility σ , see Remark 7.2.8, and the term swap rate volatility $\bar{\lambda}_S$ is given by

$$\bar{\lambda}_S = \left(\frac{1}{T_j} \int_0^{T_j} \|\lambda_S(t)\|^2 dt \right)^{1/2},$$

with $\lambda_S(t)$ defined in Proposition 14.4.3.

Proof. By Proposition 14.4.3, the approximate dynamics of $S(t)$ are given by

$$dS(t) \approx (bS(t) + (1 - b)S(0)) \|\lambda_S(t)\| dY^A(t).$$

The result then follows from Proposition 7.2.12. \square

While we do not document the performance of the approximation (14.33) in detail here, many tests are available in the literature; see e.g. Andersen and Andreasen [2000b], Glasserman and Merener [2001], and Rebonato [2002]. Suffice to say that the approximation above is virtually always accurate enough for the calibration purposes for which it is designed, particularly if we restrict ourselves to pricing swaptions with strikes close to the forward swap rate. As mentioned earlier, should further precision be desired, one can turn to the more sophisticated swaption pricing approximations that we discuss in Chapter 15. Finally, we should note the existence of models where no approximations are required to price swaptions; these so-called *swap market models* are reviewed in Section 15.4.

14.4.3 Spread Options

When calibrating LM models to market data, the standard approach is to fix the correlation structure in the model to match empirical forward rate correlations. It is, however, tempting to consider whether one alternatively could imply the correlation structure directly from traded market data, thereby avoiding the need for “backward-looking” empirical data altogether. As it turns out, the dependence of swaptions and caps on the correlation structure is, not surprisingly, typically too indirect to allow one to simultaneously back out correlations and volatilities from the prices of these types of instruments alone. To overcome this, one can consider amending the set of calibration instruments with securities that have stronger sensitivity to forward rate

correlations. A good choice would here be to use *yield curve spread options*, a type of security that we encountered earlier in Section 5.13.3. Spread options are natural candidates, not only because their prices are highly sensitive to correlation, but also because they are relatively liquid and not too difficult to value in an LM model setting.

14.4.3.1 Term Correlation

Let $S_1(t) = S_{j_1, k_1 - j_1}(t)$ and $S_2(t) = S_{j_2, k_2 - j_2}(t)$ be two forward swap rates, and assume that we work with a stochastic volatility LM model of type (14.15)–(14.16). Following the result of Proposition 14.4.3, for $i = 1, 2$ we have, to good approximation,

$$dS_i(t) \approx O(dt) + \sqrt{z(t)}\varphi(S_i(t))\lambda_{S_i}(t)^\top dW^B(t),$$

$$\lambda_{S_i}(t) \triangleq \sum_{n=j_i}^{k_i-1} w_{S_i,n}(0)\lambda_n(t),$$

where $W^B(t)$ is a vector-valued Brownian motion in the spot measure, and we use an extended notation $w_{S_i,n}$ to emphasize which swap rate a given weight relates to. Notice the presence of drift terms, of order $O(dt)$. The quadratic variation and covariation of $S_1(t)$ and $S_2(t)$ satisfy

$$d\langle S_1(t), S_2(t) \rangle = z(t)\varphi(S_1(t))\varphi(S_2(t))\lambda_{S_1}(t)^\top\lambda_{S_2}(t) dt,$$

$$d\langle S_i(t) \rangle = z(t)\varphi(S_i(t))^2 \|\lambda_{S_i}(t)\|^2 dt, \quad i = 1, 2,$$

and the instantaneous correlation is

$$\text{Corr}(dS_1(t), dS_2(t)) = \frac{\lambda_{S_1}(t)^\top\lambda_{S_2}(t)}{\|\lambda_{S_1}(t)\| \|\lambda_{S_2}(t)\|}. \quad (14.34)$$

Instead of the instantaneous correlation, in many applications we are normally more interested in an estimate for *term correlation* $\rho_{\text{term}}(T', T)$ of S_1 and S_2 on some finite interval $[T', T]$. Formally, we define this time 0 measurable quantity as

$$\rho_{\text{term}}(T', T) \triangleq \text{Corr}(S_1(T) - S_1(T'), S_2(T) - S_2(T')).$$

Ignoring drift terms and freezing the swap rates at their time 0 forward levels, to decent approximation we can write

$$\begin{aligned} \rho_{\text{term}}(T', T) &\approx \frac{\varphi(S_1(0))\varphi(S_2(0)) \int_{T'}^T \mathbb{E}^B(z(t))\lambda_{S_1}(t)^\top\lambda_{S_2}(t) dt}{\varphi(S_1(0))\varphi(S_2(0)) \prod_{i=1}^2 \sqrt{\int_{T'}^T \mathbb{E}^B(z(t)) \|\lambda_{S_i}(t)\|^2 dt}} \\ &= \frac{\int_{T'}^T \lambda_{S_1}(t)^\top\lambda_{S_2}(t) dt}{\sqrt{\int_{T'}^T \|\lambda_{S_1}(t)\|^2 dt} \sqrt{\int_{T'}^T \|\lambda_{S_2}(t)\|^2 dt}}, \end{aligned} \quad (14.35)$$

where in the second equality we have used the fact that the parameterization (14.15) implies that, for all $t \geq 0$,

$$\mathbb{E}^B(z(t)) = z_0.$$

14.4.3.2 Spread Option Pricing

Consider a spread option paying at time $T \leq \min(T_{j_1}, T_{j_2})$

$$V_{\text{spread}}(T) = (S_1(T) - S_2(T) - K)^+,$$

such that

$$V_{\text{spread}}(0) = P(0, T) \mathbb{E}^T \left((S_1(T) - S_2(T) - K)^+ \right),$$

where, as always, \mathbb{E}^T denotes expectations in measure Q^T . An accurate (analytic) evaluation of this expected value is somewhat involved, and we postpone it until Chapter 17. Here, as a preview, we consider a cruder approach which may, in fact, be adequate for calibration purposes. We assume that the spread

$$\varepsilon(T) = S_1(T) - S_2(T)$$

is a Gaussian variable with mean

$$\mathbb{E}^T(\varepsilon(T)) = \mathbb{E}^T(S_1(T)) - \mathbb{E}^T(S_2(T)).$$

In a pinch, the mean of $\varepsilon(T)$ can be approximated as $S_1(0) - S_2(0)$, which assumes that the drift terms of $S_1(t)$ and $S_2(t)$ in the T -forward measure are approximately identical. For a better approximation, see Chapter 16. As for the variance of $\varepsilon(T)$, it can be approximated in several different ways, but one approach simply writes

$$\begin{aligned} \text{Var}^T(\varepsilon(T)) &\approx \sum_{i=1}^2 \varphi(S_i(0))^2 z_0 \int_0^T \|\lambda_{S_i}(t)\|^2 dt \\ &\quad - 2\rho_{\text{term}}(0, T) z_0 \prod_{i=1}^2 \varphi(S_i(0)) \left(\int_0^T \|\lambda_{S_i}(t)\|^2 dt \right)^{1/2}. \end{aligned} \quad (14.36)$$

With these approximations, the Bachelier formula (7.16) yields

$$V_{\text{spread}}(0) = P(0, T) \sqrt{\text{Var}^T(\varepsilon(T))} (d\Phi(d) + \phi(d)), \quad d = \frac{\mathbb{E}^T(\varepsilon(T)) - K}{\sqrt{\text{Var}^T(\varepsilon(T))}}. \quad (14.37)$$

14.5 Calibration

14.5.1 Basic Principles

Suppose that we have fixed the tenor structure, have decided upon the number of factors m to be used, and have selected the basic form (e.g. DVF or SV) of the LM model that we are interested in deploying. Suppose also, for now, that any skew functions and stochastic volatility dynamics have been exogenously specified by the user. To complete our model specification, what then remains unanswered is the fundamental question of how to establish the set of m -dimensional deterministic volatility vectors $\lambda_k(\cdot)$, $k = 1, 2, \dots, N-1$, that together determine the overall correlation and volatility structure of forward rates in the model.

As evidenced by the large number of different calibration approaches proposed in the literature, there are no precise rules for calibration of LM models. Still, certain common steps are nearly always invoked:

- Prescribe the basic form of $\|\lambda_k(t)\|$, either through direct parametric assumptions, or by introduction of discrete time- and tenor-grids.
- Use correlation information to establish a map from $\|\lambda_k(t)\|$ to $\lambda_k(t)$.
- Choose the set of observable securities against which to calibrate the model.
- Establish the norm to be used for calibration.
- Recover $\lambda_k(t)$ by norm optimization.

In the next few sections, we will discuss these steps in sequence. In doing so, our primary aim is to expose a particular calibration methodology that we personally prefer for most applications, rather than give equal mention to all possible approaches that have appeared in the literature. We note up front that our discussion is tilted towards applications that ultimately involve pricing and hedging of exotic Libor securities (see e.g. Chapters 18 and 19).

14.5.2 Parameterization of $\|\lambda_k(t)\|$

For convenience, let us write

$$\lambda_k(t) = h(t, T_k - t), \quad \|\lambda_k(t)\| = g(t, T_k - t), \quad (14.38)$$

for some functions $h : \mathbb{R}_+^2 \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}_+^2 \rightarrow \mathbb{R}_+$ to be determined. We focus on g in this section, and start by noting that many ad-hoc parametric forms for this function have been proposed in the literature. A representative example is the following 4-parameter specification, due to Rebonato [1998]:

$$g(t, x) = g(x) = (a + bx)e^{-cx} + d, \quad a, b, c, d \in \mathbb{R}_+. \quad (14.39)$$

We notice that this specification is *time-stationary* in the sense that $\|\lambda_k(t)\|$ does not depend on calendar time t , but only on the remaining time to maturity ($T_k - t$) of the forward rate in question. While attractive from a perspective of smoothness of model volatilities, assumptions of perfect time stationarity will generally not allow for a sufficiently accurate fit to market prices. To address this, some authors have proposed “separable” extensions of the type

$$g(t, x) = g_1(t)g_2(x), \quad (14.40)$$

where g_1 and g_2 are to be specified separately. See Brace et al. [1997] for an early approach along these lines.

For the applications we have in mind, relying on separability or parametric forms is ultimately too inflexible, and we seek a more general approach. For this, let us introduce a rectangular grid of times and tenors $\{t_i\} \times \{x_j\}$, $i = 1, \dots, N_t$, $j = 1, \dots, N_x$; and an $(N_t \times N_x)$ -dimensional matrix G . The elements $G_{i,j}$ will be interpreted as

$$g(t_i, x_j) = G_{i,j}. \quad (14.41)$$

When dimensioning the grid $\{t_i\} \times \{x_j\}$, we would normally⁸ require that $t_1 + x_{N_x} \geq T_N$, to ensure that all forward rates on the Libor forward curve are covered by the table; beyond this, there need not be any particular relationship between the grid and the chosen tenor structure, although we find it convenient to ensure that $t_i + x_j \in \{T_n\}$ as long as $t_i + x_j \leq T_N$ — a convention we adopt from now on. Note that the bottom right-hand corner of the grid contains Libor maturities beyond that of our tenor structure and is effectively redundant.

A few further comments on the grid-based approach above are in order. First, we notice that both time-stationary and separable specifications along the lines of (14.39) and (14.40) can be emulated closely as special cases of the grid-based approach. For instance, the parametric specification (14.39) would give rise to a matrix G where

$$G_{i,j} = (a + bx_j)e^{-cx_j} + d,$$

i.e. all rows would be perfectly identical. We also point out that free parameters to be determined here equate all non-superfluous elements in G . In practice N_t and N_x would often both be around 10–15, so even after accounting for the fact that the bottom-right corner of G is redundant, the total number of free parameters to be determined is potentially quite large. To avoid overfitting, additional regularity conditions must be imposed — an important point to which we return in Section 14.5.6.

⁸An alternative would be to rely on extrapolation.

14.5.3 Interpolation on the Whole Grid

Suppose that we have somehow managed to construct the matrix G in (14.41), i.e. we have uncovered $\|\lambda_k(t)\| = g(t, T_k - t)$ for the values of t and $x = T_k - t$ on the grid $\{t_i\} \times \{x_j\}$. The next step is to construct $\|\lambda_k(t)\|$ for all values of t and k , $k = 1, \dots, N - 1$.

It is common⁹ to assume that for each k , the function $\|\lambda_k(t)\|$ is piecewise constant in t , with discontinuities at T_n , $n = 1, \dots, k - 1$,

$$\|\lambda_k(t)\| = \sum_{n=1}^k 1_{\{T_{n-1} \leq t < T_n\}} \|\lambda_{n,k}\| = \sum_{n=1}^k 1_{\{q(t)=n\}} \|\lambda_{n,k}\|. \quad (14.42)$$

In this case, we are left with constructing the matrix $\|\lambda_{n,k}\|$ from G , for all $1 \leq n \leq k \leq N - 1$. This is essentially a problem of two-dimensional interpolation (and, perhaps, extrapolation if the $\{t_i\} \times \{x_j\}$ grid does not cover the whole tenor structure). Simple, robust schemes such as separate t - and x -interpolation of low order seem to perform well, whereas high-order interpolation (cubic or beyond) may lead to undesirable effects during risk calculations. In practice, one would normally use either piecewise constant or piecewise linear interpolation.

Suppose, for concreteness, that linear interpolation in both dimensions of G is chosen. Then for each n, k ($1 \leq n \leq k \leq N - 1$) we have the following scheme

$$\|\lambda_{n,k}\| = w_{++} G_{i,j} + w_{+-} G_{i,j-1} + w_{-+} G_{i-1,j} + w_{--} G_{i-1,j-1}, \quad (14.43)$$

where, denoting $\tau_{n,k} = T_k - T_{n-1}$, we have

$$\begin{aligned} i &= \min \{a : t_a \geq T_{n-1}\}, \quad j = \min \{b : x_b \geq \tau_{n,k}\}, \\ w_{++} &= \frac{(T_{n-1} - t_{i-1})(\tau_{n,k} - x_{j-1})}{(t_i - t_{i-1})(x_j - x_{j-1})}, \quad w_{+-} = \frac{(T_{n-1} - t_{i-1})(x_j - \tau_{n,k})}{(t_i - t_{i-1})(x_j - x_{j-1})}, \\ w_{-+} &= \frac{(t_i - T_{n-1})(\tau_{n,k} - x_{j-1})}{(t_i - t_{i-1})(x_j - x_{j-1})}, \quad w_{--} = \frac{(t_i - T_{n-1})(x_j - \tau_{n,k})}{(t_i - t_{i-1})(x_j - x_{j-1})}. \end{aligned}$$

Apart from the order of interpolation, we can also choose which type of volatilities we want to interpolate. To explain, let us recall from Chapters 7 and 8 that we often normalize the local volatility function in such a way that $\varphi(L_n(0)) \approx L_n(0)$. Then, $\|\lambda_k(\cdot)\|$'s have the dimensionality of log-normal, or percentage, volatilities, and (14.43) defines interpolation in *log-normal* Libor volatilities. This is not the only choice, and using volatilities that are scaled differently in the interpolation will sometimes lead to smoother and more robust results, as was the case during much of the 2007–2009 financial

⁹A more refined approach, especially for low values of time-to-maturity, is advisable for some applications where the fine structure of short-term volatilities is important. See the discussion in Remark 15.1.1.

crisis. To demonstrate the basic idea, let us fix p , $0 \leq p \leq 1$. Then we can replace (14.43) with

$$\begin{aligned} L_k(0)^{1-p} \|\lambda_{n,k}\| &= w_{++} L_{n(i,j)}(0)^{1-p} G_{i,j} + w_{+-} L_{n(i,j-1)}(0)^{1-p} G_{i,j-1} \\ &\quad + w_{-+} L_{n(i-1,j)}(0)^{1-p} G_{i-1,j} + w_{--} L_{n(i-1,j-1)}(0)^{1-p} G_{i-1,j-1}, \end{aligned} \quad (14.44)$$

where the indexing function $n(i,j)$ is defined by $T_{n(i,j)} = t_i + x_j$. For $p = 0$, this can be interpreted as interpolation in Gaussian volatilities (see Remark 7.2.9). For arbitrary p , the formula (14.44) specifies interpolation in “CEV” volatilities.

Finally, note that even if we use linear interpolation between the knot points (either in t or x or both), it is normally better to use *constant* extrapolation before the initial t_1 and x_1 and after the final t_{N_t} and x_{N_x} .

14.5.4 Construction of $\lambda_k(t)$ from $\|\lambda_k(t)\|$

Suppose the values of volatility norm $\|\lambda_{n,k}\|$ are known on the full grid $1 \leq n \leq k \leq N-1$. For each T_n , the components of the m -dimensional $\lambda_k(T_n)$ vectors may now be obtained from instantaneous Libor rate volatilities $\|\lambda_{n,k}\|$ for $k \geq n$, and instantaneous correlations of Libor rates fixing on or after T_n . The procedure is similar in spirit to the one we employed previously for parameterizing multi-factor Gaussian short rate models in Section 12.1.7. So, with the calendar time fixed at some value T_n , we introduce an $(N-n) \times (N-n)$ instantaneous correlation matrix $R(T_n)$, with elements

$$(R(T_n))_{i,j} = \text{Corr}(dL_i(T_n-), dL_j(T_n-)), \quad i, j = n, \dots, N-1.$$

The correlation matrix would, in many applications, be computed from an estimated parametric form, such as those covered in Section 14.3.2. Furthermore, we define a diagonal volatility matrix $c(T_n)$ with elements $\|\lambda_{n,n}\|, \|\lambda_{n,n+1}\|, \dots, \|\lambda_{n,N-1}\|$ along its diagonal. That is,

$$(c(T_n))_{j,j} = \|\lambda_{n,n+j-1}\|, \quad j = 1, \dots, N-n,$$

with all other elements set to zero. Given $R(T_n)$ and $c(T_n)$, an instantaneous covariance matrix¹⁰ $C(T_n)$ for forward rates on the grid can now be computed as

$$C(T_n) = c(T_n)R(T_n)c(T_n). \quad (14.45)$$

Let us define $H(T_n)$ to be an $(N-n) \times m$ matrix composed by stacking each dimension of $h(T_n, T_{n+j-1} - T_n)$ (see 14.38) side by side, with j running on the grid:

¹⁰Earlier results show that the true instantaneous covariance matrix for forward rates may involve DVF- or SV-type scales on the elements of c . For the purposes of calibration of λ_k , we omit these scales.

$$(H(T_n))_{j,i} = h_i(T_n, T_{n+j-1} - T_n), \quad j = 1, \dots, N-n, \quad i = 1, \dots, m.$$

Then, it follows that we should have

$$C(T_n) = H(T_n)H(T_n)^\top. \quad (14.46)$$

Equations (14.45) and (14.46) specify two different representations of the covariance matrix, and we want them to be identical, i.e.

$$H(T_n)H(T_n)^\top = c(T_n)R(T_n)c(T_n), \quad (14.47)$$

which gives us a way to construct the $H(T_n)$ matrix, and thereby the vectors $h(T_n, T_{n+j-1})$ for all values of n, j on the full grid $1 \leq n \leq N-1$, $1 \leq j \leq N-n$. Assuming, as before, piecewise constant interpolation of $\lambda_k(t)$ for t between knot dates $\{T_i\}$, the full set of factor volatilities $\lambda_k(t)$ can be constructed for all t and T_k .

As written, equation (14.47) will normally *not* have a solution as the left-hand side is rank-deficient, whereas the right-hand side will typically have full rank. To get around this, we can proceed to apply PCA methodology, in several different ways. We discuss two methods below, but first quickly note that for n close to N (in particular for $N-n < m$), the equation (14.47) will have *too many* solutions. A pragmatic approach here is to zero out the last few (namely, $m - (N-n)$) columns of the matrix $H(T_n)$ before solving the equation, in effect “forbidding” Brownian motions with high index affecting remaining Libor rates. We trust the reader can fill in the details of this scheme, and will ignore this slight complication going forward.

14.5.4.1 Covariance PCA

In this approach, we apply PCA decomposition to the entire right-hand side of (14.47), writing

$$c(T_n)R(T_n)c(T_n) \approx e_m(T_n)\Lambda_m(T_n)e_m(T_n)^\top,$$

where $\Lambda_m(T_n)$ is an $m \times m$ diagonal matrix of the m largest eigenvalues of $c(T_n)R(T_n)c(T_n)$, and $e_m(T_n)$ is an $(N-n) \times m$ matrix of eigenvectors corresponding to these eigenvalues. Inserting this result into (14.47) leads to

$$H(T_n) = e_m(T_n)\sqrt{\Lambda_m(T_n)}. \quad (14.48)$$

As discussed in Chapter 3, this approximation is optimal in the sense of minimizing the Frobenius norm of the covariance matrix errors.

14.5.4.2 Correlation PCA

An attractive alternative to the approach in Section 14.5.4.1 uses the correlation PCA decomposition discussed in Section 14.3.4. Here we write

$$R(T_n) = D(T_n)D(T_n)^\top, \quad (14.49)$$

for an $(N - n) \times m$ matrix D found by the techniques in Section 14.3.4. Inserting this into (14.47) yields

$$H(T_n) = c(T_n)D(T_n). \quad (14.50)$$

In computing the matrix D , we would normally use the result from Proposition 14.3.2, which would minimize the Frobenius norm on correlation matrix errors.

14.5.4.3 Discussion and Recommendation

Several papers in the literature focus on the method in Section 14.5.4.1 (e.g. Sidenius [2000], and Pedersen [1998]), but we nevertheless strongly prefer the approach in Section 14.5.4.2 for calibration applications. Although performing the PCA decomposition (as in Proposition 14.3.2) of a correlation matrix is technically more difficult than the same operation on a covariance matrix, the correlation PCA is independent of the c matrix and as such will not have to be updated when we update guesses for the G matrix (on which c depends) in a calibration search loop. When the correlation matrix $R(T_n)$ originates from a parametric form independent of calendar time (which we recommend), the matrix D in (14.49) will, in fact, need estimation only once per tenor date¹¹ T_n , at a minimal computational overhead cost. In comparison, the covariance PCA operation will have to be computed at each T_n every time G gets updated in the calibration loop. We also notice that $D(T_n)D(T_n)^\top$ having a unit diagonal will automatically ensure that the total forward rate volatility will be preserved if m is changed; this is *not* the case for covariance PCA, where the total volatility of forward rates will normally increase as m is increased, *ceteris paribus*.

If the complexity of the optimal PCA algorithm in Proposition 14.3.2 of Section 14.3.4 is deemed too egregious, the simplified approach of Section 14.3.4.2 could be used instead. It shares the performance advantages of the “true” correlation PCA as it only needs to be run once outside the calibration loop, but its theoretical deficiencies suggest that its use should, in most circumstances, be limited to the case where the correlations are themselves calibrated, rather than exogenously specified by the user. We return to the concept of correlation calibration in Section 14.5.9.

14.5.5 Choice of Calibration Instruments

In a standard LM model calibration, we choose a set of swaptions and caps (and perhaps Eurodollar options) with market-observable prices; these prices

¹¹Since the matrix D shrinks in T_n , we need to repeat the PCA analysis at each tenor date. Alternatively, but suboptimally, we can do the PCA analysis only once at time 0, pruning the results as needed for other values of T_n .

serve as calibration targets for our model. The problem of determining precisely which caps and swaptions should be included in the calibration is a difficult and contentious one, with several opposing schools of thought represented in the literature. We shall spend this section¹² outlining the major arguments offered in the literature as well as our own opinion on the subject. Before commencing on this, we emphasize that the calibration algorithm we develop in this book accommodates arbitrary sets of calibration instruments and as such will work with any selection philosophy.

One school of thought — the *fully calibrated* or *global* approach — advocates calibrating an LM model to a large set of available interest options, including both caps and swaptions in the calibration set. When using grid-based calibration, this camp would typically recommend using at-the-money swaptions with maturities and tenors chosen to coincide with each point in the grid. That is, if T_s is the maturity of a swaption and T_e is the end date of its underlying swap, then we would let T_s take on all values in the time grid $\{t_i\}$, while at the same time letting $T_e - T_s$ progress through all values¹³ of the tenor grid $\{x_j\}$. On top of this, one would often add at-the-money caps at expiries ranging from $T = t_1$ to $T = t_{N_t}$.

The primary advantage of the fully calibrated approach is that a large number of liquid volatility instruments are consistently priced within the model. This, in turn, gives us some confidence that the vanilla option market is appropriately “spanned” and that the calibrated model can be used on a diverse set of exotic securities. In vega hedging (see Section 8.9.1 for definition and Chapter 26 for much more on vega hedging in LM models) of an exotic derivative, one will undoubtedly turn to swaptions and caps, so mispricing these securities in the model would be highly problematic.

Another school of thought — the *parsimonious* or *local* approach — judiciously chooses a small subset of caps and swaptions in the market, and puts significant emphasis on specification of smooth and realistic term structures of forward rate volatilities. Typically this will involve imposing strong time-homogeneity assumptions, or observed statistical relationships, on the $\lambda_k(\cdot)$ vectors. The driving philosophy behind the parsimonious approach (besides the desire for calibration speed) is the observation that, fundamentally, the price of a security in a model is equal to the model-predicted cost of hedging the security over its lifetime. Hedging profits in the future as specified by the model are, in turn, directly related to the forward rate volatility structures that the model predicts for the future. For these model-predicted hedging profits to have any semblance to the actual realized hedging profits, the dynamics of the volatility structure in the model should be a reasonable

¹²We also revisit the subject in the context of callable Libor exotics in Section 18.1.

¹³One would here limit T_s to be no larger than T_N , so the total number of swaptions would be less than $N_t \cdot N_x$. See our discussion of redundant grid entries in Section 14.5.2.

estimate of the actual dynamics. In many cases, however, our best estimate of future volatility structures might be today's volatility structures (or those we have estimated historically), suggesting that the evolution of volatility should be as close to being time-homogeneous as possible. This can be accomplished, for instance, by using time-homogeneous mappings such as (14.39) or similar.

The strong points of the parsimonious approach are, of course, weak ones of the fully calibrated approach. Forward rate volatilities produced by the fully calibrated model can easily exhibit excessively non-stationary behavior, impairing the performance of dynamic hedging. On the other hand, the inevitable mispricings of certain swaptions and/or caps in the parsimonious approach are troublesome. In a pragmatic view of a model as a (sophisticated, hopefully) interpolator that computes prices of complex instruments from prices of simple ones, mispricing of simple instruments obviously does not inspire confidence in the prices returned for complex instruments. As discussed, the parsimonious approach involves an attempt to identify a small enough set of "relevant" swaptions and caps that even a time-homogeneous model with a low number of free parameters can fit reasonably well, but it can often be very hard to judge which swaption and cap volatilities are important for a particular exotic security. In that sense, a fully calibrated model is more universally applicable, as the need to perform trade-specific identification of a calibration set is greatly reduced. Notice also that the risk profile of a given security may change greatly over time as market rates move around, potentially necessitating the use of *different* calibration instruments over time. Changing the calibration instrument set will obviously trigger a discontinuity in the hedge strategy, which is not ideal.

It is easy to imagine taking both approaches to the extremes to generate results that would convincingly demonstrate the perils of using either of them. To avoid such pitfalls we recommend looking for an equilibrium between the two. While we overall favor the fully calibrated approach, it is clear that, at the very least, it should be supplemented by an explicit mechanism to balance price precision versus regularity (e.g. smoothness and time-homogeneity) of the forward rate volatility functions. In addition, one should always perform rigorous checks of the effects of calibration assumptions on pricing and hedging results produced by the model. These checks should cover, at a minimum, result variations due to changes in

- Number of factors used (m).
- Relative importance of recovering all cap/swaption prices vs. time-homogeneity of the resulting volatility structure.
- Correlation assumptions.

A final question deserves a brief mention: should one calibrate to either swaptions or caps, or should one calibrate to both simultaneously? Followers of the parsimonious approach will often argue that there is a persistent

basis between cap and swaption markets, and any attempt to calibrate to both markets simultaneously is bound to distort the model dynamics. Instead, it is argued, one should only calibrate to one of the two markets, based on an analysis of whether the security to be priced is more cap- or swaption-like. Presumably this analysis would involve judging whether either caps or swaptions will provide better vega hedges for the security in question. The drawback of this approach is obvious: many complicated interest rates securities depend on the evolution of both Libor rates as well as swap rates and will simultaneously embed “cap-like” and “swaption-like” features.

To avoid discarding potentially valuable information from either swaption or cap markets, we generally recommend that both markets be considered in the calibration of the LM model. However, we do not necessarily advocate that both types of instruments receive equal weighting in the calibration objective function; rather, the user should be allowed some mechanism to affect the relative importance of the two markets. We return to this idea in the next section.

14.5.6 Calibration Objective Function

As discussed above, several issues should be considered in the choice of a calibration norm, including the smoothness and time-stationarity of the $\lambda_k(\cdot)$ functions; the precision to which the model can replicate the chosen set of calibration instruments; and the relative weighting of caps and swaptions. To formally state a calibration norm that will properly encompass these requirements, assume that we have chosen calibration targets that include N_S swaptions, $V_{\text{swaption},1}, V_{\text{swaption},2}, \dots, V_{\text{swaption},N_S}$, and N_C caps, $V_{\text{cap},1}, V_{\text{cap},2}, \dots, V_{\text{cap},N_C}$. Strategies for selecting these instruments were discussed in the previous section. We let \bar{V} denote their quoted market prices and, adopting the grid-based framework from Section 14.5.2, we let $\bar{V}(G)$ denote their model-generated prices as functions of the volatility grid G defined in Section 14.5.2. We introduce a calibration objective function \mathcal{I} as

$$\begin{aligned} \mathcal{I}(G) = & \frac{w_S}{N_S} \sum_{i=1}^{N_S} \left(\bar{V}_{\text{swaption},i}(G) - \hat{V}_{\text{swaption},i} \right)^2 \\ & + \frac{w_C}{N_C} \sum_{i=1}^{N_C} \left(\bar{V}_{\text{cap},i}(G) - \hat{V}_{\text{cap},i} \right)^2 \\ & + \frac{w_{\partial t}}{N_x N_t} \sum_{i=1}^{N_t} \sum_{j=1}^{N_x} \left(\frac{\partial G_{i,j}}{\partial t_i} \right)^2 + \frac{w_{\partial x}}{N_x N_t} \sum_{i=1}^{N_t} \sum_{j=1}^{N_x} \left(\frac{\partial G_{i,j}}{\partial x_j} \right)^2 \\ & + \frac{w_{\partial t^2}}{N_x N_t} \sum_{i=1}^{N_t} \sum_{j=1}^{N_x} \left(\frac{\partial^2 G_{i,j}}{\partial t_i^2} \right)^2 + \frac{w_{\partial x^2}}{N_x N_t} \sum_{i=1}^{N_t} \sum_{j=1}^{N_x} \left(\frac{\partial^2 G_{i,j}}{\partial x_j^2} \right)^2, \quad (14.51) \end{aligned}$$

where $w_S, w_C, w_{\partial t}, w_{\partial x}, w_{\partial t^2}, w_{\partial x^2} \in \mathbb{R}_+$ are exogenously specified weights. In (14.51) the various derivatives of the elements in the table G are, in practice, to be interpreted as discrete difference coefficients on neighboring table elements — see (14.52) below for an example definition¹⁴.

As we have defined it, $\mathcal{I}(G)$ is a weighted sum of i) the mean-squared swaption price error; ii) the mean-squared cap price error; iii) the mean-squared average of the derivatives of G with respect to calendar time; iv) the mean-squared average of the second derivatives of G with respect to calendar time; v) the mean-squared average of the derivatives of G with respect to forward rate tenor; and vi) the mean-squared average of the second derivatives of G with respect to forward rate tenor. The terms in i) and ii) obviously measure how well the model is capable of reproducing the supplied market prices, whereas the remaining four terms are all related to regularity. The term iii) measures the degree of volatility term structure time homogeneity and penalizes volatility functions that vary too much over calendar time. The term iv) measures the smoothness of the calendar time evolution of volatilities and penalizes deviations from linear evolution (a straight line being perfectly smooth). Terms v) and vi) are similar to iii) and iv) and measure constancy and smoothness in the tenor direction. In (14.51), the six weights $w_S, w_C, w_{\partial t}, w_{\partial x}, w_{\partial t^2}, w_{\partial x^2}$ determine the trade-off between volatility smoothness and price accuracy, and are normally to be supplied by the user based on his or her preferences. In typical applications, the most important regularity terms are those scaled by the weights $w_{\partial t}$ and $w_{\partial x^2}$ which together determine the degree of time homogeneity and tenor smoothness in the resulting model.

We should note that there are multiple ways to specify smoothness criteria, with (14.51) being one of many. For example, as we generalized the basic log-normal interpolation scheme (14.43) to allow for interpolation in “CEV” volatilities in (14.44), we can adjust the definition of smoothness to be in terms of compatible quantities. In particular, instead of using

$$\frac{w_{\partial x}}{N_x N_t} \sum_{i=1}^{N_t} \sum_{j=2}^{N_x} \left(\frac{G_{i,j} - G_{i,j-1}}{x_j - x_{j-1}} \right)^2 \quad (14.52)$$

as implicit in (14.51) for the tenor-smoothness term, we could use

$$\frac{w_{\partial x}}{N_x N_t} \sum_{i=1}^{N_t} \sum_{j=2}^{N_x} \left(\frac{L_{n(i,j)}(0)^{1-p} G_{i,j} - L_{n(i,j-1)}(0)^{1-p} G_{i,j-1}}{x_j - x_{j-1}} \right)^2 \quad (14.53)$$

for some p , $0 \leq p \leq 1$. The case of $p = 0$ would then correspond to smoothing basis-point Libor volatilities rather than log-normal Libor volatilities.

As written, the terms of the calibration norm that measure precision in cap and swaption pricing involve mean-squared errors directly on prices.

¹⁴Depending on how table boundary elements are treated, notice that the range for i and j may not always be as stated in (14.51).

In practice, however, the error function is often applied to some transform of outright prices, e.g. implied volatilities. For an SV-type LM model, for instance, we could institute a pre-processing step where the market price of each swaption $\widehat{V}_{\text{swaption},i}$ would be converted into a constant implied volatility $\widehat{\lambda}_{S_i}$, in such a way that the scalar SDE for the swap rate S_i underlying the swaption $V_{\text{swaption},i}$,

$$dS_i(t) = \sqrt{z(t)\widehat{\lambda}_{S_i}}\varphi(S_i(t))dY_i(t),$$

would reproduce the observed swaption market price. Denoting by $\bar{\lambda}_{S_i}(G)$ the corresponding model volatility of the swap rate S_i (as given by, for example, Proposition 14.4.4) and repeating this exercise for all caps and swaptions in the calibration set, we obtain an alternative calibration norm definition where the cap and swaption terms in (14.51) are modified as follows:

$$\mathcal{I}(G) = \frac{w_S}{N_S} \sum_{i=1}^{N_S} \left(\bar{\lambda}_{S_i}(G) - \widehat{\lambda}_{S_i} \right)^2 + \frac{w_C}{N_C} \sum_{i=1}^{N_C} \left(\bar{\lambda}_{C_i}(G) - \widehat{\lambda}_{C_i} \right)^2 + \dots \quad (14.54)$$

The advantage of working with implied volatilities in the precision norm is two-fold. First, the relative scaling of individual swaptions and caps is more natural; when working directly with prices, high-value (long-dated) trades would tend to be overweighed relative to low-value (short-dated) trades¹⁵. Second, in many models computation of the implied volatility terms $\bar{\lambda}_{S_i}$ and $\bar{\lambda}_{C_i}$ can often be done by simple time integration of (combinations of) $\lambda_k(\cdot)$ (see e.g. Proposition 14.4.4) avoiding the need to apply a possibly time-consuming option pricing formula to compute the prices $\overline{V}_{\text{swaption},i}$ and $\overline{V}_{\text{cap},i}$. Considerable attention to this particular issue was paid in Section 9.3 (for SV models) and Section 7.6.2 (for DVF models), and we review relevant results and apply them to LM models in Chapter 15.

The quality-of-fit objective can be expressed in terms of *scaled* volatilities, which sometimes improves performance. Following the ideas developed for interpolation (14.44) and smoothing (14.53), we could express the fit objective as

$$\mathcal{I}(G) = \frac{w_S}{N_S} \sum_{i=1}^{N_S} \left(S_i(0)^{1-p} \left(\bar{\lambda}_{S_i}(G) - \widehat{\lambda}_{S_i} \right) \right)^2 + \dots,$$

for a given p , $0 \leq p \leq 1$. Taking this idea further we note that a more refined structure of mean-squared weights in the definition of calibration norm is possible. For instance, rather than weighting all swaptions equally with the term w_S , one could use different weights for each swaption in the calibration set. Similarly, by using node-specific weights on the derivatives of the entries in G one may, say, express the view that time homogeneity is more important for large t than for small t .

¹⁵Another approach to producing more equitable scaling involves using relative (=percentage) price errors, rather than absolute price errors.

14.5.7 Sample Calibration Algorithm

At this point, we are ready to state our full grid-based calibration algorithm. We assume that a tenor structure and a time/tenor grid $\{t_i\} \times \{x_j\}$ have been selected, as have the number of Brownian motions (m), a correlation matrix R , and the set of calibration swaptions and caps. In addition, the user must select the weights in the calibration norm \mathcal{I} in (14.51) or (14.54). Starting from some guess for G , we then run the following iterative algorithm:

1. Given G , interpolate using (14.43) or (14.44) to obtain the full norm volatility grid $\|\lambda_{n,k}\|$ for all Libor indices $k = 1, \dots, N - 1$ and all expiry indices $n = 1, \dots, k$.
2. For each $n = 1, \dots, N - 1$, compute the matrix $H(T_n)$, and ultimately volatility loadings $\lambda_k(T_n)$, from $\|\lambda_{n,k}\|$, $k \geq n$, by PCA methodology, using either (14.48) or (14.50).
3. Given $\lambda_k(\cdot)$ for all $k = 1, \dots, N - 1$, use the formulas in Sections 14.4.1 and 14.4.2 to compute model prices for all swaptions and caps in the calibration set.
4. Establish the value of $\mathcal{I}(G)$ by direct computation of either (14.51) or (14.54).
5. Update G and repeat Steps 1–4 until $\mathcal{I}(G)$ is minimized.

Step 5 in the above algorithm calls for the use of a robust high-dimensional numerical optimizer. Good results can, in our experience, be achieved with several algorithms, including the Spellucci algorithm¹⁶, the Levenberg-Marquardt algorithm, and the downhill simplex method (the last two can be found in Press et al. [1992]). These, and many alternative algorithms, are available in standard numerical packages, such as IMSL¹⁷ and NAG¹⁸. On a standard PC, a well-implemented calibration algorithm should generally complete in about 10 seconds from a cold start (i.e. where we do not have a good initial guess for G) for, say, a 40 year model with quarterly Libor rolls.

14.5.8 Speed-Up Through Sub-Problem Splitting

An LM model calibration problem involves a substantial number of free input variables to optimize over, namely all elements of the matrix G . In a typical setup, the number of such variables may range from a few dozen to a few hundred. As the number of terms, or “targets”, in the calibration norm is of the same order of magnitude, we are dealing with a fairly sizable optimization problem. While modern optimization algorithms implemented on modern hardware can successfully handle the full-blown problem, it is still

¹⁶donlp2 SQP/ECQP algorithm, available on www.mathematik.tu-darmstadt.de:8080/ags/ag8/Mitglieder/spellucci_de.html.

¹⁷www.imsl.com.

¹⁸www.nag.com.

of interest to examine whether there are ways of to improve computational efficiency. For instance, if we could split the optimization problem into a sequence of smaller sub-problems solved separately and sequentially, the performance of the algorithm would typically improve. Indeed, imagine for illustrative purposes that we have an optimization problem with $m = m_1 m_2$ variables and computational complexity of the order¹⁹ $O(m^2) = O(m_1^2 m_2^2)$. However, if we could find the solution by sequentially solving m_1 problems of m_2 variables each, then the computational cost would be $m_1 O(m_2^2)$, yielding savings of the order $O(m_1)$.

Our ability to split the problem into sub-problems typically relies on exploring its particular structure, i.e. the relationship between input variables and targets. If, for example, target 1 depends on variable 1 but not — or only mildly — on other variables, then it makes sense to find the optimal value for the variable 1 by optimizing for target 1 while keeping other variables constant, and so on. Fortunately, the LM model optimization problem presents good opportunities for this type of analysis. First, recall that the main calibration targets for the problem are the differences in market and model prices (or implied volatilities) of caps and swaptions. Let us consider a swaption with expiry T_j and final payment date T_n ; let i be such that $T_j = t_i$. Then, as follows from the swaption approximation formula (14.33), the model volatility for this swaption depends on $\lambda_k(t)$'s for $t \in [0, t_i]$ and for $k = j, \dots, n - 1$. Hence, the part of the calibration norm associated with the price fit of the swaption will depend on the first i rows of the matrix G *only*. This observation suggests splitting the calibration problem into a collection of “row by row” calibration problems.

To simplify notations, we assume that the set of fit targets consists of *all* swaptions with expiries t_i and tenors x_l , $i = 1, \dots, N_t$, $l = 1, \dots, N_x$. In a row-by-row calibration algorithm, the first row of the matrix G is calibrated to all N_x swaptions with expiry t_1 , then the second row of G is calibrated to the swaptions with expiry t_2 , and so on.

As we emphasized earlier, having regularity terms in the calibration norm is important to ensure a smooth solution. Fortunately, regularity terms can generally be organized in the same row-by-row format as the precision terms. For instance, the regularity terms in the tenor direction naturally group into row-specific collections. As for the terms controlling the regularity of the matrix G in calendar time t , when optimizing on time slice t_i , we would only include in the norm the terms that involve rows of G with an index less than or equal to i . We trust that the reader can see how to arrange this, and omit straightforward details.

¹⁹As many of the algorithms we have in mind compute, at the very least, the sensitivity of each calibration target to each input variable, the computational complexity is at least $O(m^2)$; if the order of complexity is higher, the case for problem splitting is even more compelling.

Computational savings from the row by row scheme could be substantial — for a 40 year model with quarterly Libor rolls, a well-tuned algorithm should converge in less than a second or two. There are, however, certain drawbacks associated with problem splitting. In particular, as the calibration proceeds from one row to the next, the optimizer does not have the flexibility to adjust previous rows of the matrix G to the current row of swaption volatilities. This may result in a tell-tale “ringing” pattern of the Libor volatilities in the time direction, as the optimizer attempts to match each row of price targets through excessively large moves in the elements of G , in alternating up and down directions. Judicious application of regularity terms in the optimization norm can, however, help control this behavior, and overall the row-by-row scheme performs well. We recommend it as the default method for most applications, but note that sometimes a combination of full-blown and row-by-row calibration is the best choice.

Returning to the row-by-row calibration idea, one can try to take it further and split the calibration to an ever-finer level, eventually fitting each individual price target — a given caplet or swaption volatility, say — separately, by moving just a *single* element of the matrix G . This should seemingly work because the (t_i, x_{l+1}) -swaption volatility depends on the same elements of matrix G as the (t_i, x_l) -swaption volatility *plus* $G_{i,l+1}$. (This is not entirely true due to some grid interpolation effects, but the general idea is correct). So, in principle, $G_{i,l+1}$ can be found by just solving a quadratic equation, i.e. in closed form. For full details we refer the reader to Brigo and Mercurio [2001] where this *bootstrap*, or *cascade*, algorithm is described in detail. While this may appear to be a strong contender for practical LM calibration — full calibration is performed by just solving $N_t N_x$ quadratic equations — the scheme generally does not work for practically-sized problems. The cascade calibration suffers strongly from the ringing problem discussed above, and the quadratic equations typically fail to have a solution for swaption targets with just a few years of total maturity (i.e. from today to the final payment date). While it is possible to include regularity terms that preserve the closed-form nature of the solution, the ringing problem is difficult to remedy and calibration to long-dated options is rarely feasible. We find this to be true, even if one applies ad-hoc remediation methods proposed by various authors (see e.g. Brigo and Morini [2006]).

We should note that bootstrap calibration does have certain uses. For instance, one could use full-blown (or row by row) optimization to fundamentally calibrate G , and then use some version of bootstrap calibration to examine the effect of making small perturbations to input prices, e.g. when computing vegas. We discuss this idea in Chapter 26.

14.5.9 Correlation Calibration to Spread Options

In the calibration algorithm in Section 14.5.7, the matrix R was specified exogenously and would typically originate from an empirical analysis similar

to that in Section 14.3.2. As we discussed in Section 14.4.3, an alternative approach attempts to imply R directly from market data for spread options. Less is known about the robustness of calibrations based on this approach, but this shall not prevent us from listing a possible algorithm.

First, to make the problem tractable, we assume that the matrix R is time-homogeneous and specified as some parametric function of a low-dimension parameter-vector ξ ,

$$R = R(\xi).$$

Possible parameterizations include those listed in Section 14.3.2. We treat ξ as an unknown vector, to be determined in the calibration procedure along with the elements of the volatility matrix G . For this, we introduce a set of market-observable spread option prices $\widehat{V}_{\text{spread},1}, \widehat{V}_{\text{spread},2}, \dots, \widehat{V}_{\text{spread},N_{SP}}$, their corresponding model-based prices $\overline{V}_{\text{spread},1}(G, \xi), \overline{V}_{\text{spread},2}(G, \xi), \dots, \overline{V}_{\text{spread},N_{SP}}(G, \xi)$, and update the norm \mathcal{I} in (14.51) (or (14.54)) to the norm $\mathcal{I}^*(G, \xi)$, where²⁰

$$\mathcal{I}^*(G, \xi) = \mathcal{I}(G, \xi) + \frac{w_{SP}}{N_{SP}} \sum_{i=1}^{N_{SP}} \left(\overline{V}_{\text{spread},i}(G, \xi) - \widehat{V}_{\text{spread},i} \right)^2. \quad (14.55)$$

The algorithm in Section 14.5.7 proceeds as before with a few obvious changes; we list the full algorithm here for completeness.

1. Given G , interpolate using (14.43) or (14.44) to obtain the full norm volatility grid $\|\lambda_{n,k}\|$ for all Libor indices $k = 1, \dots, N-1$ and all expiry indices $n = 1, \dots, k$.
2. Given ξ , compute $R = R(\xi)$.
3. For each $n = 1, \dots, N-1$ and using $R(\xi)$, compute the matrix H , and ultimately volatility loadings $\lambda_k(T_n)$, from $\|\lambda_{n,k}\|$, $k \geq n$, by PCA methodology, using either (14.48) or (14.50).
4. Given $\lambda_k(\cdot)$ for all $k = 1, \dots, N-1$, use the formulas in Sections 14.4.1, 14.4.2 and 14.4.3 to compute model prices for all swaptions, caps and spread options in the calibration set.
5. Establish the value of $\mathcal{I}^*(G, \xi)$ by direct computation of (14.55).
6. Update G and ξ and repeat Steps 1–5 until $\mathcal{I}(G, \xi)$ is minimized.

When using a correlation PCA algorithm in Step 3, in practice one may find that it is most efficient to use the “poor man’s” approach in Section 14.3.4.2, rather than the slower expression listed in Proposition 14.3.2. Indeed, as long as the spread option prices ultimately are well-matched, we can be confident that our model has a reasonable correlation structure, irrespective of which PCA technique was used.

As was the case for our basic algorithm, let us note that it may be useful to transform spread option prices into implied volatilities or, even better,

²⁰Note that our original norm \mathcal{I} now also is a function of ξ , since cap and swaption prices depend on the correlation matrix R .

into implied term correlations²¹ when evaluating the mean-squared error. For spread options, a definition of implied term correlation can be extracted from the simple Gaussian spread approach in Section 14.4.3, equations (14.36) and (14.37) or, for more accurate formulas, using the results of Chapter 17 and in particular Sections 17.4.2 and 17.9.1.

Finally, we should note that the optimization problem embedded in the algorithm above can be quite challenging to solve in practice. To stabilize the numerical solution, it may be beneficial to employ a split calibration approach, where we first freeze correlation parameters ξ and then optimize G over the parts of the calibration norm that do not involve spread options. Then we freeze G at its optimum and optimize ξ over the parts of the calibration norm that do not involve caps and swaptions. This alternating volatility- and correlation-calibration is then repeated iteratively until (hopefully) convergence. A similar idea can be employed when calibrating models to a volatility smile; see Section 15.2.3 for LM model applications and Section 16.2.3 for applications to vanilla models.

14.5.10 Volatility Skew Calibration

The calibration algorithm we have discussed so far will normally take at-the-money options as calibration targets when establishing the $\lambda_k(t)$ functions. Establishing the volatility smile away from at-the-money strikes must be done in a separate step, through specification of a DVF skew function φ and, possibly, a stochastic volatility process $z(t)$. For the time-stationary specifications of these two mechanisms that we considered in Section 14.2.5, best-fitting to the volatility skew can be done relatively easily — in fact, it is probably often best to leave the parameters²² of the skew function φ as a free parameter for trader's input. We study the problem of volatility skew calibration for LM models in more detail in Chapter 15.

14.6 Monte Carlo Simulation

Once an LM model has been calibrated to market data, we can proceed to use the parameterized model for the pricing and risk management of non-vanilla options. In virtually all cases, pricing of such options will involve numerical methods. As the LM model involves a very large number of Markov state

²¹By representing spread options through implied term correlations, the information extracted from spread options is more “orthogonal” to that extracted from caps and swaptions, something that can help improve the numerical properties of the calibration algorithm, particularly if split calibration approach is used.

²²Assuming that there are only a few parameters that define the shape of the function. We generally recommend using simple skew functions that can be described by a single-parameter family, such as linear or power functions.

variables — namely the full number of Libor forward rates on the yield curve plus any additional variables used to model stochastic volatility — finite difference methods are rarely applicable (but see the brief discussion in Section 15.3 for a special case), and we nearly always have to rely on Monte Carlo methods. As we discussed in Chapter 3, the main idea of Monte Carlo pricing is straightforward: i) simulate independent paths of the collection of Libor rates through time; ii) for each path, sum the numeraire-deflated values of all cash flows generated by the specific interest rate dependent security at hand; iii) repeat i)-ii) many times and form the average. Proper execution of step i) is obviously key to this algorithm, and begs an answer to the following question: given a probability measure and the state of the Libor forward curve at time t , how do we move the entire Libor curve (and the numeraire) forward to time $t + \Delta$, $\Delta > 0$, in a manner that is consistent with the LM model dynamics? We address this question here.

14.6.1 Euler-Type Schemes

Assume that we stand at time t , and have knowledge of forward Libor rates maturing at all dates in the tenor structure after time t . We wish to devise a scheme to advance time to $t + \Delta$ and construct a sample of $L_{q(t+\Delta)}(t + \Delta), \dots, L_{N-1}(t + \Delta)$. Notice that $q(t + \Delta)$ may or may not exceed $q(t)$; if it does, some of the front-end forward rates expire and “drop off” the curve as we move to $t + \Delta$.

For concreteness, assume for now that we work in the spot measure Q^B in which case Lemma 14.2.3 tells us that general LM model dynamics are of the form

$$dL_n(t) = \sigma_n(t)^\top (\mu_n(t) dt + dW^B(t)), \quad \mu_n(t) = \sum_{j=q(t)}^n \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)}, \quad (14.56)$$

where the $\sigma_n(t)$ are adapted vector-valued volatility functions and $W^B(t)$ is an m -dimensional Brownian motion in measure Q^B . The simplest way of drawing an approximate sample $\widehat{L}_n(t + \Delta)$ for $L_n(t + \Delta)$ would be to apply a first-order Euler-type scheme. Applying results from Section 3.2.3, Euler (14.57) and log-Euler (14.58) schemes for (14.56) are, for $n = q(t + \Delta), \dots, N - 1$,

$$\widehat{L}_n(t + \Delta) = \widehat{L}_n(t) + \sigma_n(t)^\top (\mu_n(t)\Delta + \sqrt{\Delta}Z), \quad (14.57)$$

$$\widehat{L}_n(t + \Delta) = \widehat{L}_n(t) \exp \left\{ \frac{\sigma_n(t)^\top}{\widehat{L}_n(t)} \left(\left(\mu_n(t) - \frac{1}{2} \frac{\sigma_n(t)}{\widehat{L}_n(t)} \right) \Delta + \sqrt{\Delta}Z \right) \right\}, \quad (14.58)$$

where Z is a vector of m independent $\mathcal{N}(0, 1)$ Gaussian draws²³. For specifications of $\sigma_n(t)^\top$ that are close to proportional in $L_n(t)$ (e.g. the log-normal LM model), we would expect the log-Euler scheme (14.58) to produce lower biases than the Euler scheme (14.57). As discussed in Chapter 3, the log-Euler scheme will keep forward rates positive, whereas the Euler scheme will not.

As shown, both schemes (14.57), (14.58) advance time only by a single time step, but creation of a full path of forward curve evolution through time is merely a matter of repeated application of the single-period stepping schemes on a (possibly non-equidistant) time line t_0, t_1, \dots . When working in the spot measure, it is preferable to have the tenor structure dates T_1, T_2, \dots, T_{N-1} among the simulation dates, in order to keep track of the spot numeraire $B(\cdot)$ without having to resort to extrapolations. In fact, it is common in practice to set $t_i = T_i$, which, unless accrual periods τ_i are unusually long or volatilities unusually high, will normally produce an acceptable discretization error for many types of LM models. See e.g. Andersen and Andreasen [2000b] and Glasserman and Zhao [2000] for some numerical investigations of the Euler bias.

Remark 14.6.1. When t coincides with a date in the tenor structure, $t = T_k$, say, $q(t)$ will equal T_{k+1} due to our definition of q being right-continuous. As a result, when stepping forward from time $t = T_k$, $\hat{L}_k(T_k)$ will *not* be included in the computation of the drifts μ_n , $n \geq k + 1$. As it turns out, this convention reduces discretization bias, a result that makes sense when we consider that the contribution from $L_k(t)$ to the drifts drops to zero at time $T_k + dt$ in a continuous-time setting.

While Euler-type schemes such as (14.57) and (14.58) are not very sophisticated and, as we recall from Chapter 3, result in rather slow convergence of the discretization bias ($O(\Delta)$), these schemes are appealing in their straightforwardness and universal applicability. Further, they serve to highlight the basic structure of an LM simulation and the computational effort in advancing the forward curve.

14.6.1.1 Analysis of Computational Effort

Focusing on the straight Euler scheme (14.57), a bit of contemplation reveals that the computational effort involved in advancing L_n is dominated by the computation of $\mu_n(\cdot)$ which, in a direct implementation of (14.56), involves

$$m \cdot (n - q(t) + 1) = O(mn)$$

²³In addition to these time-stepping schemes for the forward rates, it may be necessary to simultaneously evolve stochastic volatility variables if one works with models such as those in Section 14.2.5.

operations for a given value of n . To advance all $N - q(t + \Delta)$ forward rates, it follows that the computational effort is $O(mN^2)$ for a single time step. Assuming that our simulation time line coincides with the tenor structure dates, generation of a full path of forward curve scenarios from time 0 to time T_{N-1} will thus require a total computational effort of $O(mN^3)$. As N is often big (e.g., a 25 year curve of quarterly forward rates will have $N = 100$), a naive application of the Euler scheme will often require considerable computing resources.

As should be rather obvious, however, the computational order of $O(mN^3)$ is easy to improve on, as there is no need to spend $O(mN)$ operations on the computation of each μ_n . Instead, we can invoke the recursive relationship

$$\mu_n(t) = \mu_{n-1}(t) + \frac{\tau_n \sigma_n(t)}{1 + \tau_n \hat{L}_n(t)}, \quad (14.59)$$

which allows us to compute all μ_n , $n = q(t + \Delta), \dots, N - 1$, by an $O(mN)$ -step iteration starting from

$$\mu_{q(t+\Delta)}(t) = \sum_{j=q(t)}^{q(t+\Delta)} \frac{\tau_j \sigma_j(t)}{1 + \tau_j \hat{L}_j(t)}.$$

In total, the computational effort of advancing the full curve one time step will be $O(mN)$, and the cost of taking N such time steps will be $O(mN^2)$ — and not $O(mN^3)$.

We summarize this result in a lemma.

Lemma 14.6.2. *Assume that we wish to simulate the entire Libor forward curve on a time line that contains the dates in the tenor structure and has $O(N)$ points. The computational effort of Euler-type schemes — such as (14.57) and (14.58) — is $O(mN^2)$.*

Remark 14.6.3. The results of the lemma can be verified to hold for any of the probability measures we examined in Section 14.2.2.

We note that when simulating in other measures, the starting point of the iteration for μ_n will be measure-dependent. For instance, in the terminal measure,

$$\mu_n(t) = - \sum_{j=n+1}^{N-1} \frac{\tau_j \sigma_j(t)}{1 + \tau_j \hat{L}_j(t)}$$

and the equation (14.59) still holds. Now, however, the iteration starts at

$$\mu_{N-1}(t) = 0,$$

and proceeds *backwards* through $\mu_{N-2}, \mu_{N-3}, \dots, \mu_{q(t+\Delta)}$. We leave it to the reader to carry out the analysis for other probability measures.

14.6.1.2 Long Time Steps

Most exotic interest rate derivatives involve revolving cash flows paid on a tightly spaced schedule (e.g. quarterly). As our simulation time line should always include dates on which cash flows take place, the average time spacing used in path generation will thus normally, by necessity, be quite small. In certain cases, however, there may be large gaps between cash flow dates, e.g. when a security is forward-starting or has an initial lock-out period. When simulating across large gaps, we may always choose to sub-divide the gap into smaller time steps, thereby retaining a tightly spaced simulation time line. To save computational time, however, it is often tempting to cover large gaps in a small number of coarse time steps, in order to lower overall computation effort. Whether such coarse stepping is possible is, in large part, a question of how well we can keep the discretization bias under control as we increase the time step, something that is quite dependent on the magnitude of volatility and the particular formulation of the LM model under consideration. Section 14.6.2 below deals with this question and offers strategies to improve on the basic Euler scheme. Here, we instead consider the pure mechanics of taking large time steps, i.e. steps that skip past several dates in the tenor structure.

Assume that we stand at the j -th date in the tenor structure, $t = T_j$, and wish to simulate the forward curve to time T_k , $k > j + 1$, in a single step. As noted earlier, the mere notion of skipping over dates in the tenor structure makes usage of the spot measure Q^B inconvenient, as the numeraire $B(T_k)$ cannot be constructed without knowledge of the realizations of $L_{j+1}(T_{j+1}), L_{j+2}(T_{j+2}), \dots, L_{k-1}(T_{k-1})$; in turn, numeraire-deflation of cash flows is not possible and derivatives cannot be priced. Circumventing this issue, however, is merely a matter of changing the numeraire from $B(t)$ to the price of an asset that involves no roll-over in the interval $[T_j, T_k]$. One such asset price is $P(t, T_N)$, the choice of which corresponds to running our simulated paths in the terminal measure. In particular, we recognize that

$$P(T_k, T_N) = \prod_{n=k}^{N-1} \frac{1}{1 + \tau_n L_n(T_k)}, \quad (14.60)$$

which depends only on the state of the forward curve at time T_k . Another valid numeraire asset would be $\bar{P}_j(t)$, as defined in Section 14.2.2:

$$\bar{P}_j(t) = \begin{cases} B(t), & t \leq T_j, \\ B(T_j)P(t, T_N)/P(T_j, T_N), & t > T_j. \end{cases}$$

The numeraire $\bar{P}_j(T_k)$ can always be computed without knowledge of $L_{j+1}(T_{j+1}), \dots, L_{k-1}(T_{k-1})$, as long as $B(T_j)$ is known²⁴. In the measure induced by this asset, the LM model dynamics are

²⁴This precludes the existence of other large gaps in the simulation time line prior to time T_j . When using a hybrid measure such as \bar{P}_j , we would need to

$$dL_n(t) = \begin{cases} \sigma_n(t)^\top \left(-\sum_{l=n+1}^{N-1} \frac{\tau_l \sigma_l(t)}{1+\tau_l L_l(t)} dt + d\bar{W}^j(t) \right), & t > T_j, \\ \sigma_n(t)^\top \left(\sum_{l=q(t)}^n \frac{\tau_l \sigma_l(t)}{1+\tau_l L_l(t)} dt + d\bar{W}^j(t) \right), & t \leq T_j. \end{cases}$$

14.6.1.3 Notes on the Choice of Numeraire

Given our discussion above, the terminal measure may strike the reader as an obvious first choice for simulating the LM model — after all, simulations in the terminal measure will never fail to be meaningful, irrespective of the coarseness of the simulation time line. Other issues, however, come in play here as well. For instance, updating the numeraire $P(t, T_N)$ from one time step to the next is generally a more elaborate operation than updating the spot numeraire $B(t)$: the former requires multiplying together $O(N)$ terms (see (14.60)), whereas the latter only involves multiplying $B(t)$ at the previous time step with a single discount bond price. Also, the statistical sample properties of price estimators in the terminal measure may be inferior to those in the spot measure, in the sense that the Monte Carlo noise is larger in the terminal measure. Glasserman and Zhao [2000] list empirical results indicating that this is, indeed, often the case for many common interest rate derivatives. A formal analysis of this observation is complex, but we can justify it by considering the pricing of a very simple derivative security, namely a discount bond maturing at some arbitrary time T_k in the tenor structure. In the spot measure, we would estimate the price of this security by forming the sample average of random variables

$$P(T_k, T_k)/B(T_k) = B(T_k)^{-1} = \frac{1}{\prod_{n=0}^{k-1} (1 + \tau_n L_n(T_n))}, \quad (14.61)$$

whereas in the terminal measure we would form the sample average of random variables

$$P(T_k, T_k)/P(T_k, T_N) = P(T_k, T_N)^{-1} = \prod_{n=k}^{N-1} (1 + \tau_n L_n(T_k)). \quad (14.62)$$

Assuming that Libor rates stay positive, the important thing to notice is that the right-hand side of (14.61) is bounded from above by 1, whereas the right-hand side of (14.62) can grow arbitrarily large. For moderate to high Libor rate volatilities, we would thus intuitively expect price estimators based on (14.62) to have higher sample error.

As discussed in Section 14.6.1.2, sometimes it is mechanically inconvenient to simulate in the spot measure, due to a desire to take large time steps. In these cases, usage of a hybrid numeraire $\bar{P}(t)$ that switches from $B(t)$ to $P(t, T_N)$ at the latest possible date may be a useful strategy.

position T_j at the start of the first simulation time step that spans multiple dates in the tenor structure.

14.6.2 Other Simulation Schemes

When simulating on a reasonably tight time schedule, the accuracy of the Euler or log-Euler schemes is adequate for most applications. However, as discussed above, we may occasionally be interested in using coarse time steps in some parts of the path generation algorithm, requiring us to pay more attention to the discretization scheme. Generic techniques for these purposes were introduced in detail in Chapter 3; we proceed to discuss a few of these in the context of LM models. We also consider the case where special-purpose schemes happen to exist for the discretization of the stochastic integral in the forward rate dynamics.

14.6.2.1 Special-Purpose Schemes with Drift Predictor-Corrector

In integrated form, the general LM dynamics in (14.56) become

$$\begin{aligned} L_n(t + \Delta) &= L_n(t) + \int_t^{t+\Delta} \sigma_n(u)^\top \mu_n(u) du + \int_t^{t+\Delta} \sigma_n(u)^\top dW^B(u) \\ &\triangleq L_n(t) + D_n(t, t + \Delta) + M_n(t, t + \Delta), \end{aligned} \quad (14.63)$$

where $M_n(t, t + \Delta)$ is a zero-mean martingale increment and $D_n(t, t + \Delta)$ is the increment of a predictable process. In many cases of practical interest, high-performance special-purpose schemes exist for simulation of $M_n(t, t + \Delta)$. This, for instance, is the case for the SV-LM model specification (Section 14.2.5), as discussed in detail in Section 9.5. In such cases, we obviously will choose to generate $M_n(t, t + \Delta)$ from the special-purpose scheme, and it thus suffices to focus on the term $D_n(t, t + \Delta)$. A simple approach is to use Euler stepping:

$$\widehat{L}_n(t + \Delta) = \widehat{L}_n(t) + \sigma_n(t)^\top \mu_n(t) \Delta + \widehat{M}_n(t, t + \Delta), \quad (14.64)$$

where $\widehat{M}_n(t, t + \Delta)$ is generated by a special-purpose scheme.

The drift adjustments in (14.64) are explicit in nature, as they are based only on the forward curve at time t . To incorporate information from time $t + \Delta$, we can use the *predictor-corrector* scheme from Section 3.2.5, which for (14.64) will take the two-step form

$$\overline{L}_n(t + \Delta) = \widehat{L}_n(t) + \sigma_n(t, \widehat{\mathbf{L}}(t))^\top \mu_n(t, \widehat{\mathbf{L}}(t)) \Delta + \widehat{M}_n(t, t + \Delta), \quad (14.65)$$

$$\begin{aligned} \widehat{L}_n(t + \Delta) &= \widehat{L}_n(t) + \theta_{PC} \sigma_n(t, \widehat{\mathbf{L}}(t))^\top \mu_n(t, \widehat{\mathbf{L}}(t)) \Delta \\ &\quad + (1 - \theta_{PC}) \sigma_n(t + \Delta, \overline{\mathbf{L}}(t + \Delta))^\top \mu_n(t + \Delta, \overline{\mathbf{L}}(t + \Delta)) \Delta \\ &\quad + \widehat{M}_n(t, t + \Delta), \end{aligned} \quad (14.66)$$

where θ_{PC} is a parameter in $[0, 1]$ that determines the amount of implicitness we want in our scheme ($\theta_{\text{PC}} = 1$: fully explicit; $\theta_{\text{PC}} = 0$: fully implicit). In practice, we would nearly always go for the balanced choice of $\theta_{\text{PC}} = 1/2$. In (14.65)–(14.66), \mathbf{L} denotes the vector of all Libor rates, $\mathbf{L}(t) = (L_1(t), \dots, L_{N-1}(t))^{\top}$ (with the convention that $L_i(t) \equiv L_i(T_i)$ for $i < q(t)$), and $\widehat{\mathbf{L}}, \overline{\mathbf{L}}$ defined accordingly. In particular, the short-hand notation $\mu_n(t, \widehat{\mathbf{L}}(t))$ is used to indicate that μ_n (and σ_n) may depend on the state of the entire forward curve at time t .

The technique above is based on a standard (additive) Euler scheme. If one is more inclined to use a multiplicative scheme in the vein of (14.58), we may replace the explicit scheme (14.64) with

$$\widehat{L}_n(t + \Delta) = \widehat{L}_n(t) \exp \left\{ \frac{\sigma_n(t)^{\top}}{\widehat{L}_n(t)} \mu_n(t) \Delta \right\} \widehat{M}_n(t, t + \Delta), \quad (14.67)$$

where $\widehat{M}_n(t, t + \Delta)$ now has been redefined to be a unit-mean positive random variable, often a discretized multiplicative increment of an exponential martingale. The construction of a predictor-corrector extension of (14.67) follows closely the steps above, and is left for the reader.

While the weak convergence order of simulation schemes may not be affected by predictor-corrector schemes (Section 3.2.5), experiments show that (14.65)–(14.66) often will reduce the bias significantly relative to a fully explicit Euler scheme. Some results for the simple log-normal LM model can be found in Hunter et al. [2001] and Rebonato [2002]. As the computational effort of applying the predictor step is not insignificant, the speed-accuracy trade-off must be evaluated on a case-by-case basis. Section 14.6.2.3 below discusses a possible modification of the predictor-corrector scheme to improve efficiency.

14.6.2.2 Euler Scheme with Predictor-Corrector

In simulating the term $M_n(t, t + \Delta)$ in the predictor-corrector scheme above, we can always use the Euler scheme, i.e. in (14.64) we set

$$\widehat{M}_n(t, t + \Delta) = \sigma_n(t)^{\top} \sqrt{\Delta} Z,$$

where Z is an m -dimensional vector of standard Gaussian draws. As we recall from Chapter 3, however, it may also be useful to apply the predictor-corrector principle to the martingale part of the forward rate evolution itself, although this would involve the evaluation of derivatives of the LM volatility term with respect to the forward Libor rates; see Chapter 3 for details.

14.6.2.3 Lagging Predictor-Corrector Scheme

Drift calculations, as was pointed out earlier, are the most computationally expensive part of any Monte Carlo scheme for a Libor market model. The

predictor-corrector scheme of (14.65)–(14.66) requires *two* calculations of the drift and is thus considerably more expensive than the standard Euler scheme. We often prefer to use a “lagging” modified predictor-corrector scheme which, as it turns out, allows us to realize most of the benefits of the predictor-corrector scheme, while keeping computational costs comparable to the standard Euler scheme.

Recall the definition of the drift of the n -th Libor rate under the spot measure,

$$\mu_n(t) = \sum_{j=q(t)}^n \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)}.$$

Note that the drift depends on the values of the Libor rates of indices less than or equal to n . Let us split the contributions coming from Libor rates with an index strictly less than n , and the n -th Libor rate,

$$\mu_n(t) = \sum_{j=q(t)}^{n-1} \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)} + \frac{\tau_n \sigma_n(t)}{1 + \tau_n L_n(t)}.$$

Denoting $t' = t + \Delta$, we observe that if we simulate the Libor rates in the order of increasing index, then by the time we need to simulate $L_n(t')$, we have already simulated $L_j(t')$, $j = q(t), \dots, n-1$. Hence, it is natural to use the predictor-corrector technique for the part of the drift that depends on Libor rates maturing strictly before T_n , while treating the part of the drift depending on the n -th Libor rate explicitly. This idea leads to the following scheme (compare to (14.64) or (14.65)–(14.66) with $\theta_{PC} = 1/2$),

$$\begin{aligned} \widehat{L}_n(t') &= \widehat{L}_n(t) + \sigma_n(t)^\top \\ &\times \left(\frac{1}{2} \sum_{j=q(t)}^{n-1} \left(\frac{\tau_j \sigma_j(t)}{1 + \tau_j \widehat{L}_j(t)} + \frac{\tau_j \sigma_j(t')}{1 + \tau_j \widehat{L}_j(t')} \right) + \frac{\tau_n \sigma_n(t)}{1 + \tau_n \widehat{L}_n(t)} \right) \Delta + \widehat{M}_n(t, t'). \end{aligned} \quad (14.68)$$

Importantly, the drifts required for this scheme also satisfy a recursive relationship, allowing for an efficient update. Defining

$$\widehat{\alpha}_n(t') = \sum_{j=q(t)}^n \left(\frac{\tau_j \sigma_j(t)}{1 + \tau_j \widehat{L}_j(t)} + \frac{\tau_j \sigma_j(t')}{1 + \tau_j \widehat{L}_j(t')} \right),$$

we see that, clearly,

$$\widehat{\alpha}_n(t') = \widehat{\alpha}_{n-1}(t') + \frac{\tau_n \sigma_n(t)}{1 + \tau_n \widehat{L}_n(t)} + \frac{\tau_n \sigma_n(t')}{1 + \tau_n \widehat{L}_n(t')},$$

and (14.68) can be rewritten as

$$\widehat{L}_n(t') = \widehat{L}_n(t) + \sigma_n(t)^\top \left(\frac{1}{2} \widehat{\alpha}_{n-1}(t') + \frac{\tau_n \sigma_n(t)}{1 + \tau_n \widehat{L}_n(t)} \right) \Delta + \widehat{M}_n(t, t'). \quad (14.69)$$

The scheme above can easily be applied to other probability measures. In fact, since in the terminal measure the drift

$$\mu_n(t) = - \sum_{j=n+1}^{N-1} \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)},$$

does not depend on $L_n(t)$ in the first place, no “lag” is required in this measure. Indeed, we simply redefine

$$\widehat{\alpha}_n(t') = - \sum_{j=n+1}^{N-1} \left(\frac{\tau_j \sigma_j(t)}{1 + \tau_j \widehat{L}_j(t)} + \frac{\tau_j \sigma_j(t')}{1 + \tau_j \widehat{L}_j(t')} \right)$$

and, starting from $n = N - 1$ and working backwards, use the scheme

$$\widehat{L}_n(t') = \widehat{L}_n(t) + \sigma_n(t)^\top \frac{1}{2} \widehat{\alpha}_n(t') \Delta + \widehat{M}_n(t, t'). \quad (14.70)$$

Notice that $\widehat{\alpha}_n$ now satisfies the recursion

$$\widehat{\alpha}_{n-1}(t') = \widehat{\alpha}_n(t') - \frac{\tau_n \sigma_n(t)}{1 + \tau_n \widehat{L}_n(t)} - \frac{\tau_n \sigma_n(t')}{1 + \tau_n \widehat{L}_n(t')},$$

to be started at $\widehat{\alpha}_{N-1}(t) = 0$.

The modifications of (14.69) and (14.70) to accommodate log-Euler stepping are trivial and left to the reader to explore. The lagging predictor-corrector scheme in the spot Libor measure has, as far as we know, not appeared in the literature, and its theoretical properties are not well-known (although the terminal measure version was studied in Joshi and Stacey [2008]). Still, its practical performance is very good and we do not hesitate recommending it as the default choice for many applications.

14.6.2.4 Further Refinements of Drift Estimation

For large time steps, it may be useful to explicitly integrate the time-dependent parts of the drift, rather than rely on pure Euler-type approximations. Focusing on, say, (14.63), assume that we can write

$$\sigma_n(u)^\top \mu_n(u) \approx g(u, \mathbf{L}(t)), \quad u \geq t, \quad (14.71)$$

for a function g that depends on time as well as the state of the forward rates frozen at time t . Then,

$$D_n(t, t + \Delta) = \int_t^{t+\Delta} \sigma_n(u)^\top \mu_n(u) du \approx \int_t^{t+\Delta} g(u, \mathbf{L}(t)) du. \quad (14.72)$$

As g evolves deterministically for $u > t$, the integral on the right-hand side can be evaluated either analytically (if g is simple enough) or by numerical quadrature. If doing the integral numerically, a decision must be made on the spacing of the integration grid. For volatility functions that are piecewise flat on the tenor-structure — which is a common assumption in model calibration — it is natural to align the grid with dates in the tenor structure.

To give an example, consider the DVF LM model, where we get (in the terminal measure, for a change)

$$\begin{aligned}\sigma_n(u)^\top \mu_n(u) &= -\varphi(L_n(u)) \lambda_n(u)^\top \sum_{j=n+1}^{N-1} \frac{\tau_j \lambda_j(u) \varphi(L_j(u))}{1 + \tau_j L_j(u)} \\ &\approx -\varphi(L_n(t)) \lambda_n(u)^\top \sum_{j=n+1}^{N-1} \frac{\tau_j \lambda_j(u) \varphi(L_j(t))}{1 + \tau_j L_j(t)}, \quad u \geq t,\end{aligned}$$

which is of the form (14.71). For stochastic volatility models we might, say, additionally assume that the process $z(t)$ would stay on its expected path, i.e. $z(u) \approx E_t^N(z(u))$ which can often be computed in closed form for models of interest. For instance, for the SV model in (14.15) we have

$$E_t^N(z(u)) = z_0 + (z(t) - z_0)e^{-\theta(u-t)}.$$

The approach in (14.72) easily combines with predictor-corrector logic, i.e. we could write

$$\begin{aligned}D_n(t, t + \Delta) &\approx \theta_{PC} \int_t^{t+\Delta} g(u, \mathbf{L}(t)) du \\ &\quad + (1 - \theta_{PC}) \int_t^{t+\Delta} g(u, \bar{\mathbf{L}}(t + \Delta)) du, \quad (14.73)\end{aligned}$$

where $\bar{\mathbf{L}}_i(t + \Delta)$ has been found in a predictor step using (14.72) in (14.64). The “lagged” schemes in Section 14.6.2.3 work equally well. Formula (14.72) also applies to exponential-type schemes such as (14.67), with or without predictor-corrector adjustment; we leave details to the reader.

14.6.2.5 Brownian-Bridge Schemes and Other Ideas

As a variation on the predictor-corrector scheme, we could attempt a further refinement of taking into account variance of the Libor curve between the sampling dates t and $t + \Delta$. Schemes attempting to do so by application of *Brownian bridge techniques*²⁵ were proposed in Andersen [2000b] and Pietersz et al. [2004], among others. While performance of these schemes is

²⁵See Section 3.2.9 for an introduction to the Brownian bridge, albeit for a somewhat different application.

mixed — tests in Joshi and Stacey [2008] show rather unimpressive results in comparison to simpler predictor-corrector schemes — the basic idea is sufficiently simple and instructive to merit a brief mention. In a nutshell, the Brownian bridge approach aims to replace in (14.72) all forward rates $\mathbf{L}(t)$ with the expectation of $\mathbf{L}(u)$, *conditional* upon the forward rates ending up at $\bar{\mathbf{L}}(t + \Delta)$, where $\bar{\mathbf{L}}(t + \Delta)$ is generated in a predictor step. Under simplifying assumptions on the dynamics of $L_n(t)$, a closed-form expression is possible for this expectation.

Proposition 14.6.4. *Assume that*

$$dL_n(t) \approx \sigma_n(t)^\top dW(t),$$

where $\sigma_n(t)$ is deterministic and $W(t)$ is an m -dimensional Brownian motion in some probability measure P . Let

$$v_n(t, T) = \int_t^T \|\sigma_n(s)\|^2 ds, \quad T \geq t.$$

Then, for $u \in [t, t + \Delta]$,

$$\mathbb{E}(L_n(u)|L_n(t), L_n(t + \Delta)) = L_n(t) + \frac{v_n(t, u)}{v_n(t, t + \Delta)} (L_n(t + \Delta) - L_n(t)).$$

Proof. We first state a very useful general result for multi-variate Gaussian variables.

Lemma 14.6.5. *Let $X = (X_1, X_2)^\top$ be a partitioned vector of Gaussian variables, where X_1 and X_2 are themselves vectors. Assume that the covariance matrix between X_i and X_j is $\Sigma_{i,j}$ such that the total covariance matrix of X is*

$$\Sigma = \begin{pmatrix} \Sigma_{1,1} & \Sigma_{1,2} \\ \Sigma_{2,1} & \Sigma_{2,2} \end{pmatrix}$$

(where, of course, $\Sigma_{2,1} = \Sigma_{1,2}^\top$). Let the vector means of X_i be μ_i , $i = 1, 2$, and assume that $\Sigma_{2,2}$ is invertible. Then $X_1|X_2 = x$ is Gaussian:

$$(X_1|X_2 = x) \sim \mathcal{N}(\mu_1 + \Sigma_{1,2}\Sigma_{2,2}^{-1}(x - \mu_2), \Sigma_{1,1} - \Sigma_{1,2}\Sigma_{2,2}^{-1}\Sigma_{2,1}).$$

In Lemma 14.6.5, now set $X_1 = L_n(u) - L_n(t)$ and $X_2 = L_n(t + \Delta) - L_n(t)$. Note that $\mu_1 = \mu_2 = 0$ and

$$\Sigma_{1,2} = \Sigma_{2,1} = \Sigma_{1,1} = v_n(t, u), \quad \Sigma_{2,2} = v_n(t, t + \Delta).$$

The result of Proposition 14.6.4 follows. \square

We can use the result of Proposition 14.6.4 in place of the ordinary corrector step. For instance, in (14.73) we write

$$D_n(t, t + \Delta) \approx \int_t^{t+\Delta} g(u, \mathbf{m}(u)) du,$$

where, for $\mathbf{m}(u) = (m_1(u), \dots, m_{N-1}(u))$,

$$m_i(u) = \mathbb{E}^B(L_i(u)|L_i(t), \bar{L}_i(t + \Delta))$$

is computed according to Proposition 14.6.4 once $\bar{L}_i(t + \Delta)$ has been sampled in a predictor step.

In some cases, it may be more appropriate to assume that L_n is roughly log-normal, in which case Proposition 14.6.4 must be altered slightly.

Lemma 14.6.6. *Assume that*

$$dL_n(t)/L_n(t) \approx \sigma_n(t)^\top dW(t),$$

where $\sigma_n(t)$ is deterministic and $W(t)$ is an m -dimensional Brownian motion in some probability measure P . Then, for $u \in [t, t + \Delta]$,

$$\begin{aligned} \mathbb{E}(L_n(u)|L_n(t), L_n(t + \Delta)) &= L_n(t) \left(\frac{L_n(t + \Delta)}{L_n(t)} \right)^{v_n(t, u)/v_n(t, t + \Delta)} \\ &\times \exp \left(\frac{v_n(t, u)(v_n(t, t + \Delta) - v_n(t, u))}{2v_n(t, t + \Delta)} \right), \end{aligned}$$

where $v_n(t, T)$ is given in Proposition 14.6.4.

Proof. Apply Lemma 14.6.5 to $X_1 = \ln L_n(u) - \ln L_n(t)$ and $X_2 = \ln L_n(t + \Delta) - \ln L_n(t)$. To translate back to find the conditional mean of e^{X_1} , one may use the fact that $\mathbb{E}(e^{a+bY}) = e^{a+b^2/2}$ if Y is Gaussian $\mathcal{N}(0, 1)$. \square

Joshi and Stacey [2008] investigate a number of other possible discretization schemes for the drift term in the LM model, including ones that attempt to incorporate information about the correlation between various forward rates. In general, many of these schemes will result in some improvement of the discretization error, but at the cost of more computational complexity and effort. All things considered, we hesitate to recommend any of these methods (and this goes for the Brownian bridge scheme above) for general-purpose use, as the bias produced by simpler methods is often adequate. If not, it may, in fact, often be the case that we can insert a few extra simulation dates inside large gaps to bring down the bias, yet still spend less computational time than we would if using a more complex method of bridging the gap in a single step. Finally, we should note that most authors (including Joshi and Stacey [2008]) exclusively examine simple log-normal models where the martingale component (M_n in the notation of Section 14.6.2.1) can be simulated completely bias-free. When using more realistic models, this will not always be the case, in which case high-precision simulation of the drift term D_n will likely be a waste of time.

14.6.2.6 High-Order Schemes

Even with predictor-corrector adjustment, all Euler-type discretization schemes are limited to a convergence order of Δ . To raise this, one possibility is to consider higher-order schemes, such as the Milstein scheme and similar Taylor-based approaches; see Section 3.2.6 for details. Many high-order schemes unfortunately become quite cumbersome to deal with for the type of high-dimensional vector-SDE that arises in the context of LM models and, possibly as a consequence of this, there are currently very few empirical results in the literature to lean on. One exception is Brotherton-Ratcliffe (Brotherton-Ratcliffe [1997]) where a Milstein scheme was developed for the basic log-normal LM model with piecewise flat volatilities. The efficacy of this, and similar high-order schemes, in the context of the generalized LM model would obviously depend strongly on the particular choice of model formulation.

A simple alternative to classical Taylor-based high-order schemes involves Richardson extrapolation based on prices found by simulating on two separate time lines, one coarser than the other (see Section 3.2.7 for details). Andersen and Andreasen [2000b] list some results for this scheme, the efficacy of which seems to be rather modest.

14.6.3 Martingale Discretization

Consider again the hybrid measure induced by the numeraire \tilde{P}_{n+1} , defined in Section 14.2.2. As discussed, one effect of using this measure is to render the process for the n -th forward Libor rate $L_n(t)$ a martingale. When time-discretizing the LM model using, say, an Euler scheme, the martingale property of $L_n(t)$ is automatically preserved, ensuring that the expectation of the discretized approximation $\hat{L}_n(t)$ will have expectation $L_n(0)$, with no discretization bias. Also, when using Monte Carlo to estimate the price of the zero-coupon bond maturing at time T_{n+1} , we get

$$P(0, T_{n+1}) = \tilde{P}_{n+1}(0)E^{n+1}(1),$$

which will (obviously) be estimated bias-free as well.

As the discussion above highlights, it is possible to select a measure such that a particular zero-coupon bond and a particular FRA will be priced bias-free²⁶ by Monte Carlo simulation, even when using a simple Euler scheme. While we are obviously rarely interested in pricing zero-coupon bonds by Monte Carlo methods, this observation can nevertheless occasionally help guide the choice of simulation measure, particularly if, say, a security can be argued to depend primarily on a single forward rate (e.g. caplet-like securities). In practice, matters are rarely this clear-cut, and one wonders

²⁶But not error-free, of course — there will still be a statistical zero-mean error on the simulation results. See Section 14.6.4 below.

whether perhaps simulation schemes exist that will simultaneously price all zero-coupon bonds $P(t, T_1), P(t, T_2), \dots, P(t, T_N)$ bias-free. It should be obvious that this cannot be accomplished by a simple measure-shift, but will require a more fundamental change in simulation strategy.

14.6.3.1 Deflated Bond Price Discretization

Fundamentally, we are interested in a simulation scheme that by construction will ensure that all numeraire-deflated bond prices are martingales. The easiest way to accomplish this is to follow a suggestion offered by Glasserman and Zhao [2000]: instead of discretizing the dynamics for Libor rates directly, simply discretize the deflated bond prices themselves. To demonstrate, let us consider the spot measure, and define

$$U(t, T_{n+1}) = \frac{P(t, T_{n+1})}{B(t)}. \quad (14.74)$$

Lemma 14.6.7. *Let dynamics in the spot measure Q^B be as in Lemma 14.2.3. The dynamics for deflated zero-coupon bond prices (14.74) are given by*

$$\frac{dU(t, T_{n+1})}{U(t, T_{n+1})} = - \sum_{j=q(t)}^n \tau_j \frac{U(t, T_{j+1})}{U(t, T_j)} \sigma_j(t)^\top dW^B(t), \quad n = q(t), \dots, N-1. \quad (14.75)$$

Proof. We note that, by definition,

$$U(t, T_{n+1}) = \frac{P(t, T_{q(t)}) \prod_{j=q(t)}^n \frac{1}{1+\tau_j L_j(t)}}{P(t, T_{q(t)}) B(T_{q(t)-1})} = \frac{\prod_{j=q(t)}^n \frac{1}{1+\tau_j L_j(t)}}{B(T_{q(t)-1})},$$

where $B(T_{q(t)-1})$ is non-random at time t . We have that $U(t, T_{n+1})$ must, by construction, be a martingale in Q^B . An application of Ito's lemma to the diffusion term of U gives

$$dU(t, T_{n+1}) = -U(t, T_{n+1}) \sum_{j=q(t)}^n \frac{\tau_j \sigma_j(t)^\top}{1 + \tau_j L_j(t)} dW^B(t),$$

and the lemma follows once we note that

$$\frac{U(t, T_{j+1})}{U(t, T_j)} = \frac{1}{1 + \tau_j L_j(t)}.$$

□

Discretization schemes for (14.75) that preserve the martingale property are easy to construct. For instance, we could use the log-Euler scheme

$$\widehat{U}(t + \Delta, T_{n+1}) = \widehat{U}(t, T_{n+1}) \exp \left(-\frac{1}{2} \|\gamma_{n+1}(t)\|^2 \Delta + \gamma_{n+1}(t)^\top Z \sqrt{\Delta} \right). \quad (14.76)$$

where, as before, Z is an m -dimensional standard Gaussian random variable and

$$\gamma_{n+1}(t) \triangleq - \sum_{j=q(t)}^n \tau_j \frac{\widehat{U}(t, T_{j+1})}{\widehat{U}(t, T_j)} \sigma_j(t). \quad (14.77)$$

We have several remarks to the log-Euler scheme (14.76). First, for models where interest rates cannot become negative, $U(t, T_{n+1})/U(t, T_n) = P(t, T_{n+1})/P(t, T_n)$ cannot exceed 1 in a continuous-time model, so it might be advantageous to replace (14.77) with

$$\gamma_{n+1}(t) \triangleq - \sum_{j=q(t)}^n \tau_j \min \left(\frac{\widehat{U}(t, T_{j+1})}{\widehat{U}(t, T_j)}, 1 \right) \sigma_j(t),$$

as recommended in Glasserman and Zhao [2000]. Second, for computational efficiency we should rely on iterative updating,

$$\gamma_{n+1}(t) = \gamma_n(t) - \tau_n \min \left(\frac{\widehat{U}(t, T_{n+1})}{\widehat{U}(t, T_n)}, 1 \right) \sigma_n(t),$$

using the same arguments as those presented in Section 14.6.1.1. Third, once $\widehat{U}(t + \Delta, T_n)$ has been drawn for all possible n , we can reconstitute the Libor curve from the relation

$$\widehat{L}_n(t + \Delta) = \frac{\widehat{U}(t + \Delta, T_n) - \widehat{U}(t + \Delta, T_{n+1})}{\tau_n \widehat{U}(t + \Delta, T_{n+1})}, \quad n = q(t + \Delta), \dots, N - 1. \quad (14.78)$$

For completeness, we note that dynamics of the deflated bond prices in the terminal measure Q^{T_N} can easily be derived to be

$$\frac{dU(t, T_{n+1})}{U(t, T_{n+1})} = \sum_{j=n+1}^{N-1} \tau_j \frac{U(t, T_{j+1})}{U(t, T_j)} \sigma_j(t)^\top dW^N(t), \quad (14.79)$$

where we must now (re-)define $U(t, T_n)$ as

$$U(t, T_n) = P(t, T_n)/P(t, T_N).$$

Equation (14.79) can form the basis of a discretization scheme in much the same manner as above.

14.6.3.2 Comments and Alternatives

The discretization scheme presented above will preserve the martingale property of all deflated bonds maturing in the tenor structure, and in this

sense can be considered arbitrage-free. The resulting lack of bias on bond prices, however, does not necessarily translate into a lack of bias on any other derivative security price, e.g. a caplet or a swaption. In particular, we notice that nothing in the scheme above will ensure that bond price moments of any order other than one will be simulated accurately.

The extent of the bias induced by the scheme in Section 14.6.3.1 is specific to the security and model under consideration. For instance, using a log-Euler scheme for deflated bonds might work well in an LM model with rates that are approximately Gaussian, but might work less well in a model where rates are approximately log-normal. If results are disappointing, we can replace (14.76) with another discretization of (14.75) (see Chapter 3 for many examples), or we can try to discretize a quantity other than the deflated bonds $U(t, T_n)$. The latter idea is pursued in Glasserman and Zhao [2000], where several suggestions for discretization variables are considered. For instance, one can consider the differences

$$U(t, T_n) - U(t, T_{n+1}) \quad (14.80)$$

which are martingales since the U 's are. As follows from (14.78), discretizing $U(t, T_n) - U(t, T_{n+1})$ is, in a sense, close to discretizing $L_n(t)$ itself which may be advantageous. Joshi and Stacey [2008] contains some tests of discretization schemes based on (14.80), but, again, only in a log-normal setting. Additional tests in Beveridge et al. [2008] (in a displaced log-normal case) find (14.80) inferior to the standard predictor-corrector scheme and demonstrate that the scheme can lead to negative (path realizations of) bond prices.

14.6.4 Variance Reduction

We recall from the discussion in Chapter 3 that the errors involved in Monte Carlo pricing of derivatives can be split into two sources: the statistical Monte Carlo error (the standard error), and a bias unique to the discretization scheme employed. So far, our discussion has centered exclusively on the latter of these two types of errors and we now wish to provide some observations about the former. We should note, however, that it is difficult to provide generic prescription for variance reduction techniques in the LM model, as most truly efficient schemes tend to be quite specific to the product being priced. We shall offer several such product-specific variance reduction schemes in later chapters, and here limit ourselves to rather brief suggestions.

We recall that Chapter 3 discussed three types of variance reduction techniques: i) antithetic sampling; ii) control variates; and iii) importance sampling. All have potential uses in simulation of LM models.

14.6.4.1 Antithetic Sampling

Application of antithetic sampling to LM modeling is straightforward. Using the Euler scheme as an example, each forward rate sample path generated

from the relation

$$\widehat{L}_n(t + \Delta) = \widehat{L}_n(t) + \sigma_n(t)^\top (\mu_n(t)\Delta + \sqrt{\Delta}Z)$$

is simply accompanied by a “reflected” sample path computed by flipping the vector-valued Gaussian variable Z around the origin, i.e.

$$\widehat{L}_n^{(a)}(t + \Delta) = \widehat{L}_n^{(a)}(t) + \sigma_n(t)^\top (\mu_n^{(a)}(t)\Delta - \sqrt{\Delta}Z).$$

The reflection of Z is performed at each time step, with both paths having identical starting points, $\widehat{L}_n^{(a)}(0) = \widehat{L}_n(0) = L_n(0)$. Using antithetic variates thus doubles the number of sample paths that will be generated from a fixed budget of random number draws. In practice, the variance reduction associated with antithetic variates is often relatively modest.

14.6.4.2 Control Variates

As discussed in Chapter 3, the basic (product-based) control variate method involves determining a set of securities (control variates) that i) have payouts close to that of the instrument we are trying to price; and ii) have known expected values in the probability measure in which we simulate. Obvious control variates in the LM model include (portfolios of) zero-coupon bonds and caplets. Due to discretization errors in generation of sample paths, we should note, however, that the sample means of zero-coupon bonds and caplets will deviate from their true continuous-time means with amounts that depend on the time step and the discretization scheme employed. This error will nominally cause a violation of condition ii) — we are generally able only to compute in closed-form the continuous-time expected values — but the effect is often benign and will theoretically be of the same order²⁷ as the weak convergence order of the discretization scheme employed. Swaptions can also be included in the control variate set, although additional care must be taken here due to the presence of hard-to-quantify approximation errors in the formulas in Section 14.4.2. See Jensen and Svenstrup [2003] for an example of using swaptions as control variates for Bermudan swaptions.

An alternative interpretation of the control variate idea involves pricing a particular instrument using, in effect, two different LM models, one of which allows for an efficient computation of the instrument price, and one of which is the true model we are interested in applying. We shall return to this in Chapter 25.

Finally, the dynamic control variate method, based on the idea that an (approximate) self-financed hedging strategy could be a good proxy for the value of a security, is available for LM models as well. The method was developed in Section 3.4.3.2.

²⁷ Suppose that we estimate $E(X) \approx E(X' + Y' - \mu_Y)$, where $\mu_Y = E(Y') + O(\Delta^p)$ and $E(X') = E(X) + O(\Delta^p)$. Then clearly also $E(X' + Y' - \mu_Y) = E(X) + O(\Delta^p)$.

14.6.4.3 Importance Sampling

Importance sampling techniques have so far found relatively limited use in the simulation of LM models, although Capriotti [2007] demonstrates that least-squares importance sampling (see Section 3.4.4.4) gives good results when pricing simple European options (caps and swaptions) in a three-factor log-normal LM model. As the variance reduction efficiency of importance sampling depends strongly on the payout, it is, however, unclear to what extent the results in Capriotti [2007] carry over to more complex security payouts (and models, for that matter).

Probably the most fruitful application of importance sampling in LM modeling is in the pricing of securities with a knock-out barrier. The basic idea is here that sample paths are generated conditional on a barrier not being breached, ensuring that all paths survive to maturity; this conditioning step induces a change of measure. We will expose the details of this technique in Chapter 20, where we discuss the pricing of the TARN product introduced in Section 5.15.2.

The Libor Market Model II

For the sake of cohesion, our discussion of LM models in Chapter 14 silently skipped over a number of practical issues. Chief amongst these is the problem of how to construct an entire (continuous) discount curve from knowledge of only a finite set of simply compounded Libor rates. This surprisingly subtle issue shall be discussed in this chapter, along with a select set of other advanced pricing and calibration topics in LM modeling. For instance, we provide a number of extensions to the stochastic volatility setup of Chapter 14 and also show how to construct swaption pricing formulas more accurate and more general than those in Chapter 14. We cover the generic problem of evolving separate discount and forward curves, and also include brief discussions of the so-called *swap market models* and of LM models with “near-Markov” structure. The latter topic shall be taken up again in Chapter 25.

15.1 Interpolation

The simulation schemes that we developed so far (see Section 14.6) allow us to obtain at any time t a vector of forward Libor rates on a pre-specified tenor structure. As should be obvious, and as pointed out previously in Sections 14.1.2 and 14.2.3, this information is not sufficient to recover the full interest rate yield curve at time t . At the very least, to be able to compute $P(t, T_n)$ for all n (see (14.3)), we need to additionally establish the *front stub* discount factor $P(t, T_{q(t)})$. In addition, as many actual security payoffs dictate that we calculate $P(t, T)$ for an arbitrary T , the *back stub* (forward) discount factor $P(t, T, T_{q(T)}) = P(t, T_{q(T)})/P(t, T)$ will also be required, since

$$P(t, T) = P(t, T_{q(t)}) \times \left(\prod_{i=q(t)}^{q(T)-1} (1 + \tau_i L_i(t))^{-1} \right) / P(t, T, T_{q(T)}).$$

Both the front and back stubs cannot, in general, be expressed as a function of Libor rates on a fixed tenor structure.

There are a number of approaches that could be employed to obtain the front and back stub in a simulation. We start with the back stub as it is somewhat easier to handle.

15.1.1 Back Stub, Simple Interpolation

Let us fix the discount factor maturity time T , and set $m = q(T)$ such that

$$T_{m-1} \leq T < T_m.$$

Observe that as $T \rightarrow T_{m-1}$ the back stub $P(t, T, T_m) = P(t, T_m)/P(t, T)$ converges to $P(t, T_{m-1}, T_m)$, a discount factor that can be calculated from Libor rates as $P(t, T_{m-1}, T_m) = (1 + L_{m-1}(t)\tau_{m-1})^{-1}$. At the other limit, when $T \rightarrow T_m$, the back stub converges to 1. Hence, it seems reasonable to approximate $P(t, T, T_m)$ by interpolating between these two known extremes. This idea gives rise to a number of plausible schemes that we now proceed to describe.

A particularly simple idea is to apply linear interpolation directly to bond prices, resulting in the scheme

$$P(t, T, T_m) = \frac{T - T_{m-1}}{T_m - T_{m-1}} + \frac{T_m - T}{T_m - T_{m-1}} P(t, T_{m-1}, T_m). \quad (15.1)$$

Using $P(t, T_{m-1}, T) = P(t, T_{m-1}, T_m)/P(t, T, T_m)$ as the interpolation variable instead, another linear interpolation scheme arises:

$$P(t, T_{m-1}, T) = \frac{T_m - T}{T_m - T_{m-1}} + \frac{T - T_{m-1}}{T_m - T_{m-1}} P(t, T_{m-1}, T_m),$$

or

$$P(t, T, T_m) = \frac{P(t, T_{m-1}, T_m)}{\frac{T_m - T}{T_m - T_{m-1}} + \frac{T - T_{m-1}}{T_m - T_{m-1}} P(t, T_{m-1}, T_m)}. \quad (15.2)$$

Alternatively, we can apply piecewise constant interpolation to continuously compounded instantaneous forward rates $f(t, u)$ for $u \in [T_{m-1}, T_m]$, yielding

$$-\frac{1}{T_m - T} \ln P(t, T, T_m) = -\frac{1}{T_m - T_{m-1}} \ln P(t, T_{m-1}, T_m),$$

or, in terms of forward bond prices,

$$P(t, T, T_m) = P(t, T_{m-1}, T_m)^{\frac{T_m - T}{T_m - T_{m-1}}}. \quad (15.3)$$

Yet another interpolation scheme is obtained by constant interpolation of simply compounded rates,

$$\frac{1}{T_m - T} \left(\frac{1}{P(t, T, T_m)} - 1 \right) = \frac{1}{T_m - T_{m-1}} \left(\frac{1}{P(t, T_{m-1}, T_m)} - 1 \right),$$

resulting in the scheme

$$P(t, T, T_m) = \left(\frac{T - T_{m-1}}{T_m - T_{m-1}} + \frac{T_m - T}{T_m - T_{m-1}} \frac{1}{P(t, T_{m-1}, T_m)} \right)^{-1}. \quad (15.4)$$

While the interpolation schemes (15.1), (15.2), (15.3), and (15.4) are all simple to understand and to apply, they are ultimately flawed as they violate the basic no-arbitrage conditions

$$P(0, T_{m-1}, T) = E^{T_{m-1}}(P(t, T_{m-1}, T)), \quad (15.5)$$

$$P(0, T_{m-1}, T_m) = E^{T_{m-1}}(P(t, T_{m-1}, T_m)). \quad (15.6)$$

Any interpolation scheme, once we apply the expected value operator, imposes a certain relationship on discount bonds at time 0, a relationship that, in general, will not be satisfied by actual market prices. Taking as an example (15.2), applying the expectation operator $E^{T_{m-1}}$ and using (15.5), we obtain

$$P(0, T_{m-1}, T) = \frac{T_m - T}{T_m - T_{m-1}} + \frac{T - T_{m-1}}{T_m - T_{m-1}} P(0, T_{m-1}, T_m).$$

a relationship between time 0 discount bond prices that is unlikely to be satisfied *a priori*.

15.1.2 Back Stub, Arbitrage-Free Interpolation

To ensure that observable relationships between time 0 discount bond prices are respected, consider choosing an arbitrary constant $\alpha(T)$ and setting

$$P(t, T_{m-1}, T) = P(0, T_{m-1}, T) + \alpha(T) (P(t, T_{m-1}, T_m) - P(0, T_{m-1}, T_m)). \quad (15.7)$$

Clearly, the additive scheme (15.7) will preserve (15.5) as long as (15.6) is satisfied — and (15.6) is essentially a no-arbitrage condition for discount bonds maturing on the tenor structure and is guaranteed by the LM model construction itself.

We can regard $\alpha(\cdot)$ as a function of maturity time¹; for consistency $\alpha(T_m)$ must equal 1 and $\alpha(T_{m-1}) = 0$, but beyond this there are few restrictions on $\alpha(T)$. Yet it is not advisable to specify $\alpha(T)$ arbitrarily, as this may affect model dynamics in unintended ways. To devise a reasonable approach to the definition of $\alpha(T)$, we note from (15.7) that

$$dP(t, T_{m-1}, T) = O(dt) + \alpha(T) dP(t, T_{m-1}, T_m).$$

¹And, implicitly, calendar time, which we ignore here as we work with a fixed t .

On the other hand, in an HJM model we have

$$\begin{aligned} dP(t, T_{m-1}, T) / P(t, T_{m-1}, T) \\ = O(dt) + (\sigma_P(t, T_{m-1}) - \sigma_P(t, T))^{\top} dW(t), \end{aligned}$$

and

$$\begin{aligned} dP(t, T_{m-1}, T_m) / P(t, T_{m-1}, T_m) \\ = O(dt) + (\sigma_P(t, T_{m-1}) - \sigma_P(t, T_m))^{\top} dW(t), \end{aligned}$$

from which we conclude that $\alpha(T)$ is linked to the ratio of forward discount bond volatilities. Exploiting this link, we may define $\alpha(T)$ from, for instance, the equation

$$\begin{aligned} P(t, T_{m-1}, T) \|\sigma_P(t, T_{m-1}) - \sigma_P(t, T)\| \\ = \alpha(T) P(t, T_{m-1}, T_m) \|\sigma_P(t, T_{m-1}) - \sigma_P(t, T_m)\|. \end{aligned}$$

Then²

$$\begin{aligned} \alpha(T) &= \frac{P(t, T_{m-1}, T)}{P(t, T_{m-1}, T_m)} \frac{\|\sigma_P(t, T_{m-1}) - \sigma_P(t, T)\|}{\|\sigma_P(t, T_{m-1}) - \sigma_P(t, T_m)\|} \\ &\approx \frac{P(0, T_{m-1}, T)}{P(0, T_{m-1}, T_m)} \frac{\|\sigma_P(t, T_{m-1}) - \sigma_P(t, T)\|}{\|\sigma_P(t, T_{m-1}) - \sigma_P(t, T_m)\|}. \end{aligned} \quad (15.8)$$

As we have seen in Chapter 14, the LM model does not uniquely define all the bond volatilities in (15.8) and we would need to interpolate available volatilities to compute (15.8). Note that (15.8) turns the problem of interpolating bond prices into a problem of interpolating bond *volatilities* instead. This point of view is advantageous as there are few, if any, arbitrage restrictions on interpolating the volatilities of bonds, as opposed to the bonds themselves. One can, for example, choose a linear interpolation to obtain $\sigma_P(t, T)$ from $\sigma_P(t, T_{m-1})$ and $\sigma_P(t, T_m)$ or, for a more sophisticated scheme, draw inspiration from the shape of forward volatilities in a mean-reverting one-factor Gaussian model. To explore the latter idea, recall that in the one-dimensional Gaussian model with constant mean reversion (Section 10.1.2),

$$\sigma_P(t, T) = \sigma(t) \frac{1 - e^{-\kappa(T-t)}}{\kappa}, \quad (15.9)$$

where κ is the mean reversion and $\sigma(t)$ is the short rate volatility. Plugging into (15.8) and rearranging we get

²As discount bond volatility vectors $\sigma_P(t, T)$ are generally non-deterministic, a suitable approximation is required to make sense of this formula. This is most easily done by freezing any state variables appearing in bond volatilities, such as Libor rates, at their time 0 values.

$$\alpha(T) \approx \frac{P(0, T_{m-1}, T)}{P(0, T_{m-1}, T_m)} \frac{1 - e^{-\kappa(T-T_{m-1})}}{1 - e^{-\kappa(T_m-T_{m-1})}}.$$

The mean reversion parameter κ could be either set as part of the user input, or obtained by best-fitting the Gaussian parametric form (15.9) to the volatility structure of the LM model.

15.1.3 Back Stub, Interpolation Inspired by the Gaussian Model

Above, we used a volatility parameterization inspired by the Gaussian model to construct an interpolation scheme for bond prices. An alternative, more direct, approach to extracting information from a Gaussian model in the interpolation exercise is to let the bond reconstitution formulas from the Gaussian model form the basis for interpolation. To demonstrate, recall that in the one-factor Gaussian model with constant mean reversion κ and short rate volatility $\sigma(t)$ (see Section 10.1.2.2),

$$\begin{aligned} P(t, T_{m-1}, T) &= P(0, T_{m-1}, T) \exp \left(-G(T_{m-1}, T) e^{-\kappa(T_{m-1}-t)} x(t) \right) \\ &\quad \times \exp \left(-\frac{1}{2} \left(G(t, T)^2 - G(t, T_{m-1})^2 \right) y(t) \right), \end{aligned} \quad (15.10)$$

and

$$\begin{aligned} P(t, T_{m-1}, T_m) &= P(0, T_{m-1}, T_m) \exp \left(-G(T_{m-1}, T_m) e^{-\kappa(T_{m-1}-t)} x(t) \right) \\ &\quad \times \exp \left(-\frac{1}{2} \left(G(t, T_m)^2 - G(t, T_{m-1})^2 \right) y(t) \right), \end{aligned} \quad (15.11)$$

where

$$G(t, T) = \frac{1 - e^{-\kappa(T-t)}}{\kappa}, \quad (15.12)$$

$$y(t) = e^{-2\kappa t} \int_0^t e^{2\kappa s} \sigma(s)^2 ds, \quad (15.13)$$

and $x(t)$ is the short rate state. Eliminating $x(t)$ in (15.10) and (15.11) defines a relationship between bond prices,

$$\begin{aligned} &\frac{1}{G(T_{m-1}, T)} \left(\ln \frac{P(t, T_{m-1}, T)}{P(0, T_{m-1}, T)} + (G(t, T)^2 - G(t, T_{m-1})^2) \frac{y(t)}{2} \right) \\ &= \frac{1}{G(T_{m-1}, T_m)} \left(\ln \frac{P(t, T_{m-1}, T_m)}{P(0, T_{m-1}, T_m)} + (G(t, T_m)^2 - G(t, T_{m-1})^2) \frac{y(t)}{2} \right) \end{aligned}$$

or, after a few additional manipulations,

$$\begin{aligned}
P(t, T_{m-1}, T) &= P(0, T_{m-1}, T) \left(\frac{P(t, T_{m-1}, T_m)}{P(0, T_{m-1}, T_m)} \right)^{\frac{G(T_{m-1}, T)}{G(T_{m-1}, T_m)}} \\
&\quad \times \exp \left(\frac{1}{2} \frac{G(T_{m-1}, T)}{G(T_{m-1}, T_m)} \left(G(t, T_m)^2 - G(t, T_{m-1})^2 \right) y(t) \right) \\
&\quad \times \exp \left(-\frac{1}{2} \left(G(t, T)^2 - G(t, T_{m-1})^2 \right) y(t) \right). \tag{15.14}
\end{aligned}$$

In (15.14), the volatility $\sigma(t)$ and the mean reversion κ can be obtained, for example, by fitting the Gaussian volatility structure to the volatilities of Libor rates generated by the LM model itself. High level of precision is not required here; we can take, say, $\kappa = 0$ and $\sigma(t) = \|\sigma_m(t)\|$ where $\sigma_m(t)$ is the vector of volatilities for the m -th Libor rate $L_m(t)$ (the comment of footnote 2 applies here as well, to Libor volatilities).

The scheme (15.14) would be arbitrage-free in the context of a Gaussian model, i.e. if the expected value in (15.5) were computed in the Gaussian model. In the LM model the equality (15.5) would not hold exactly, but would be a good approximation as long as the choice of the volatility/mean reversion in the scheme is reasonably consistent with the actual LM model volatility structure. While we have no strong opinions on the matter, on the whole (15.14) tends to be our preferred choice for the back-stub interpolation as it is nearly arbitrage-free and is perhaps somewhat easier to implement and maintain than other schemes from Section 15.1.2.

15.1.4 Front Stub, Zero Volatility

Having considered various options for the back stub, let us now focus on the front stub. The simplest way to “complete” the LM model volatility specification is to specify that the front stub bond has no volatility, i.e.

$$\sigma_P(t, T_{q(t)}) \equiv 0 \tag{15.15}$$

in the notation of Section 14.2.3. This automatically specifies the front stub interpolation scheme, as

$$P(t, T_{q(t)}) = P(T_{q(t)-1}, t, T_{q(t)}),$$

where the right-hand side can be computed from the previous results on *back* stub interpolation.

The choice (15.15) was first proposed in Brace et al. [1997] and implies that the continuously rolling money market account $\beta(t)$ coincides with the discrete numeraire $B(t)$, whereby the risk-neutral measure is identical to the spot Libor measure. The bond volatilities for “core” bonds — that is, the bonds paying on the dates in the tenor structure — are explicitly given by

$$\sigma_P(t, T_n) = \sum_{j=q(t)}^{n-1} \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)}.$$

While certainly technically convenient, (15.15) leads to unrealistic dynamics of the yield curve. To give an example, assume that the LM model is specified with a 6 month tenor; i.e. τ_n 's are all approximately 0.5 and $T_1 = 0.5$. For a 3 month option on a 3 month rate, i.e. a caplet with the payoff

$$(L(0.25, 0.25, 0.5) - K)^+, \quad (15.16)$$

(15.15) implies that the rate $L(t, 0.25, 0.5)$ has no volatility for $t \in [0, 0.25]$, i.e. the option with the payoff (15.16) will be priced at its intrinsic value. Clearly, this is unrealistic and will not be consistent with an actual market price. This example is not as contrived as one might think, as many complex derivatives have at least some exposure to short-expiry, short-tenor volatility. The specification (15.15) forces such volatility to zero, and cannot be recommended.

The analysis of the shortcomings of (15.15) puts into focus the requirements that one would wish to impose on all “good” front-stub interpolation schemes. In particular, we will be looking for schemes that recover both market-implied forward rates, as well as volatilities for rates with non-standard expiries and tenors. With “non-standard”, we here refer to rates that are not aligned with the tenor structure. For example, for a 6 month tenor LM model, we can examine the dynamics and values of 3 month tenor forward rates maturing in 1 month, 2 months, ..., or, in 3 months time, we can look at rates with 3 month, 6 month, 9 month, ..., tenors. Note that such an analysis will also help uncover problems with the *back*-stub interpolation, should there be any.

15.1.5 Front Stub, Exogenous Volatility

At a fundamental level, the volatilities $\sigma_P(t, T_{q(t)})$ for all t constitute extra information that is required to specify the model behavior between tenor dates. This information would define the dynamics of the stub bonds $P(t, T_{q(t)})$ (for all t) which, in principle, is all that is required to define the values of *all* bonds, as for any $0 \leq s < t$,

$$P(s, t) = P(s, T_{q(s)}) P(s, T_{q(s)}, T_{q(t)}) E_s^{T_{q(t)}} \left(\frac{1}{P(t, T_{q(t)})} \right). \quad (15.17)$$

As this formula specifies the values of all discount bonds, in principle it makes back-stub interpolation methods from earlier in the section redundant. In practice, however, calculation of the expected value in (15.17) is non-trivial, especially for $s \ll t$, so the application of (15.17) is best left for the case of $t - s$ being small, when accurate approximations are easier to obtain.

The information on the stub bond volatilities $\sigma_P(t, T_{q(t)})$ should be supplied *in addition* to the basic Libor rate volatility structure of the LM model. To make such exogenous volatility specification easier on the model user, we recommend that some additional structure is added through

simplifying assumptions. For instance, an obvious simplification is to assume that the stub volatilities are both time-homogeneous and identical for all periods in the tenor structure; in other words, $\sigma_P(t, u) = \sigma_{\text{stub}}(u - t)$ for some chosen $\sigma_{\text{stub}}(\tau)$, $t \leq u < T_{q(t)}$. This reduces the problem to specifying a stub volatility function of only one argument; this function could be calibrated to short-dated options on short-tenor rates.

In the scheme above, we notice that σ_{stub} is vector-valued. To reduce the specification burden further, we can inherit the correlation structure from that of the core Libor forward rates, and set

$$\sigma_{\text{stub}}(u - t) = \frac{\|\sigma_{\text{stub}}(u - t)\|}{\|\sigma_1(t)\|} \sigma_1(t),$$

where $\sigma_1(t)$ is the vector of volatilities for the front Libor rate $L_1(t)$. With this scheme, we only need to specify a *scalar* function $\|\sigma_{\text{stub}}(\tau)\|$.

The dimensionality of the additional model inputs could be further reduced by using a particular parametric form. Drawing inspiration — yet again — from the one-factor Gaussian model, we could for instance specify

$$\|\sigma_{\text{stub}}(\tau)\| = \bar{\sigma}_{\text{stub}} \frac{1 - e^{-\kappa_{\text{stub}}\tau}}{\kappa_{\text{stub}}} \quad (15.18)$$

for given constants $\bar{\sigma}_{\text{stub}}$ and κ_{stub} . With this, we have reduced the specification of the front-stub volatility function to just two constants.

Once the stub volatility function is specified, it needs to be incorporated into a Monte Carlo simulation. Suppose we need to perform a time step from t to t' where t is one of the tenor dates, $t = T_{m-1}$ and $t' < T_m$. Then, in addition to the standard MC step (see (14.64))

$$\widehat{L}_n(t') = \widehat{L}_n(t) + \sigma_n(t)^\top \mu_n(t) (t' - t) + \widehat{M}_n(t, t'), \quad n = m, \dots, N - 1,$$

we also need an update equation for $P(t', T_m)$. Over the period $s \in [T_{m-1}, T_m]$, the spot numeraire coincides, up to a constant, with the bond price $P(t, T_m)$ (see (14.8)). Hence, in spot measure, the forward bond price $P(s, t', T_m)$, $t \leq s \leq t'$, satisfies

$$\begin{aligned} dP(s, t', T_m) / P(s, t', T_m) &= \|\sigma_{\text{stub}}(s, t', T_m)\|^2 ds \\ &\quad + \sigma_{\text{stub}}(s, t', T_m)^\top dW^B(s), \end{aligned}$$

where we have defined

$$\sigma_{\text{stub}}(s, t', T_m) \triangleq \sigma_{\text{stub}}(T_m - s) - \sigma_{\text{stub}}(t' - s).$$

A simple log-Euler scheme for the bond is given by

$$\begin{aligned} \widehat{P}(t', T_m) &= \widehat{P}(t, t', T_m) \exp \left(\sigma_{\text{stub}}(t, t', T_m)^\top \sqrt{t' - t} Z \right) \\ &\quad \times \exp \left(\frac{1}{2} \|\sigma_{\text{stub}}(t, t', T_m)\|^2 (t' - t) \right), \quad (15.19) \end{aligned}$$

where Z is a draw from the standard one-dimensional Gaussian distribution. This scheme, together with the update equations for $\{\widehat{L}_n(t)\}_{n=m}^{N-1}$, defines the Monte Carlo step. Extensions to more sophisticated discretization schemes for $\widehat{P}(t', T_m)$ are straightforward, but rarely needed.

Since we assumed that $t = T_{m-1}$, the term $\widehat{P}(t, t', T_m)$ in (15.19) is available at time t from Libor rates using *back-stub* interpolation methods discussed earlier in the section. If we, in addition, need to calculate $\widehat{P}(t'', T_m)$ for some t'' , $t' < t'' < T_m$, we can either apply (15.19) to step from t to t'' directly, or to suitably modify (15.19) to step from t' to t'' . In the latter case we will, however, need to be able to simulate other short-dated discount bond values, i.e. $P(t', u)$ for $u < T_m$ (something that may be required for other purposes as well); this can be incorporated into the time-stepping algorithm in the same way as for $P(t', T_m)$.

More pragmatically, instead of simulating a whole collection of short-tenor bonds $P(\cdot, u_1), \dots, P(\cdot, u_k)$, $t' < u_1 < \dots < u_k \leq T_m$, that may be required at time t' , we may choose to propagate a single discount bond with the shortest time to maturity, i.e. $P(\cdot, u_1)$, with other bond prices subsequently obtained by one of the *back-stub* interpolation schemes discussed earlier. Alternatively we could propagate the stub bond $P(\cdot, T_m)$ and use (15.17) to approximate all $P(t', u_1), \dots, P(t', u_{k-1})$ (with, perhaps, some approximation to calculate the expected values):

$$P(t', u_l) = P(t', T_m) E_{t'}^{T_m} \left(\frac{1}{P(u_l, T_m)} \right).$$

Proceeding by either of these methods, the yield curve at time t' , as well as the spot numeraire, are then fully defined by the augmented state vector $\{P(t', T_{q(t')}), L_{q(t')}(t'), \dots, L_{N-1}(t')\}$. Moreover, by construction, the Monte Carlo scheme recovers the expected values of short-tenor rates and their volatilities, as specified by the stub volatility function (within Monte Carlo error and up to discretization bias, of course). The disadvantage of the scheme is the need to carry an extra state variable for each time step.

Remark 15.1.1. Parameterizations such as (15.18) introduce volatility functions (of stub bonds) that are not constant between tenor dates. A similar idea could be applied to core Libor rate volatilities as well, as an extension of the simple piecewise constant interpolation we discussed in Section 14.5.3. This could be required to be able to match shorter-dated options on Libor rates, e.g. a 3 month option on a 6 month Libor rate in an LM model with 6 months between dates in the tenor structure. For example, drawing inspiration from a Gaussian model we could amend (14.42) by a final-period Libor volatility specification of the form

$$\|\lambda_k(t)\| = e^{\varkappa_{\text{short}}(t-T_{k-1})} \|\lambda_k(T_{k-1})\| = e^{\varkappa_{\text{short}}(t-T_{k-1})} \|\lambda_{k,k}\|, \quad t \in [T_{k-1}, T_k],$$

where \varkappa_{short} is user-specified or calibrated to short-dated options.

15.1.6 Front Stub, Simple Interpolation

As an alternative to the scheme in the previous section, one could contemplate avoiding simulation of the stub bond altogether by interpolating it from the original state vector of Libor rates. To demonstrate, let us assume the same setup as in the previous section. Then, at time t' , simulated forward bond prices $P(t', T_n, T_{n+1})$, $n \geq m = q(t')$, are available; the goal is to construct $P(t', T_m)$ from these. A simple interpolation scheme could, for example, specify that the simulated instantaneous forward rates $f(t', u)$ are constant over $u \in [t', T_{m+1}]$. This would link $P(t', T_m)$ to $P(t', T_m, T_{m+1})$ in the following way:

$$P(t', T_m) = P(t', T_m, T_{m+1})^{\frac{T_m - t'}{T_{m+1} - T_m}}. \quad (15.20)$$

This scheme, while simple, does not satisfy our notion of a “good” one. In particular, it does not recover the time 0 forward discount bond value, as generally

$$P(0, t', T_m) \neq E^{t'} \left(P(t', T_m, T_{m+1})^{\frac{T_m - t'}{T_{m+1} - T_m}} \right).$$

Nor does the scheme recover the market volatility of $P(t', T_m)$, as now

$$\ln P(t', T_m) = \frac{T_m - t'}{T_{m+1} - T_m} \ln P(t', T_m, T_{m+1}),$$

and the volatilities of $\ln P(t', T_m)$ and $\ln P(t', T_m, T_{m+1})$ are then linked in a specific way that does not necessarily hold in the market.

Both of the shortcomings above can be addressed in a more sophisticated interpolation scheme, to be discussed in the next section. Before proceeding, however, we note that a similar, albeit somewhat more general, scheme has been proposed by Schlogl [2002]:

$$\frac{1}{P(t', T_m)} = 1 + (T_m - t') (\xi(t') L_{m-1}(T_{m-1}) + (1 - \xi(t')) L_m(t')),$$

where an essentially arbitrary deterministic function $\xi(t)$ is chosen so that $1 = \xi(T_{m-1}) \geq \xi(t') \geq \xi(T_m) = 0$. While the volatility of the stub bond $P(t', T_m)$ could be manipulated by the choice of the function $\xi(t)$, the forward discount bond value is still not recovered by the model.

15.1.7 Front Stub, Interpolation Inspired by the Gaussian Model

The idea of employing a bond reconstitution formula from a Gaussian model, already used for the back stub, can be applied to the front stub as well. Here we assume that for short tenors, the LM model can be locally approximated by a one-factor Gaussian model. Recall that in the latter,

$$P(t', T_m) = P(0, t', T_m) \exp \left(-G(t', T_m) x(t') - \frac{1}{2} G(t', T_m)^2 y(t') \right)$$

and

$$\begin{aligned} P(t', T_m, T_{m+1}) &= P(0, T_m, T_{m+1}) \\ &\quad \times \exp(-(G(t', T_{m+1}) - G(t', T_m)) x(t')) \\ &\quad \times \exp \left(-\frac{1}{2} (G(t', T_{m+1})^2 - G(t', T_m)^2) y(t') \right). \end{aligned}$$

Solving the second equation for $x(t')$ and substituting into the first one, we obtain

$$\begin{aligned} \ln \frac{P(t', T_m)}{P(0, t', T_m)} + \frac{1}{2} G(t', T_m)^2 y(t') &= \frac{G(t', T_m)}{G(t', T_{m+1}) - G(t', T_m)} \\ &\quad \times \left(\ln \frac{P(t', T_m, T_{m+1})}{P(0, T_m, T_{m+1})} + \frac{1}{2} (G(t', T_{m+1})^2 - G(t', T_m)^2) y(t') \right), \end{aligned}$$

so that

$$\begin{aligned} P(t', T_m) &= P(0, t', T_m) \left(\frac{P(t', T_m, T_{m+1})}{P(0, T_m, T_{m+1})} \right)^{\frac{G(t', T_m)}{G(t', T_{m+1}) - G(t', T_m)}} \\ &\quad \times \exp \left(\frac{1}{2} G(t', T_m) G(t', T_{m+1}) y(t') \right). \quad (15.21) \end{aligned}$$

The interpolation scheme depends on two parameters, the mean reversion \varkappa in (15.12) and the short rate volatility $\sigma(t)$ in (15.13) of the Gaussian model. The latter could be approximated by the volatilities of the front Libor rate, ensuring that the forward price of the front-stub discount bond is approximately recovered by the model. The mean reversion plays an important role here as it defines the relative magnitude of the volatility of (log) $P(t', T_m)$ in relationship to (log) $P(t', T_m, T_{m+1})$. As such, \varkappa can be used to set the stub bond volatility to, or near, its market-implied value.

Empirical evidence shows a close fit between time 0 market-observed short-tenor rates and those computed from the model in the manner described above. Also, with the right choice of mean reversion, the market-implied front-stub volatility is recovered as well. We recommend this scheme for most applications.

15.2 Advanced Approximations for Swaption Pricing via Markovian Projection

Having wrapped up the topic of interpolation, let us now go back to the problem of approximate pricing of vanilla options — such as swaptions

— in LM models. While the pricing approximation for swaptions derived in Section 14.4.2 has proved to be remarkably successful for calibration purposes, it does have limitations, especially for longer-dated options on longer-dated swaps and for swaptions that are not at-the-money. There are several improvements that could be made to the basic approximation. For instance, in (14.33) the skew of a swap rate is taken to be the same as the (common) skew of all Libor rates, yet numerical simulation shows that this is not the case. For instance, in a log-normal LM model, a long-expiry, long-tenor swap rate would have a “super log-normal” skew, i.e. implied log-normal swaption volatilities would trend up with the strike. Hence, a more accurate estimate of the swap rate skew from Libor rate skews would be useful. Moreover, the accuracy of the swap rate volatility calculations can be improved by a more careful analysis than that of Section 14.4.2.

Before proceeding, let us first enlarge the model setup from Section 14.2.5 somewhat. Specifically, we wish to address the fact that the specification of the model in Section 14.2.4 (or Section 14.2.5) uses the same time-homogeneous local volatility function for all Libor rates. Such a setup implies that swaptions of all tenors and expiries have essentially the same volatility skew, a model feature that is inconsistent with current market reality. For the Libor market model to be able to match the swaption volatilities for all expiries, tenors *and strikes*, it is necessary to assume that different Libor rates have different local volatility functions, and that those functions explicitly depend on time. In the stochastic volatility model, it may also be necessary to assume that the volatility of variance is time-dependent, so that the curvatures of smiles of different swaptions are allowed to differ. More advanced methods are then required to derive approximations to swap rate volatilities in such a, more generic, specification. In total, we therefore consider the following generalization³ of the specification (14.15)–(14.16):

$$dz(t) = \theta(z_0 - z(t)) dt + \eta(t)\psi(z(t)) dZ(t), \quad (15.22)$$

$$dL_n(t) = \sqrt{z(t)}\varphi_n(t, L_n(t)) \lambda_n(t)^\top dW^{T_{n+1}}(t), \quad n = 1, \dots, N-1. \quad (15.23)$$

In particular, the volatility of variance parameter $\eta(t)$ depends on t , and the DVFs $\varphi_n(t, x)$ are now specific to each Libor rate (i.e. depend on n) and may also depend on time t . Without loss of generality, we assume

$$\varphi_n(t, L_n(0)) = 1.$$

As in other applications of DVF modeling, we assume that the functions $\varphi_n(t, x)$ are well-approximated by their first-order expansions,

³As before, we here assume zero correlation between $Z(t)$ and $W(t)$, but relax this assumption in Section 15.6.

$$\begin{aligned}\varphi_n(t, x) &\approx 1 + b_n(t)(x - L_n(0)), \\ b_n(t) &\triangleq \frac{\partial}{\partial x} \varphi_n(t, L_n(0)).\end{aligned}\tag{15.24}$$

In practical applications, this usually means that $\varphi_n(t, x)$'s are either linear or power functions, see Table 14.1.

As we did in Section 14.6.2.1, let us denote by $\mathbf{L}(t)$ the vector of all Libor rates, i.e. $\mathbf{L}(t) = (L_1(t), \dots, L_{N-1}(t))^\top$, with the convention that $L_i(t) \equiv L_i(T_i)$ for $i < q(t)$. Continuing with the notations of Section 14.4.2, let $S(t) = S_{j,k-j}(t)$ be a particular swap rate. The dynamics of $S(t)$ in the model (15.22)–(15.23) are easy to write down,

$$dS(t) = \sqrt{z(t)} \lambda_S(t, \mathbf{L}(t))^\top dW^A(t),\tag{15.25}$$

where

$$\lambda_S(t, \mathbf{L}(t)) = \sum_{n=j}^{k-1} \frac{\partial S(t)}{\partial L_n(t)} \varphi_n(t, L_n(t)) \lambda_n(t),\tag{15.26}$$

and $W^A(t)$ is a Brownian motion in the annuity measure Q^A for $S(t)$. Moreover, the dynamics can be written in a one-dimensional form,

$$dS(t) = \sqrt{z(t)} \|\lambda_S(t, \mathbf{L}(t))\| dY^A(t),\tag{15.27}$$

where $Y^A(t)$ is a one-dimensional Brownian motion in Q^A .

The following result serves as a starting point for various useful approximations.

Proposition 15.2.1. *For the purposes of European swaption valuation, the dynamics of the swap rate $S(t)$ in Q^A are approximately given by the following displaced log-normal stochastic volatility SDE*

$$dS(t) \approx \sqrt{z(t)} p_S(t) (1 + b_S(t)(S(t) - S(0))) dY^A(t),\tag{15.28}$$

with

$$p_S(t) = \|\lambda_S(t, \mathbf{E}^A(\mathbf{L}(t)))\|,\tag{15.29}$$

$$b_S(t) = \frac{1}{p_S(t)} \sum_{n=j}^{k-1} \frac{\partial \|\lambda_S(t, \mathbf{E}^A(\mathbf{L}(t)))\|}{\partial L_n(t)} \frac{\int_0^t (\lambda_n(u)^\top \lambda_S(u, \mathbf{L}(0))) du}{\int_0^t \|\lambda_S(u, \mathbf{L}(0))\|^2 du}.\tag{15.30}$$

Proof. The proof relies on standard results on Markovian projection, see Appendix A. From (A.18), the European options on $S(t)$ in the model (15.27) have the same values as in the Markovian model

$$dS(t) = \sqrt{z(t)} \left(\mathbf{E}^A \left(\|\lambda_S(t, \mathbf{L}(t))\|^2 \middle| S(t) \right) \right)^{1/2} dY^A(t).$$

First, we approximate

$$\left(\mathbb{E}^A \left(\| \lambda_S(t, \mathbf{L}(t)) \|^2 \mid S(t) \right) \right)^{1/2} \approx \mathbb{E}^A (\| \lambda_S(t, \mathbf{L}(t)) \| \mid S(t)),$$

and then linearize $\| \lambda_S(t, \mathbf{L}(t)) \|$ around $\mathbb{E}^A(\mathbf{L}(t))$,

$$\begin{aligned} \| \lambda_S(t, \mathbf{L}(t)) \| &\approx \| \lambda_S(t, \mathbb{E}^A(\mathbf{L}(t))) \| \\ &\quad + (\nabla \| \lambda_S(t, \mathbb{E}^A(\mathbf{L}(t))) \|) (\mathbf{L}(t) - \mathbb{E}^A(\mathbf{L}(t))), \end{aligned}$$

where $\nabla = (\frac{\partial}{\partial L_1(t)}, \dots, \frac{\partial}{\partial L_{N-1}(t)})$ is the (row-vector) gradient.

Let us introduce the Gaussian approximation

$$d\hat{L}_n(t) = \lambda_n(t)^\top dW^A(t), \quad d\hat{S}(t) = \lambda_S(t, \mathbf{L}(0))^\top dW^A(t), \quad (15.31)$$

so that we can approximate

$$\begin{aligned} \mathbb{E}^A (\mathbf{L}(t) - \mathbb{E}^A(\mathbf{L}(t)) \mid S(t) = s) &\approx \mathbb{E}^A (\hat{\mathbf{L}}(t) - \mathbf{L}(0) \mid \hat{S}(t) = s) \\ &= \frac{\mathbb{E}^A ((\hat{\mathbf{L}}(t) - \mathbf{L}(0)) (\hat{S}(t) - S(0)))}{\mathbb{E}^A ((\hat{S}(t) - S(0)) (\hat{S}(t) - S(0)))} (s - S(0)), \end{aligned}$$

where we have defined $\hat{\mathbf{L}}(t) = (\hat{L}_1, \dots, \hat{L}_{N-1})^\top$. The result follows. \square

The price of a European swaption is given by the value of the option on $S(T_j)$ (times the annuity). The model (15.22) and (15.28) is a stochastic volatility model with time-dependent parameters. Using the methods of Chapter 9, effective, time-constant parameters can be derived, to facilitate fast pricing of European swaptions, as well as calibration of the model parameters to their market values. We do not repeat the relevant formulas here, but simply note that the total volatility, skew and volatility of variance of any swap rate are available as functions of the model parameters.

15.2.1 Advanced Formula for Swap Rate Volatility

In this section, we use the results of Proposition 15.2.1 to derive useful formulas for the swaption volatility. Recall (15.28), the approximate SDE for the swap rate used for European option pricing. The function $p_S(t)$ in (15.29) is the (time-dependent) swap rate volatility,

$$\begin{aligned} p_S(t) &= \| \lambda_S(t, \mathbb{E}^A(\mathbf{L}(t))) \| \\ &= \left\| \sum_{n=j}^{k-1} \frac{\partial S(t)}{\partial L_n(t)} \Big|_{\mathbf{L}(t)=\mathbb{E}^A(\mathbf{L}(t))} \varphi_n(t, \mathbb{E}^A(L_n(t))) \lambda_n(t) \right\|, \end{aligned}$$

where we have used (15.26). The equivalent quantity in the standard approximation of Section 14.4.2 is given by (compare to (14.33))

$$p_{S,\text{standard}}(t) = \left\| \sum_{n=j}^{k-1} \frac{\partial S(t)}{\partial L_n(t)} \Big|_{\mathbf{L}(t)=\mathbf{L}(0)} \varphi_n(t, L_n(0)) \lambda_n(t) \right\|.$$

Hence, the improvements over the standard approximation come from evaluating the actual volatility function of the swap rate ($\lambda_S(t, \mathbf{L}(t))$) at $\mathbf{L}(t) = E^A(\mathbf{L}(t))$ rather than at $\mathbf{L}(t) = \mathbf{L}(0)$; this is similar to the improvements obtained for the quasi-Gaussian model in Section 13.1.4.2. Clearly $E^A(\mathbf{L}(t)) \neq \mathbf{L}(0)$; the difference can be approximated with the help of the following proposition.

Proposition 15.2.2. *For $j \leq n \leq k-1$, the expected value of the n -th Libor rate in the annuity measure is approximately given by*

$$E^A(L_n(t)) = L_n(0)(1 + c_n(t)),$$

where

$$c_n(t) = \frac{1}{L_n(0)Q_n(0)} \sum_{i=j}^{k-1} \frac{\partial Q_n(0)}{\partial L_i(0)} \int_0^t (\lambda_i^\top(s) \lambda_n(s)) ds, \quad Q_n(t) = \frac{A(t)}{P(t, T_{n+1})},$$

with $A(t)$ defined in (14.29).

Proof. We have,

$$E^A(L_n(t)) = Q_n^{-1}(0) E^{T_{n+1}}(Q_n(t)L_n(t)).$$

Both $Q_n(t)$ and $L_n(t)$ are martingales in $Q^{T_{n+1}}$. Applying Gaussian approximations,

$$d\widehat{Q}_n(t) \approx \lambda_{Q_n}(t)^\top dW^{T_{n+1}}(t), \quad d\widehat{L}_n(t) \approx \lambda_n^\top(t) dW^{T_{n+1}}(t),$$

where

$$\lambda_{Q_n}(t) = \sum_{i=j}^{k-1} \frac{\partial Q_n(0)}{\partial L_i(0)} \lambda_i(t),$$

we obtain

$$E^{T_{n+1}}(Q_n(t)L_n(t)) \approx Q_n(0)L_n(0) + \int_0^t (\lambda_{Q_n}(s)^\top \lambda_n(s)) ds,$$

and the statement follows. \square

The idea of employing $E^A(\mathbf{L}(t))$ instead of $\mathbf{L}(0)$ in the standard swap rate volatility approximation can be used independently of any considerations involving time-dependent skews. In particular, the more accurate approximation can be applied directly in the statement of Proposition 14.4.3 when defining the value of $\lambda_S(t)$ in (14.33). The differences between the two formulas are small, but become noticeable for swaptions of longer-dated expiries and tenors.

15.2.2 Advanced Formula for Swap Rate Skew

In the approximate SDE (15.28) for the swap rate used for European option pricing, the parameter $b_S(t)$ controls the skew of the volatility smile. Define

$$v_n(t) = \frac{\partial S(t)}{\partial L_n(t)}, \quad n = j, \dots, k-1,$$

(compare to (14.31)), and

$$v_{n,n'}(t) = \frac{\partial^2 S(t)}{\partial L_n(t) \partial L_{n'}(t)}, \quad n, n' = j, \dots, k-1.$$

Proposition 15.2.3. *The time-dependent swaption skew $b_S(t)$ is approximately given by*

$$b_S(t) = \sum_{i,n=j}^{k-1} \frac{r_{S,n}(t)}{r_{S,i}(t)} v_{i,n}(0) \xi_i(t) + \sum_{i=j}^{k-1} b_i(t) v_i(0) \xi_i(t),$$

where

$$\begin{aligned} r_{S,i}(t) &= \lambda_i(t)^\top \lambda_S(t, \mathbf{L}(0)), \quad r_S(t) = \|\lambda_S(t, \mathbf{L}(0))\|^2, \\ \xi_i(t) &= \frac{r_{S,i}(t) \int_0^t r_{S,i}(u) du}{r_S(t) \int_0^t r_S(u) du}. \end{aligned}$$

Proof. Recall from Proposition 15.2.1,

$$\begin{aligned} b_S(t) &= \frac{1}{p_S(t)} \sum_{i=j}^{k-1} \frac{\partial \|\lambda_S(t, \mathbf{L}(t))\|}{\partial L_i(t)} \Big|_{\mathbf{L}(t)=\mathbf{E}^A(\mathbf{L}(t))} \frac{\int_0^t r_{S,i}(u) du}{\int_0^t r_S(u) du} \\ &= \sum_{i=j}^{k-1} \frac{\partial \ln \|\lambda_S(t, \mathbf{L}(t))\|}{\partial L_i(t)} \Big|_{\mathbf{L}(t)=\mathbf{E}^A(\mathbf{L}(t))} \frac{\int_0^t r_{S,i}(u) du}{\int_0^t r_S(u) du}. \end{aligned}$$

We have

$$\begin{aligned} \frac{\partial \ln \|\lambda_S(t, \mathbf{L}(t))\|}{\partial L_i(t)} &= \frac{1}{\|\lambda_S(t, \mathbf{L}(t))\|^2} \\ &\times \sum_{n,n'=j}^{k-1} (\lambda_n(t)^\top \lambda_{n'}(t)) (v_n(t) \varphi_n(t, L_n(t))) \frac{\partial}{\partial L_i(t)} (v_{n'}(t) \varphi_{n'}(t, L_{n'}(t))). \end{aligned}$$

While using $\mathbf{E}^A(\mathbf{L}(t))$ instead of $\mathbf{L}(0)$ in the calculations of the swaption volatility (see the previous section) results in noticeable improvements in the approximation quality, this turns out to not be the case for approximations to the skew. Furthermore, as $\varphi_n(t, L_n(0)) = 1$ for any n , the formulas resulting

from evaluating the expression for $b_S(t)$ above at $L_n(0)$ are compact and convenient, which in practice will justify any slight deterioration in precision. With this in mind, we note that

$$\frac{\partial}{\partial L_i(t)} (v_{n'}(t) \varphi_{n'}(t, L_{n'}(t))) \Big|_{\mathbf{L}(t)=\mathbf{L}(0)} = v_{i,n'}(0) + 1_{\{i=n'\}} v_{n'}(0) b_{n'}(t)$$

(recall (15.24) for the definition of b_n 's). Hence,

$$\begin{aligned} \frac{\partial \ln \|\lambda_S(t, \mathbf{L}(t))\|}{\partial L_i(t)} \Big|_{\mathbf{L}(t)=\mathbf{L}(0)} &= \frac{1}{r_S(t)} \sum_{n,n'=j}^{k-1} (\lambda_n(t)^\top \lambda_{n'}(t)) v_n(0) v_{i,n'}(0) \\ &\quad + \frac{1}{r_S(t)} b_i(t) v_i(0) \sum_{n=j}^{k-1} (\lambda_n(t)^\top \lambda_i(t)) v_n(0). \end{aligned}$$

We recognize

$$\sum_{n=j}^{k-1} (\lambda_n(t)^\top \lambda_{n'}(t)) v_n(0) = \lambda_{n'}(t)^\top \lambda_S(t, \mathbf{L}(0)) = r_{S,n'}(t),$$

so that

$$\begin{aligned} \frac{\partial \ln \|\lambda_S(t, \mathbf{L}(t))\|}{\partial L_i(t)} \Big|_{\mathbf{L}(t)=\mathbf{L}(0)} &= \left(\sum_{n'=j}^{k-1} v_{i,n'}(0) \frac{r_{S,n'}(t)}{r_S(t)} \right) + b_i(t) v_i(0) \frac{r_{S,i}(t)}{r_S(t)}. \end{aligned}$$

The result follows. \square

Remark 15.2.4. The swaption skew consists of two parts, one that involves second-order derivatives of the swap rate with respect to Libor rates and captures overall convexity of a swap rate with respect to Libor rates, and the other being a weighted average of Libor skews.

Remark 15.2.5. Even with Libor rates sharing a common skew, i.e. in the “classic” LM model of Sections 14.2.4, 14.2.5, the swaption skew is not exactly equal to the (shared) Libor skews. If $b_i(t) = b$ for all i and t , then

$$b_S(t) = \sum_{i,n=j}^{k-1} \frac{r_{S,n}(t)}{r_{S,i}(t)} v_{i,n}(0) \xi_i(t) + b \sum_{i=j}^{k-1} v_i(0) \xi_i(t).$$

Even if the convexity term is ignored, we have

$$b_S(t) = b \sum_{i=j}^{k-1} v_i(0) \xi_i(t) \neq b.$$

15.2.3 Skew and Smile Calibration in LM Models

With pricing issues out of the way, let us now turn our attention to calibration, and see what modifications to the calibration algorithm of Section 14.5 are required by the more general model specification (15.22)–(15.23). We assume a typical set of swaption volatilities is given, and we have available a collection of market-implied volatility smiles across strikes for swaptions of different expiries and tenors. The model (15.22)–(15.23) has enough parameter flexibility to match

- At-the-money volatilities of all European swaptions, using the volatility structure $\|\lambda_n(t)\|$.
- Slopes of volatility smiles (skews) of all European swaptions, using the skew structure $b_n(t)$ (see (15.24)).
- Curvatures of smiles for swaptions of different expiries, for a given tenor, using term structure of volatilities of variance $\eta(t)$.

Of course, all these parameters are in addition to the correlation structure of Libor rates, as in the standard LM model case. Note that, for the last point, technically there is no flexibility in the model to change the volatility smile curvature for swaptions of the same expiry but different tenors. This is not really a serious limitation as the curvature of the smile tends to be constant across tenors for a given expiry.

Assume, as in Section 14.5, that we have chosen calibration targets that include N_S swaptions, $V_{\text{swaption},1}, V_{\text{swaption},2}, \dots, V_{\text{swaption},N_S}$, and let us ignore caps; the considerations below extend trivially to cover them as well. Unlike previously, however, let us assume that each target includes not one swaption, but a collection of them of different strikes. Hence, we redefine the calibration as the goal to match volatility smiles of N_S swaptions.

Having chosen a (vanilla) stochastic volatility model for European swaptions such as (8.3)–(8.4), the target volatility smiles can be summarized by a collection of market-implied SV parameters, namely volatilities $\widehat{\lambda}_{S_i}$, skews \widehat{b}_{S_i} and volatilities of variance $\widehat{\eta}_{S_i}$ for $i = 1, \dots, N_S$ (a common mean reversion of volatility parameter is assumed). Recall that in Section 14.5.2 we denoted by G a grid of the Libor volatilities $\|\lambda_n(t)\|$. To be able to generalize, redefine $G_\lambda = G$ and, in the same spirit, define G_b to be the grid of Libor skews $b_n(t)$, and G_η (a vector) to be the discretized term structure of volatilities of variance $\eta(t)$. The formulas from the previous section allow us to compute term volatilities, skews and volatilities of variance from the model. We denote them by

$$\bar{\lambda}_{S_i}(G_\lambda, G_b, G_\eta), \quad \bar{b}_{S_i}(G_\lambda, G_b, G_\eta), \quad \bar{\eta}_{S_i}(G_\lambda, G_b, G_\eta).$$

One can incorporate the skew and smile calibration in the algorithm of Section 14.5 by adding extra terms to the calibration norm, replacing (14.54) with

$$\begin{aligned}
\mathcal{I}(G_\lambda, G_b, G_\eta) = & \frac{w_{S,\lambda}}{N_S} \sum_{i=1}^{N_S} \left(\bar{\lambda}_{S_i}(G_\lambda, G_b, G_\eta) - \hat{\lambda}_{S_i} \right)^2 \\
& + \frac{w_{S,b}}{N_S} \sum_{i=1}^{N_S} \left(\bar{b}_{S_i}(G_\lambda, G_b, G_\eta) - \hat{b}_{S_i} \right)^2 \\
& + \frac{w_{S,\eta}}{N_S} \sum_{i=1}^{N_S} \left(\bar{\eta}_{S_i}(G_\lambda, G_b, G_\eta) - \hat{\eta}_{S_i} \right)^2 \\
& + \dots,
\end{aligned} \tag{15.32}$$

with dots denoting various regularization terms for G_λ , G_b , and G_η . This, however, is not necessarily the best approach, as it increases the number of degrees of freedom in the non-linear optimization problem quite substantially, thus potentially significantly reducing the speed at which it could be solved. It is much better to take advantage of the structure of the problem and solve for volatilities, skews and volatilities of variance *separately* rather than all at the same time.

To motivate the method, we note that the impact of changes in the Libor skews G_b on term swaption volatilities, or their volatilities of variance, is rather small. Likewise, changes in Libor volatilities G_λ have only a small impact on term swaption skews, and so on. This near-orthogonality allows us to solve for various parameters sequentially. To facilitate this approach, we define three norms

$$\begin{aligned}
\mathcal{I}_\lambda(G_\lambda, G_b, G_\eta) &= \frac{w_{S,\lambda}}{N_S} \sum_{i=1}^{N_S} \left(\bar{\lambda}_{S_i}(G_\lambda, G_b, G_\eta) - \hat{\lambda}_{S_i} \right)^2 + \dots, \\
\mathcal{I}_b(G_\lambda, G_b, G_\eta) &= \frac{w_{S,b}}{N_S} \sum_{i=1}^{N_S} \left(\bar{b}_{S_i}(G_\lambda, G_b, G_\eta) - \hat{b}_{S_i} \right)^2 + \dots, \\
\mathcal{I}_\eta(G_\lambda, G_b, G_\eta) &= \frac{w_{S,\eta}}{N_S} \sum_{i=1}^{N_S} \left(\bar{\eta}_{S_i}(G_\lambda, G_b, G_\eta) - \hat{\eta}_{S_i} \right)^2 + \dots,
\end{aligned}$$

and modify the algorithm from Section 14.5.7 as such.

1. First, make a guess for G_b and G_η , denoted by G_b^0 and G_η^0 . The guesses could be quite approximate. For example the grid G_b^0 could be set to the same value, an average of swaption term skews $N_S^{-1} \sum_{i=1}^{N_S} \hat{b}_{S_i}$, and the same for G_η^0 .
2. Perform steps 1–5 from Section 14.5.7 with $\mathcal{I} = \mathcal{I}_\lambda(G_\lambda, G_b^0, G_\eta^0)$ until G_λ^1 , the solution, is found. Note that we keep the skew and volatility of variance grids constant throughout
3. Minimize $\mathcal{I}_b(G_\lambda^1, G_b, G_\eta^0)$ by iterating over G_b ; denote the solution by G_b^1 .
4. Minimize $\mathcal{I}_\eta(G_\lambda^1, G_b^1, G_\eta)$ by iterating over G_η ; denote the solution by G_η^1 .

Typically, the triple of parameters $(G_\lambda^1, G_b^1, G_\eta^1)$ provides a good overall solution. If a better fit is desired, the steps could be iterated, starting with (G_b^1, G_η^1) on Step 1. If the number of iterations is more than 1, it could be beneficial to stop after Step 2 (of the second or subsequent iterations) in order to have the best possible fit to the volatilities, usually the most important target.

15.3 Near-Markov LM Models

The LM model, even if driven by a single Brownian motion, is not Markovian in a low number of state variables. Rather, the state vector comprises all Libor rates and whatever state variables are needed for the stochastic volatility, if present. This is readily seen from the expression for the drift of each Libor rate (see e.g. (14.7)) which involves multiple other rates.

This state of affairs, however, has not stopped some researchers from attempting to *approximate* an LM model with a low-dimensional Markovian one. These attempts mostly involve i) restricting the volatility structure to a “separable” form, similar to that used to specify Markovian models in Chapters 12 and 13; and ii) removing path dependence in the Libor drifts by either “freezing” them at time 0 values or by employing various tricks not unlike those discussed previously in the context of drift estimation for large time steps in Monte Carlo.

The practical usefulness of such approximations for pricing and risk managing of derivatives is limited. The imposed restrictions on the volatility structure remove many of the main LM advantages in terms of calibration flexibility, and the necessity of approximating the Libor rate drifts makes the model arbitrageable and problematic to use for anything other than short-dated derivatives. Fundamentally, there is also something misplaced about trying to force LM models into a low-dimensional Markovian setting: if such a setting is desired, there are really no appealing reasons to use an LM model in the first place. Instead, one should from the outset pick one of the many perfectly good low-dimensional Markovian models we have covered earlier in the book. The models in Chapter 13 are particularly appealing, we think.

That said, a low-dimensional Markov approximation to an LM model can, however, still find use as a type of a model-based control variate. We discuss this application in more detail in Chapter 25; this chapter also covers the mechanics of the various steps involved in creating near-Markov LM models.

15.4 Swap Market Models

In a nutshell, the LM modeling principle revolves around using models for simple forward rates (Libor) that become tractable in properly selected

martingale measures. Longer-dated swap rates can be constructed iteratively by, in effect, adding up the individual forward rates that constitute the LM model primitives. As we have seen in Section 14.4.2, the swap rates constructed in this manner rarely have tractable dynamics, and swaption pricing formulas will nearly always involve approximations and/or numerical methods.

An alternative modeling approach due to Jamshidian [1997] turns the LM philosophy on its head by using forward swap rates as the fundamental model primitives, and constructing individual Libor rates as differences of such swap rates. This is similar to the idea behind swap Markov-functional model construction that we briefly discussed in Appendix 11.A.3 of Chapter 11.

By specifying tractable dynamics for forward swap rates, swaption pricing formulas now often can be done exactly, whereas cap pricing must be done by approximation. While this so-called *swap market model* framework has some proponents (e.g. Galluccio and Hunter [2003], Galluccio et al. [2005]), it is fair to say that the approach has attracted much less interest from practitioners and academics than has the LM model. Our treatment of swap market models shall therefore be brief.

Given a tenor structure $0 = T_0 < T_1 < \dots < T_N$, our market primitives are now the swap rates

$$S_j(t) = S_{j,N-j}(t), \quad 1 \leq j < N,$$

where the notation is similar to that used in Section 14.4.2:

$$S_j(t) = \frac{P(t, T_j) - P(t, T_N)}{A_j(t)}, \quad A_j(t) = A_{j,N-j}(t) = \sum_{n=j}^{N-1} P(t, T_{n+1})\tau_n.$$

We emphasize that all forward swaps here have identical terminal maturity, namely T_N . We also emphasize that the information content of the primitives of a swap market model is identical to that of an LM model, as the forward Libor rates curve can be uniquely constructed from knowledge of $S_i(t)$, for all $i \in [q(t), N-1]$. To show this, write for any j

$$\begin{aligned} L_{j-1}(t) &= \frac{S_{j-1}(t)A_{j-1}(t) - S_j(t)A_j(t)}{P(t, T_j)\tau_{j-1}} \\ &= S_{j-1}(t) \frac{A_{j-1}(t)}{P(t, T_j)\tau_{j-1}} - S_j(t) \frac{A_{j-1}(t) - P(t, T_j)\tau_{j-1}}{P(t, T_j)\tau_{j-1}} \\ &= S_j(t) + (S_{j-1}(t) - S_j(t)) \frac{A_{j-1}(t)}{P(t, T_j)\tau_{j-1}}. \end{aligned}$$

As

$$\frac{P(t, T_{n+1})}{P(t, T_j)} = \frac{P(t, T_{j+1})}{P(t, T_j)} \frac{P(t, T_{j+2})}{P(t, T_{j+1})} \dots \frac{P(t, T_{n+1})}{P(t, T_n)} = \prod_{k=j}^n (1 + L_k(t)\tau_k)^{-1},$$

it follows from the definition of $A_{j-1}(t)$ that

$$L_{j-1}(t) = S_j(t) + (S_{j-1}(t) - S_j(t)) \left(1 + \sum_{n=j}^{N-1} \prod_{k=j}^n (1 + L_k(t)\tau_k)^{-1} \frac{\tau_n}{\tau_{j-1}} \right). \quad (15.33)$$

Starting from $L_{N-1}(t) = S_{N-1}(t)$, equation (15.33) gives us an iterative formula to construct the Libor forward curve at time t , as claimed.

In the swap market model framework, the modeler specifies dynamics on the swap forward rates in their respective annuity measures. Let $W^{A_j}(t)$ be an m -dimensional Brownian motion in the annuity measure induced by $A_j(t)$. As shown earlier, $S_j(t)$ is a martingale in this measure, and we can write

$$dS_j(t) = \sigma_{S_j}(t)^\top dW^{A_j}(t), \quad j = q(t), \dots, N-1, \quad (15.34)$$

for some adapted volatility function $\sigma_{S_j}(t)$ specific to the j -th swap rate. As for an LM model, we could now proceed to specialize the model by using DVF or SV type specifications for $\sigma_{S_j}(t)$, generally keeping an eye out for models that allow construction of easy-to-compute expressions for payer swaption prices

$$V_{\text{swaption},j}(t) = A_j(t) E_t^{A_j} \left((S_j(T_j) - c)^+ \right).$$

As we mentioned earlier, no exact pricing formulas for caplets will normally exist (as should be obvious from the complicated form of (15.33)), a situation that also holds for any swaption that does not involve a swap maturing at time T_N . For these instruments, approximative formulas must be devised if a quick calibration of the model is desired. See Galluccio et al. [2005] for some details on this.

In the context of the LM models, Section 14.2.2 derived relations between the different forward martingale measures, allowing us to describe the dynamics of all forward rates in a single measure, as required in Monte Carlo simulations, say. For the swap market models, starting with the specification (15.34) we can derive similar relations between the different annuity measures, ultimately allowing for simulation of all swap rates S_1, S_2, \dots, S_{N-1} in a common measure. For details, the reader is referred to Jamshidian [1997] and Section 14.4 of Musiela and Rutkowski [1997].

15.5 Evolving Separate Discount and Forward Rate Curves

A single yield curve for discounting and calculating Libor rates is not always compatible with no-arbitrage constraints of cross-currency markets, nor is it

particularly realistic in stressed market conditions, such as those experienced during the subprime crisis of 2007–2009. As explained in Section 6.5.2.2, separating the discounting curve from the forward, or index, curve will ensure that linear instruments (i.e. swaps and bonds) are correctly priced at time 0. In this section we consider how to incorporate the idea of curve separation into a *dynamic* model of interest rates. While we use an LM setting for some of the material in this section, the basic ideas are generic and can be applied to virtually all models in this book.

15.5.1 Basic Ideas

Suppose two yield curves are given at time 0, the *discount* curve $P(0, T)$, $T \geq 0$, and the *index* curve $\tilde{P}(0, T)$, $T \geq 0$ (note the change of notation from $P^{(L)}$ in Section 6.5.2.2 to \tilde{P} for convenience). The index curve corresponds to a particular Libor tenor τ , and is defined through the requirement that forward Libor rates of tenor τ must equal the conditional expected values of the future spot Libor rates,

$$\tilde{L}(t, T, T + \tau) = E_t^{T+\tau} (L(T, T, T + \tau)), \quad (15.35)$$

where $E_t^{T+\tau}$ denotes expectation in the $(T + \tau)$ -forward measure and, by definition,

$$\tilde{L}(t, T, T + \tau) = \frac{1}{\tau} \left(\frac{1}{\tilde{P}(t, T, T + \tau)} - 1 \right). \quad (15.36)$$

For emphasis, we have added a tilde to the regular Libor rate notation to highlight the link to the index curve. Note that we take (15.35) to hold for arbitrary values of $T \geq t$.

It is important to point out that using different Libor tenors τ in the equation above would lead to *different* models of two-curve evolution, and our first choice shall focus on what tenor we actually want to use. In the next section we will be looking at extending the *HJM* instantaneous forward rate formalism, so our choice will be influenced by that fact. Later, in Section 15.5.3, we will consider Libor market models that, not surprisingly, lead to a somewhat different choice.

As we initially focus on adding a second curve evolution in the HJM setting, it is convenient to consider a specific choice of $\tau = 0$ that corresponds to instantaneous forward rates. In particular, we note that for $\tau \rightarrow 0$, $\tilde{L}(t, T, T + \tau)$ converges to $\tilde{f}(t, T)$, so according to (15.35), for the next section we shall require that

$$\tilde{f}(t, T) = E_t^T (\tilde{f}(T, T)). \quad (15.37)$$

As discussed in Chapter 6 (see (6.44)–(6.45)), it is convenient to represent the index curve through an additive spread ε in continuously compounded forward rates. Specifically, at time 0 we write

$$\tilde{P}(0, T) = P(0, T) e^{\int_0^T \varepsilon(0, u) du}, \quad T \geq 0. \quad (15.38)$$

We already hinted in Section 6.5.2.4 that a standard way of including the spread in a dynamic model is to assume that it evolves deterministically. We will ultimately make the same recommendation here, but it is still instructive to develop a generic dynamic model of two-curve evolution first. In developing such a model, we have several possible choices for the model primitives. For example, we could impose dynamics on the index and discount curves in separation. Alternatively, we could impose dynamics on the spread and just one of the two curves, deriving the remaining curve from these two primitives. It is often desirable to have direct control of the spread process — e.g. to ensure that its dynamics and range are in line with historical observations — so we adopt the latter approach here. As for the choice of which yield curve to use as a direct model ingredient, this is largely a matter of taste and convenience. We shall demonstrate both approaches below, using the discount curve as the primary curve in the HJM model setting of Section 15.5.2; and the index curve as the primary curve in the LM model setting of Section 15.5.3.

15.5.2 HJM Extension

The framework for discount bonds is unchanged by the presence of an index curve, and we can therefore start with the standard HJM dynamics for the discount factors (see Section 4.4.1) in the risk-neutral measure Q ,

$$dP(t, T)/P(t, T) = r(t) dt - \sigma_P(t, T)^\top dW(t),$$

where $W(t)$ is a d -dimensional Brownian motion and $\sigma_P(t, T)$ is a d -dimensional discount bond volatility function. Importantly, the T -forward measure Q^T is determined solely by the dynamics of the discount curve because the numeraire, the zero-coupon discount bond $P(t, T)$, is just a particular point on the time t discount curve. In particular, the T -forward measure is defined by the familiar requirement that

$$dW^T(t) = dW(t) + \sigma_P(t, T) dt$$

be a driftless Brownian motion.

As for the instantaneous forward rates $f(t, T)$, the standard HJM result (see Lemma 4.4.1) shows that their dynamics are given by

$$df(t, T) = \sigma_f(t, T)^\top \sigma_P(t, T) dt + \sigma_f(t, T)^\top dW(t). \quad (15.39)$$

Let $\varepsilon(t, u)$ be the *forward rate spread* defined by extending (15.38) to arbitrary t ,

$$\tilde{P}(t, T) = P(t, T) e^{\int_t^T \varepsilon(t, u) du}, \quad T \geq t. \quad (15.40)$$

Treating ε as a forward rate spread is justified by considering $\tilde{f}(t, T)$, the instantaneous forward rates calculated from the index curve $\tilde{f}(t, T) = -\partial \ln \tilde{P}(t, T) / \partial T$, and observing that

$$\tilde{f}(t, T) = f(t, T) - \varepsilon(t, T), \quad 0 \leq t \leq T. \quad (15.41)$$

Let us endow $\varepsilon(t, T)$ with the following dynamics⁴,

$$d\varepsilon(t, T) = \alpha_\varepsilon(t, T) dt + \sigma_\varepsilon(t, T)^\top dW(t), \quad (15.42)$$

for some, yet unspecified, adapted processes $\alpha_\varepsilon(t, T)$ and $\sigma_\varepsilon(t, T)$, the latter d -dimensional. Not surprisingly, the drift term in (15.42) cannot be set arbitrarily.

Proposition 15.5.1. *To satisfy the martingale restrictions of (15.37), the drift of $\varepsilon(t, T)$ in (15.42) must obey*

$$\alpha_\varepsilon(t, T) = \sigma_\varepsilon(t, T)^\top \sigma_P(t, T). \quad (15.43)$$

Proof. According to (15.37), $\tilde{f}(t, T)$ must be a martingale in the T -forward measure. Its dynamics in that measure follow from (15.39)–(15.42),

$$\begin{aligned} d\tilde{f}(t, T) &= (\sigma_f(t, T)^\top \sigma_P(t, T) - \alpha_\varepsilon(t, T)) dt \\ &\quad + (\sigma_f(t, T) - \sigma_\varepsilon(t, T))^\top dW^T(t) - (\sigma_f(t, T) - \sigma_\varepsilon(t, T))^\top \sigma_P(t, T) dt, \end{aligned}$$

and the result in the proposition follows by setting the dt term to zero. \square

The SDEs for various quantities related to the index curve can now be derived from the discount and spread parameters.

Proposition 15.5.2. *Define*

$$\tilde{\alpha}_f(t, T) \triangleq \tilde{\sigma}_f(t, T)^\top \sigma_P(t, T), \quad \tilde{\sigma}_f(t, T) \triangleq \sigma_f(t, T) - \sigma_\varepsilon(t, T).$$

The dynamics of $\tilde{f}(t, T)$ in the risk-neutral measure are then

$$d\tilde{f}(t, T) = \tilde{\alpha}_f(t, T) dt + \tilde{\sigma}_f(t, T)^\top dW(t).$$

Further, set $\tilde{r}(t) = \tilde{f}(t, t)$ and

⁴Arguably, diffusive dynamics do not reflect the observed movements of the spread particularly well, as the spread between discount and index curves tends to remain stable for extended period of times, followed by sometimes violent dislocations. On the other hand, our needs in capturing the distribution of the spread within a model are usually quite modest, rarely exceeding the requirement to capture an overall level of its dispersion. As long as the volatility of the spread is reasonable, this requirement will be satisfied by a diffusion model.

$$\begin{aligned}\tilde{\alpha}_P(t, T) &\triangleq \int_t^T \tilde{\sigma}_f(t, u)^\top (\tilde{\sigma}_P(t, u) - \sigma_P(t, u)) du, \\ \tilde{\sigma}_P(t, T) &\triangleq \int_t^T \tilde{\sigma}_f(t, u) du.\end{aligned}$$

Then

$$d\tilde{P}(t, T)/\tilde{P}(t, T) = \tilde{r}(t) dt + \tilde{\alpha}_P(t, T) dt - \tilde{\sigma}_P(t, T)^\top dW(t).$$

Proof. The result for $d\tilde{f}(t, T)$ follows directly from (15.39) and (15.41)–(15.43). From the equation

$$\tilde{P}(t, T) = e^{-\int_t^T \tilde{f}(t, u) du},$$

it follows that $Y(t, T) = \ln(\tilde{P}(t, T))$ satisfies

$$\begin{aligned}dY(t, T) &= \tilde{f}(t, t) dt - \int_t^T \tilde{\alpha}_f(t, u) du dt - \tilde{\sigma}_P(t, T)^\top dW(t) \\ &= \tilde{r}(t) dt - \tilde{\sigma}_P(t, T)^\top dW(t) - \int_t^T \tilde{\sigma}_f(t, u)^\top \sigma_P(t, u) du dt.\end{aligned}$$

An application of Ito's lemma (to e^Y) then shows that

$$d\tilde{P}(t, T)/\tilde{P}(t, T) = \tilde{r}(t) dt + \tilde{\alpha}_P(t, T) dt - \tilde{\sigma}_P(t, T)^\top dW(t),$$

where

$$\tilde{\alpha}_P(t, T) = - \int_t^T \tilde{\sigma}_f(t, u)^\top \sigma_P(t, u) du + \frac{1}{2} \tilde{\sigma}_P(t, T)^\top \tilde{\sigma}_P(t, T).$$

Integration by parts shows that

$$\tilde{\sigma}_P(t, T)^\top \tilde{\sigma}_P(t, T) = 2 \int_t^T \tilde{\sigma}_f(t, u)^\top \tilde{\sigma}_P(t, u) du,$$

and the result follows. \square

Remark 15.5.3. Note the presence of an “extra” drift term $\tilde{\alpha}_P(t, T)$ in the dynamics for the pseudo-bonds $\tilde{P}(t, T)$; sometimes, in an analogy with foreign exchange markets, this term is called a *quanto correction*.

For future use, let us define

$$Z(t, T) \triangleq e^{-\int_t^T \varepsilon(t, u) du}, \quad (15.44)$$

so that

$$\tilde{P}(t, T) = P(t, T)/Z(t, T), \quad 0 \leq t \leq T. \quad (15.45)$$

The drift and diffusion terms in the process for $Z(t)$,

$$dZ(t, T)/Z(t, T) = (r(t) - \tilde{r}(t)) dt + \alpha_Z(t, T) dt - \sigma_Z(t, T)^\top dW(t), \quad (15.46)$$

are given by

$$\alpha_Z(t, T) = -\tilde{\alpha}_P(t, T) - \sigma_Z(t, T)^\top \tilde{\sigma}_P(t, T), \quad \sigma_Z(t, T) = \int_t^T \sigma_\varepsilon(t, u) du. \quad (15.47)$$

Options linked directly to the spread between the index and discount curves are rarely, if ever, traded, so a high level of sophistication in basis spread dynamics is seldom required. A simple one-factor Gaussian model for the spread, say, is more than sufficient for most applications. Nevertheless, richer dynamics are possible. Recall that if the forward rate volatility function is of separable form,

$$\sigma_f(t, T) = g(t)h(T),$$

where $g(\cdot)$ is a $d \times d$ matrix-valued process and $h(\cdot)$ is a d -dimensional deterministic vector-valued function (see e.g. (4.44), (12.2) or (13.70)), then the discount curve admits a d -dimensional Markovian representation, see, for example, Proposition 13.3.1. If, in addition, the forward rate spread volatility function $\sigma_\varepsilon(t, T)$ is also of separable form with the same $g(t)$ but a different $h_\varepsilon(T)$,

$$\sigma_\varepsilon(t, T) = g(t)h_\varepsilon(T),$$

then the forward rate spread curve $\varepsilon(t, T)$ and, by extension, the index curve also admit Markovian representations with the same Markovian state variables. This fact opens up a possibility of building efficient Markovian models of joint evolution of the discount and the index curve with non-deterministic spread. We leave this line of inquiry for the reader to explore, and instead move on to the changes required to extend the LM model framework to support two-curve dynamics.

15.5.3 Applications to LM Models

We now return to the LM framework, and continue with the notations of Section 14.2. In particular, we assume that a tenor structure (14.1) has been specified. Clearly, in the LM setup, it is natural to assume that (15.36) is satisfied for the Libor tenor used in the definition of the LM model. This would correspond to extending the LM model to drive both the discounting curve and the *index curve that corresponds to the Libor tenor used in LM model definition*. We remind the reader that if multiple index curve dynamics are required, each will have to be driven by its own LM model (of a given Libor tenor).

As discussed earlier, we here choose as our primary model variables (in addition to the spread) the set of “regular” Libor rates $\tilde{L}_n(t) \triangleq \tilde{L}(t, T_n, T_{n+1})$,

$n = 0, \dots, N - 1$, as defined off the index curve in (15.36). Apart from giving ourselves a chance to present a slightly different twist on the material of the previous section, the choice of the index Libors as building blocks allows us to use the volatilities of $\{\tilde{L}_n(t)\}_{n=0}^{N-1}$ that are more directly related to the market-observable volatility information, i.e. the volatilities of swap rates of various maturities and tenors. Conveniently, for each n , $\tilde{L}_n(t)$ is a martingale in T_{n+1} -forward measure by construction, see (15.35). Therefore, in close analogy to (14.5), we can define the dynamics by

$$d\tilde{L}_n(t) = \tilde{\sigma}_n(t)^\top dW^{n+1}(t), \quad (15.48)$$

where $W^{n+1}(t)$ is an m -dimensional Brownian motion in $Q^{T_{n+1}}$ and $\tilde{\sigma}_n(t)$ is an m -dimensional adapted process.

The discrete money market account $B(t)$, see (14.8), induces a useful common measure for all Libor rates in the single-curve case. The extension to the two-curve framework is straightforward: we apply (14.8) literally, i.e. use simply compounded forward rates computed off the discount curve in the definition of $B(t)$. This corresponds exactly to the actual trading strategy of re-investing into deposits of a given tenor, with the value of the strategy solely determined by the discount curve:

$$B(t) = P(t, T_{q(t)}) \prod_{n=0}^{q(t)-1} \frac{1}{P(T_n, T_{n+1})}. \quad (15.49)$$

We define the spot measure Q^B accordingly. Interestingly, a measure change from any T -forward measure to the spot measure is determined by the dynamics of discount factors $P(t, T)$ only, and we see that if \tilde{L}_n 's satisfy (15.48), then in the spot measure Q^B , the process for \tilde{L}_n is given by

$$d\tilde{L}_n(t) = \tilde{\sigma}_n(t)^\top \left(\sum_{j=q(t)}^n \frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)} dt + dW^B(t) \right), \quad (15.50)$$

where $W^B(t)$ is an m -dimensional Brownian motion in measure Q^B . Here $L_j(t)$ are simply compounded forward rates calculated off the discount curve, and $\sigma_j(t)$ are their volatilities. It would, however, be more convenient to have the drift in terms of the Libor rate model primitives $\tilde{L}_j(t)$ and the process for the spread.

Before deriving the result for the drifts, let us decide what variables to use to represent the spread evolution. Using the instantaneous forward spread curve $\varepsilon(t, T)$ is not convenient in the LM framework, and we settle on *forward bond ratios* $Z^n(t)$, $n = 0, \dots, N - 1$, as state variables, defined by

$$Z^n(t) \triangleq Z(t, T_{n+1})/Z(t, T_n), \quad (15.51)$$

with the $Z(t, T)$'s given in (15.44). This choice of state variables keeps the LM framework in line with (but not exactly equivalent to!) the HJM extension

from Section 15.5.2, although other choices are possible; we comment on that at the end of the section.

Let us assume that forward bond ratios follow the dynamics

$$dZ^n(t)/Z^n(t) = O(dt) - \sigma_{Z^n}(t)^\top dW^B(t) \quad (15.52)$$

in the spot measure, with $\sigma_{Z^n}(t)$'s being their volatilities. Note that $\sigma_{Z^n}(t) = \sigma_Z(t, T_{n+1}) - \sigma_Z(t, T_n)$, where $\sigma_Z(t, T)$ are defined in (15.46). The drifts of forward bond ratios, not surprisingly, are not free parameters. We have the following result on the dynamics of Libor rates and forward bond ratios.

Proposition 15.5.4. *If the Libor rates $\tilde{L}_n(t)$ satisfy (15.48), and forward bond ratios $Z^n(t)$ follow the dynamics of (15.52), then in the spot measure Q^B , the process for the state variables of the model is given by*

$$d\tilde{L}_n(t) = \tilde{\sigma}_n(t)^\top (\tilde{\mu}_n(t) dt + dW^B(t)), \quad (15.53)$$

$$\begin{aligned} dZ^n(t)/Z^n(t) &= \sigma_{Z^n}(t)^\top \\ &\times \left(\left(\frac{\tau_n \tilde{\sigma}_n(t)}{1 + \tau_n \tilde{L}_n(t)} + \sigma_{Z^n}(t) - \tilde{\mu}_n(t) \right) dt - dW^B(t) \right), \end{aligned} \quad (15.54)$$

where

$$\tilde{\mu}_n(t) = \sum_{j=q(t)}^n \left(\frac{\tau_j \tilde{\sigma}_j(t)}{1 + \tau_j \tilde{L}_j(t)} + \sigma_{Z^j}(t) \right)$$

and $W^B(t)$ is an m -dimensional Brownian motion in measure Q^B .

Proof. Clearly

$$\frac{1}{1 + \tau_j L_j(t)} = \frac{1}{1 + \tau_j \tilde{L}_j(t)} Z^j(t). \quad (15.55)$$

Differentiating (15.55) and matching the dW terms, we obtain that

$$\frac{\tau_j \sigma_j(t)}{1 + \tau_j L_j(t)} = \frac{\tau_j \tilde{\sigma}_j(t)}{1 + \tau_j \tilde{L}_j(t)} + \sigma_{Z^j}(t),$$

and (15.53) follows.

To derive the drift in (15.54), we note that $(1 + \tau_n L_n(t))Z^n(t) = 1 + \tau_n \tilde{L}_n(t)$ is a Q^{n+1} -martingale. Hence, in T^{n+1} -forward measure we have the following dynamics,

$$dZ^n(t)/Z^n(t) = \sigma_{Z^n}(t)^\top \left(\frac{\tau_n \sigma_n(t)}{1 + \tau_n L_n(t)} dt - dW^{n+1}(t) \right).$$

Switching to the spot measure, we obtain (15.54). \square

The simulation scheme for the two-curve LM model is similar to the one-curve case, with obvious modifications. The initial values of the model primitives, the Libor rates $\tilde{L}_n(0)$ and the forward bond ratios $Z^n(t)$, $n = 0, \dots, N - 1$, are derived from the initial values of the discount and index curves. Having specified the Libor rate volatilities $\tilde{\sigma}_n(t)$ and forward bond spread volatilities $\sigma_{Z^n}(t)$, $n = 1, \dots, N - 1$, we simulate the Libor rates and the bond spreads using SDEs from Proposition 15.5.4. At each time t , we may construct an index curve $\tilde{P}(t, T)$, $T \geq t$, from the simulated Libor rates $\tilde{L}_n(t)$, $n = 0, \dots, N - 1$, as in the one-curve case in Chapter 14. From the simulated index curve and the simulated bond spread curve $Z(t, T)$, $T \geq t$, the discount curve $P(t, T)$, $T \geq t$, is calculated via (15.40). The discount curve is used to update the numeraire by (15.49) and, together with the index curve, to calculate market rates needed for evaluating derivative payoffs. In particular, forward swap rates are calculated by projecting Libor rates off the index curve and discounting them off the discount curve (compare to (6.39)), so a rate fixing at T_j covering $k - j$ periods is calculated by

$$S_{j,k-j}(t) = \frac{\sum_{i=j}^{k-1} \tau_i P(t, T_{i+1}) \tilde{L}_i(t)}{\sum_{i=j}^{k-1} \tau_i P(t, T_{i+1})}. \quad (15.56)$$

In the two-curve setup, the derivation of swaption pricing approximations for model calibration can proceed as for the single-curve case. For this purpose, first rewrite (15.56) in terms of model primitives, to obtain

$$S_{j,k-j}(t) = \frac{\sum_{i=j}^{k-1} \tau_i Z(t, T_{i+1}) \tilde{P}(t, T_{i+1}) \tilde{L}_i(t)}{\sum_{i=j}^{k-1} \tau_i Z(t, T_{i+1}) \tilde{P}(t, T_{i+1})}.$$

Proceeding as in the standard single-curve case and ignoring contributions from the $Z(t, T_{i+1})$ terms in the swap rate dynamics — a respectable approximation even with stochastic spreads — we obtain the following dynamics in the appropriate annuity measure

$$dS_{j,k-j}(t) = \sum_{n=1}^{N-1} \left. \frac{\partial S_{j,k-j}(t)}{\partial \tilde{L}_n(t)} \right|_{t=0} \tilde{\sigma}_n(t)^\top dW^{A_{j,k-j}}(t). \quad (15.57)$$

This formula can be used as the starting point for the usual European swaption approximations. Compared to earlier results in Section 14.4.2, the presence of the basis spread leads to slightly altered expressions for the weights $\partial S(t)/\partial \tilde{L}_n(t)|_{t=0}$; the exact form of these weights is left for the reader to derive.

To conclude our discussion of dynamic two-curve modeling in the LM framework, we note that our decision to use forward bond spreads $Z^j(t)$ as variables driving the spread between the index and discount curves is not the only choice. For example, we could have used the spreads between the simply compounded forward rates computed off the index and discount curves,

$\widetilde{L}_j(t) - L_j(t)$, $j = 0, \dots, N - 1$, instead. This would lead to a somewhat different, but equally tractable, extension. As this area of interest rate modeling is still rapidly evolving, consensus on what are the right variables to use has not emerged yet.

15.5.4 Deterministic Spread

It is not hard to see that the values of derivative securities that do not directly reference the spread between the index and the discount curve — which is the majority of them — respond approximately linearly to this spread. For one, the spread is usually confined to a rather tight range, ensuring that a linear approximation may be adequate. Moreover, changes in the spread will normally affect the discount curve substantially more than the index curve, as the latter is predominantly calibrated to linear market instruments such as FRAs and swaps. Values of most securities, even exotic ones, tend to be approximately linear to changes in discounting.

A fairly self-evident rule states that if the dependence of a value of a security on a given market factor is approximately linear, there is little reason to model that factor with a stochastic process for the purposes of valuation and hedging of that security. As a consequence, it is not uncommon to assume deterministic evolution for the spread between the discount and index curves. In the framework we developed, this may be achieved by setting the volatility of the spread $\sigma_\varepsilon(t, T)$ to zero. The two-curve LM model simplifies accordingly. Obviously, one does not need to simulate $\varepsilon(t, T)$ anymore, as

$$\varepsilon(t, T) = \varepsilon(0, T), \quad Z(t, T) = Z(0, T) / Z(0, t)$$

for all t, T such that $0 \leq t \leq T$. Furthermore, the drift term in (15.53) reduces to the standard single-curve expression, and discount factors $P(t, T)$ are obtained from pseudo-discount ones by

$$P(t, T) = \frac{P(0, t, T)}{\widetilde{P}(0, t, T)} \widetilde{P}(t, T), \quad T \geq t.$$

Derivation of swaption pricing expressions proceeds as in Section 15.5.3 where the randomness in the spreads was already ignored in (15.57).

15.6 SV Models with Non-Zero Correlation

Let us return to the SV-type LM model we considered in Section 14.2.5. In the spot measure, we recall that forward rate dynamics are postulated to be of the form

$$dL_n(t) = \sqrt{z(t)}\varphi(L_n(t))\lambda_n(t)^\top \left(\sqrt{z(t)}\mu_n(t)dt + dW^B(t) \right),$$

$$\mu_n(t) = \sum_{j=q(t)}^n \frac{\tau_j \varphi(L_j(t)) \lambda_j(t)}{1 + \tau_j L_j(t)},$$

with

$$dz(t) = \theta(z_0 - z(t))dt + \eta\psi(z(t))dZ(t).$$

See Section 14.2.5 for further details on the notation.

In previous work, we assumed that the scalar Brownian motion $Z(t)$ was independent of all components of the m -dimensional Brownian motion $W(t)$. This assumption allowed us to switch from the spot measure into more convenient forward measures, without altering the form of the process for $z(t)$. We shall now briefly consider relaxing this assumption. In particular, let us assume a non-zero deterministic correlation vector $\rho(t)$ between $Z(t)$ and $W^B(t)$. In this case, Lemma 14.2.6 tells us the dynamics of $z(t)$ in the forward measure $Q^{T_{n+1}}$. As we defined $\langle dW^B(t), dZ(t) \rangle = \rho(t)dt$, then

$$dz(t) = \theta \left(z_0 - z(t) - \frac{\eta}{\theta} \psi(z(t)) \sqrt{z(t)} \rho(t)^\top \mu_n(t) \right) dt + \eta\psi(z(t))dZ^{n+1}(t),$$

where $Z^{n+1}(t)$ is a $Q^{T_{n+1}}$ -Brownian motion. In certain important cases, this expression can be simplified and approximated somewhat, as demonstrated in the following corollary.

Corollary 15.6.1. *When $\psi(z(t)) = \sqrt{z(t)}$, we have*

$$dz(t) = \theta(z_0 - \alpha(t)z(t))dt + \eta\sqrt{z(t)}dZ^{n+1}(t), \quad (15.58)$$

where

$$\alpha(t) = 1 + \frac{\eta}{\theta} \rho(t)^\top \mu_n(t) \approx 1 + \frac{\eta}{\theta} \rho(t)^\top \sum_{j=q(t)}^n \frac{\tau_j \varphi(L_j(0)) \lambda_j(t)}{1 + \tau_j L_j(0)}. \quad (15.59)$$

Notice that in Corollary 15.6.1 we have used an approximation under which the multiplier $\alpha(t)$ in the drift term of (15.58) is approximately deterministic. The approximation is based on the same “freeze along forwards” idea as was used in Section 14.4.2 and elsewhere, and makes the SDE for $z(t)$ affine. As such, from the results in Chapter 8 we would expect that a Fourier-based approach would allow for (approximate) caplet pricing if the dynamics for $L_n(t)$ are themselves affine, e.g. if $\varphi(x) = ax + b$. We verify this below for the log-normal case ($\varphi(x) = x$); for displaced log-normal dynamics the result follows along similar lines.

Proposition 15.6.2. *Define $X_n(t) = \ln L_n(t)$ where*

$$dL_n(t)/L_n(t) = \sqrt{z(t)}\lambda_n(t)^\top dW^{n+1}(t),$$

$$dz(t) = \theta(z_0 - \alpha(t)z(t))dt + \eta\sqrt{z(t)}dZ^{n+1}(t),$$

for a deterministic function $\alpha(t)$ and $\langle dW^{n+1}(t), dZ^{n+1}(t) \rangle = \rho(t) dt$. The moment-generating function

$$\Psi_{X_n}(u) = \mathbb{E}^{T_{n+1}} \left(e^{u X_n(T_n)} \right)$$

can be written as

$$\Psi_{X_n}(u) = e^{A(0,T)+X_n(0)B(0,T)+z(0)C(0,T)},$$

where A , B , and C solve a system of Riccati ODEs

$$0 = A'(t, T) + \theta z_0 C(t, T),$$

$$0 = B'(t, T),$$

$$\begin{aligned} 0 = C'(t, T) - \frac{1}{2} \|\lambda_n(t)\|^2 B(t, T) + \frac{1}{2} \|\lambda_n(t)\|^2 B(t, T)^2 \\ - \theta \alpha(t) C(t, T) + \frac{1}{2} \eta^2 C(t, T)^2 + \eta \lambda_n(t)^\top \rho(t) B(t, T) C(t, T), \end{aligned}$$

with the terminal conditions

$$A(T, T) = 0, \quad B(T, T) = u, \quad C(T, T) = 0.$$

Proof. Positing the form

$$\mathbb{E}^{T_{n+1}} \left(e^{u X_n(T_n)} \middle| X_n(t) = x, z(t) = v \right) = e^{A(t, T) + x B(t, T) + v C(t, T)}$$

and substituting it into the Feynman-Kac PDE that corresponds to the dynamics of $(X_n(t), z_n(t))$, we obtain the following equation

$$\begin{aligned} 0 = A'(t, T) + x B'(t, T) + v C'(t, T) \\ - \frac{1}{2} v \|\lambda_n(t)\|^2 B(t, T) + \frac{1}{2} v \|\lambda_n(t)\|^2 B(t, T)^2 \\ + \theta (z_0 - \alpha(t)v) C(t, T) + \frac{1}{2} \eta^2 v C(t, T)^2 \\ + v \eta \lambda_n(t)^\top \rho(t) B(t, T) C(t, T). \end{aligned}$$

Combining the terms in x and v together, the result follows. \square

The result of the proposition above can be combined with Theorem 8.4.1 to allow for analytical pricing of caplets by Fourier methods. Pricing of swaptions follows a similar line of attack. To start, we write down the drift of $z(t)$ under the corresponding annuity measure; while somewhat more complicated in appearance, it will still be linear in z after application of the “freezing” technique. Then, the moment-generating function for the logarithm of the swap rate is available, allowing for analytic pricing of swaptions via Fourier methods. We trust the reader to fill in missing details.

15.7 Multi-Stochastic Volatility Extensions

15.7.1 Introduction

In the specification (15.22)–(15.23) of the LM model, a single stochastic volatility process $\sqrt{z(t)}$ is used to scale the diffusion coefficients of all forward rates. As such, the volatility structure of the model is only allowed (nearly) parallel moves up and down. While sufficiently rich to introduce the volatility smile for all European swaptions, one has to wonder about the limitations of this one-factor specification.

The value of many exotic interest rate derivatives, sometimes called “first generation” exotics, is primarily linked to the overall *level* of the interest rate curve. This class includes, for example, Bermudan swaptions, callable inverse floaters, or callable range accruals on a Libor or CMS rate. For such instruments, having a single stochastic volatility factor applied to all rates is typically adequate.

On the other hand, interest rate exotics linked to the spread of two CMS rates, such as CMS spread callable swaps (see Section 5.13.3) or CMS spread TARNs (see Section 5.15.2) derive their value from the distribution of the *slope* of the interest rate curve. Just like a single-factor model is unsuitable for pricing such exotics — being unable to represent the changes in the slope of the curve — a common stochastic volatility factor applied to all rates does not always allow for sufficient control over the distribution of the slope of the interest rate curve. In particular, such a specification does not typically allow for much control over the finer features of the *volatility smile of the spread*, e.g. its slope or curvature.

We will have quite a bit more to say about modeling the smile of a CMS spread later on in Chapter 17, but let us take these observations as a rationale (or excuse) for sketching an extension of the LM model that allows for some de-correlation in stochastic volatility factors applied to different rates. While many roads could be taken, the route we shall suggest here

- Incorporates the features necessary for realistic CMS spread modeling.
- Remains relatively parsimonious.
- Can still be calibrated using analytical approximations to prices of caps and swaptions.

We must point out that the whole area of multi-dimensional stochastic volatility interest rate modeling is quite new, so we keep the discussion suitably brief, with details to be filled by future research.

15.7.2 Setup

One can take a view, admittedly not inconsistent with the philosophy of LM modeling, that each Libor rate should be driven by its own stochastic variance process. However, it should be obvious that such specification

would be quite unwieldy, and would likely not lend itself easily to closed-form approximations for swaption prices. With the point of our proposed extension focused primarily on controlling the smiles of CMS spreads, we instead consider a parsimonious extension that involves only *two* stochastic variance processes.

We define

$$dz^i(t) = \theta^i (z_0^i - z^i(t)) dt + \eta^i \sqrt{z^i(t)} dZ^i(t), \quad i = 1, 2. \quad (15.60)$$

Note the time-independence of parameters; while a more general specification is certainly possible, we keep our focus on more important details. Moreover, we require

$$\langle dZ^1(t), dZ^2(t) \rangle = 0.$$

For the applications we have in mind, we need to allow for non-zero correlation between Brownian motions driving the Libor rates and those driving the stochastic variances. Hence, (15.60) is understood to hold under Q^B , the spot Libor measure only. Under the same measure, we assume that the Libor rates follow

$$\begin{aligned} dL_n(t)/\varphi(L_n(t)) &= \sqrt{z^1(t)} \lambda_n^1(t)^\top (dW^1(t) + \mu_n^1(t) dt) \\ &\quad + \sqrt{z^2(t)} \lambda_n^2(t)^\top (dW^2(t) + \mu_n^2(t) dt), \quad n = 1, \dots, N-1. \end{aligned} \quad (15.61)$$

Here, $W^1(t)$ and $W^2(t)$ are two independent copies of a d -factor Brownian motion. Moreover, we assume that

$$\langle dW^i(t), dZ^i(t) \rangle = \chi^i, \quad i = 1, 2, \quad (15.62)$$

but also

$$\langle dW^1(t), dZ^2(t) \rangle = \langle dW^2(t), dZ^1(t) \rangle = 0. \quad (15.63)$$

Under the T_{n+1} -forward measure, the Libor rate $L_n(t)$ is a martingale,

$$\begin{aligned} dL_n(t)/\varphi(L_n(t)) &= \sqrt{z^1(t)} \lambda_n^1(t)^\top dW^{1,T_{n+1}}(t) \\ &\quad + \sqrt{z^2(t)} \lambda_n^2(t)^\top dW^{2,T_{n+1}}(t), \quad n = 1, \dots, N-1. \end{aligned} \quad (15.64)$$

15.7.3 Pricing Caplets and Swaptions

Following the same techniques used in Section 15.6 above, it is straightforward to show that in the T_{n+1} -forward measure, the dynamics of the stochastic variance processes are

$$dz^i(t) = \theta^i (z_0^i - z^i(t)) dt - \eta^i \nu^{i,n+1}(t, \mathbf{L}(t)) z^i(t) dt + \eta^i \sqrt{z^i(t)} dZ^{i,T_{n+1}}(t),$$

$i = 1, 2$, where

$$\nu^{i,n+1}(t, \mathbf{L}(t)) = (\chi^i)^\top \sum_{j=q(t)}^n \frac{\tau_j \varphi(L_j(t)) \lambda_j^i(t)}{1 + \tau_j L_j(t)}.$$

Freezing the drifts, in the same manner as in Corollary 15.6.1, we get

$$dz^i(t) = \theta^i(z_0^i - z^i(t)) dt - \eta^i \nu^{i,n+1}(t, \mathbf{L}(0)) z^i(t) dt + \eta^i \sqrt{z^i(t)} dZ^{i,T_{n+1}}(t), \quad (15.65)$$

and we obtain that (15.64), (15.65) constitute an affine specification (thanks to (15.63)). Hence, the moment-generating function could be represented in a quasi-closed form of an exponential of coefficients computable from Riccati equations, and prices of caplets could be obtained by Fourier inversion of the moment-generating function. We omit straightforward details.

As far as swaptions are concerned, as is the case in many other LM model specifications, we could derive the dynamics of a swap rate of essentially the same form as that of the Libor rates. Then, the same arguments as above can be applied to compute European swaption prices.

15.7.4 Spread Options

The intention of introducing two stochastic variance processes above is to be able to control the volatility smile⁵ of the spread option, while keeping the smiles of individual swap rates fixed and in calibration with observed swaption values. As shown later in Chapter 17, such control can be achieved if we have a mechanism to affect i) the correlation of the stochastic variance processes affecting the two rates in the spread; and ii) the “cross” correlations between a forward rate and the variance process of the other forward rate in the spread. The specification outlined above allows for this, as we shall now demonstrate.

Consider two forward Libor rates, $L_n(t)$ and $L_m(t)$, $n \neq m$. For simplicity we assume that in (15.61) $\varphi(x) \equiv 1$ and $\lambda_k^i(t) \equiv \lambda_k^i$, $i = 1, 2$, $k = n, m$. Focusing on the diffusion terms only, we have (drifts and probability measure are irrelevant),

$$dL_k(t) = O(dt) + \sqrt{z^1(t)} (\lambda_k^1)^\top dW^1(t) + \sqrt{z^2(t)} (\lambda_k^2)^\top dW^2(t), \quad k = n, m, \quad (15.66)$$

$$dz^i(t) = O(dt) + \eta^i \sqrt{z^i(t)} dZ^i(t), \quad i = 1, 2. \quad (15.67)$$

In the following easily proven result, we rewrite the dynamics in the form convenient for correlation analysis.

Proposition 15.7.1. *Assume that L_n and L_m satisfy (15.66)–(15.67) above, with the correlation structure as in (15.62) and (15.63). The joint dynamics*

⁵One way to define such a smile is in terms of the implied Bachelier volatility of the spread itself; see Sections 14.4.3 and 17.4.1 for more details.

of the two Libor rates and their stochastic variance processes can then be written as

$$dL_k(t) = O(dt) + \sqrt{u^k(t)} dU^k(t), \quad du^k(t) = O(dt) + \sqrt{\eta^k(t)} dX^k(t), \quad k = n, m,$$

where

$$\begin{aligned} u^k(t) &= z^1(t) \|\lambda_k^1\|^2 + z^2(t) \|\lambda_k^2\|^2, \\ \eta^k(t) &= z^1(t) \|\lambda_k^1\|^4 (\eta^1)^2 + z^2(t) \|\lambda_k^2\|^4 (\eta^2)^2, \end{aligned}$$

and

$$\begin{aligned} dU^k(t) &= \frac{1}{\sqrt{u^k(t)}} \left(\sqrt{z^1(t)} (\lambda_k^1)^\top dW^1(t) + \sqrt{z^2(t)} (\lambda_k^2)^\top dW^2(t) \right), \\ dX^k(t) &= \frac{1}{\sqrt{\eta^k(t)}} \left(\sqrt{z^1(t)} \|\lambda_k^1\|^2 \eta^1 dZ^1(t) + \sqrt{z^2(t)} \|\lambda_k^2\|^2 \eta^2 dZ^2(t) \right). \end{aligned}$$

It follows from Proposition 15.7.1 that our model setup allows for essentially independent control of the correlations between the variance processes u^n and u^m , as well as the correlations between L_m and u^n and L_n and u^m as should be clear from the expression for one of the correlations (others are similar):

$$\begin{aligned} \text{Corr}(dL^n(t), du^m(t)) &= \frac{1}{\sqrt{u^n(t)\eta^m(t)}} \\ &\times (z^1(t) \|\lambda_m^1\|^2 \eta^1 (\lambda_n^1)^\top \chi^1 + z^2(t) \|\lambda_m^2\|^2 \eta^2 (\lambda_n^2)^\top \chi^2). \end{aligned}$$

The results extend to swap rate spreads in a predictable, and largely mechanical, fashion, allowing us to set up an LM calibration that targets quantities linked to the shape of smile of various CMS spread options of interest. As the reader might expect, the formulas become rather cumbersome and we do not list them here.

At this point we have gathered enough results for our later discussion on spread option pricing, so we conclude the analysis here. We return to spread options in stochastic volatility models in Chapter 17.

15.7.5 Another Use of Multi-Dimensional Stochastic Volatility

In models for foreign exchange or for equity prices, multi-dimensional stochastic volatility is typically used not as a way to refine spread option pricing, but rather as i) a mechanism to induce multiple⁶ time-scales in the mean-reverting behavior of volatility; or ii) to control the evolution of the volatility smile through time. A discussion of multiple time-scales in empirical data

⁶Basically this means that κ_1 is either much larger or much smaller than κ_2 .

can be found in Perello et al. [2004], and Kainth and Saravanamuttu [2007] (among others) discuss applications specific to foreign exchange, with an emphasis on smile dynamics. We also note that Andersen and Brotherton-Ratcliffe [2005] introduce a tractable alternative to (15.61) that has perfect correlation between the variances of all Libor forwards, but still introduces multiple time-scales in the variance dynamics; this setup (the details of which we omit) is mainly useful for applications where the correlation between variances of different forwards is thought to be high.

Nevertheless, it is important to realize that different models, while possibly producing identical swaption prices, may imply different — sometimes *very* different — hedging strategies and, ultimately, P&L (*Profit-And-Loss*) of a vanilla options desk. We elaborate on this in the next two sections, and then turn to a few relevant topics associated with practicalities of swaption model calibration.

16.1.1 Smile Dynamics

Delta hedging in vanilla models is intimately linked to the model-implied *volatility smile dynamics*, i.e. how the smile moves with the underlying. To expand on this, let us focus on a T -maturity, K -strike European call option $c(t) = c(t, S(t); T, K)$, and consider the computation of its time t delta

$$\Delta(t) = \frac{\partial c(t)}{\partial S}.$$

Recall (see (7.6)) the notion of implied volatility $\sigma_B(t, S(t); T, K)$,

$$c(t) = c_B(t, S(t); T, K; \sigma_B(t, S(t); T, K)),$$

where $c_B(t, S; T, K; \sigma)$ represents the usual Black formula at a volatility level σ (see Remark 7.2.8). By the chain rule, it follows that

$$\Delta(t) = \frac{\partial c_B}{\partial S} + \frac{\partial c_B}{\partial \sigma_B} \frac{\partial \sigma_B}{\partial S} = \Delta_B(t) + \gamma_B(t) \frac{\partial \sigma_B}{\partial S}, \quad (16.1)$$

where $\Delta_B(t)$ and $\gamma_B(t)$ are, respectively, the delta and vega³ computed in a Black model at a volatility of σ_B . From the Black formula (7.6), it follows that

$$\Delta_B(t) = \Phi(d_+), \quad \gamma_B(t) = S(t)\sqrt{T-t}\phi(d_+),$$

with d_{\pm} defined by (7.6). According to (16.1), for models more sophisticated than the Black model, one can expect that a certain (possibly negative) amount of Black vega will “leak” into the Black delta, to produce a proper model-consistent delta. The amount of this leakage is controlled by the smile dynamics of the model, as given by the term $\partial \sigma_B / \partial S$. This term is often called the *backbone* of the model.

For arguments sake, suppose we use the Black model with strike-specific volatility for risk-managing options, i.e. we use the pricing formula

$$c(t, S(t); T, K) = c_B(t, S(t); T, K; \psi(T, K)) \quad (16.2)$$

where $\psi(T, K)$ is calibrated, at time $t = 0$, to market for each T and K . According to (16.1), the delta calculated in our model will be exactly $\Delta_B(t)$, i.e. the ordinary Black delta. On the other hand, this would not be the

³Vega is the volatility sensitivity, see Remark 8.9.3.

case in, say, the Heston model, which we recall (from Section 8.8) to have “sticky delta” dynamics when the variance variable $z(t)$ is kept fixed. For the common case of a downward-sloping smile (i.e. a negative correlation parameter in the Heston model), it follows that $\Delta_{\text{Heston}}(t) > \Delta_B(t)$, since $\Upsilon_B(t) > 0$ and — according to Figure 8.8 in Section 8.8 — $\partial\sigma_B/\partial S > 0$. For a downward sloping smile in a local volatility model, on the other hand, $\partial\sigma_B/\partial S < 0$ (see Figure 8.7) and the model delta would be *less* than the Black delta.

In general, the usefulness of a given swaption model to a trading desk lies perhaps not so much in its ability to fit the market — an easy feat if model parameters are allowed to depend on strike — but by how closely the model can match the realized dynamics of the volatility smile. Indeed, to the extent that volatility smile moves predicted by the model differ markedly from observations, the hedging strategies prescribed by the model will not be successful in practice, and traders will not be able to predict the P&L implications of a move in market variables. Due to the importance of the backbone, it is not uncommon for a trading desk to exogenously supply an ad-hoc rule for smile moves that overrides the model-computed value of $\partial\sigma_B/\partial S$; this practice is sometimes known as *shadow delta hedging*. Common rules include the earlier mentioned sticky delta rule, as well as the *sticky strike* rule (16.2) which assumes that the smile remains fixed as a function K when S moves⁴. While applications of ad-hoc rules when computing deltas will compromise the theoretical integrity of the underlying model, in practice the efficacy and stability of the hedging strategy may nevertheless improve.

16.1.2 Adjustable Backbone

As shadow delta hedging is a very common practice, let us proceed to elaborate on the basic idea, by describing a possible approach for exogenously controlling the backbone in a more nuanced way than through simple sticky delta/strike rules. We present our ideas in the context of a displaced log-normal model; adding stochastic volatility to the model would follow standard procedures and will not affect the backbone. Recalling the model (7.21),

$$dS(t) = \lambda(bS(t) + (1 - b)L) dW(t), \quad (16.3)$$

we first focus on calculating its backbone. To find approximately the implied Black volatility for the model, we apply the expansion method of Proposition 7.5.1 with $\beta = 0$, $\zeta = 1$. Keeping only the first term, we obtain

$$\sigma_B(0, S(0); T, K) \approx \lambda \frac{\ln(S(0)/K)}{\int_K^{S(0)} \frac{du}{bu + (1-b)L}} = \lambda b \frac{\ln(S(0)/K)}{\ln\left(\frac{bS(0)+(1-b)L}{bK+(1-b)L}\right)}.$$

The backbone for the at-the-money strike is given by

⁴It can be shown that the sticky strike rule is arbitrageable.

$$\left. \frac{\partial \sigma_B(0, S(0); T, K)}{\partial S(0)} \right|_{K=S(0)} \approx -\frac{1-b}{2} \frac{L}{S(0)^2} \lambda$$

and with $L \approx S(0)$, as is typically the case, we get

$$\left. \frac{\partial \sigma_B(0, S(0); T, K)}{\partial S(0)} \right|_{K=S(0)} \approx -\frac{1-b}{2} \frac{\lambda}{S(0)}.$$

Clearly, the backbone is controlled by b . When $b = 1$, we obtain the Black backbone ($\partial \sigma_B / \partial S = 0$) and, when $b = 0$, we obtain what we call the *Gaussian* backbone, as it is the backbone that is consistent with Gaussian dynamics — see Remark 7.2.9. The model (16.3), however, does not allow for an independent control over the backbone as b is not a free parameter, but is determined by the slope of the market-observed volatility smile.

Let us see what would happen if we had specified the model (16.3) somewhat differently:

$$dS(t) = \lambda(bS(t) + (1-b)S(0)) dW(t). \quad (16.4)$$

Then, following the same steps as above, we would get

$$\begin{aligned} \sigma_B(0, S(0); T, K) &\approx \lambda b \frac{\ln(S(0)/K)}{\ln\left(\frac{S(0)}{bK+(1-b)S(0)}\right)}, \\ \left. \frac{\partial \sigma_B(0, S(0); T, K)}{\partial S(0)} \right|_{K=S(0)} &\approx \frac{1-b}{2} \frac{\lambda}{S(0)}. \end{aligned}$$

We see that the backbone is now different — *positive* (for $b \in [0, 1]$) rather than *negative*, as in the model (16.3). This difference, of course, originates with the fact that a perturbation to $S(0)$ now affects the local volatility function in (16.4): a shock of size $\delta S(0)$ to $S(0)$ increases the local volatility function by $\lambda(1-b)\delta S(0)$, and the impact propagates into the implied volatility itself. On the other hand, the volatility smile generated by the model (16.4) has the same slope as the smile generated by the model (16.3), so we have arrived at two models with the same (static) smile but different smile dynamics. Clearly, by “mixing” the two, we can get a model where the backbone is controlled independently of the smile.

Here is how we proceed. Introducing a new parameter, “mixing” m , we specify

$$dS(t) = \lambda(bS(t) + (m-b)S(0) + (1-m)L) dW(t). \quad (16.5)$$

For L close to $S(0)$, λ still has the meaning of relative (log-normal) volatility, and the slope of the smile is still controlled by b . On the other hand, simple calculations yield

$$\sigma_B(0, S(0); T, K) \approx \lambda b \frac{\ln(S(0)/K)}{\ln\left(\frac{mS(0)+(1-m)L}{bK+(m-b)S(0)+(1-m)L}\right)},$$

$$\frac{\partial \sigma_B(0, S(0); T, K)}{\partial S(0)} \Big|_{K=S(0)} \approx \frac{1}{2} \frac{(m-b)S(0) - (1-m)L}{S(0)} \frac{\lambda}{S(0)}$$

and, with $L \approx S(0)$,

$$\frac{\partial \sigma_B(0, S(0); T, K)}{\partial S(0)} \Big|_{K=S(0)} \approx \left(m - \frac{1+b}{2}\right) \frac{\lambda}{S(0)}.$$

Clearly, for any b , we can adjust the backbone by choosing a suitable m . For example, we can always obtain the Black backbone by setting $m = (1+b)/2$, or the Gaussian backbone by setting $m = b/2$.

To check that m does not, indeed, have an impact on the (static) smile, let us calculate the slope of the implied volatility smile in the model (16.5). We have

$$\frac{\partial \sigma_B(0, S(0); T, K)}{\partial K} \Big|_{K=S(0)} \approx \frac{1}{2} \frac{-mS(0) + Lm + bS(0) - L}{S(0)} \frac{\lambda}{S(0)}$$

and, with $L \approx S(0)$,

$$\frac{\partial \sigma_B(0, S(0); T, K)}{\partial K} \Big|_{K=S(0)} \approx -\frac{1-b}{2} \frac{\lambda}{S(0)}. \quad (16.6)$$

Thus, as claimed, the slope of the smile is independent of m and, in particular, is the same for the models (16.3), (16.4) and (16.5).

We should note that sometimes a different definition of the backbone is used, a definition that we call the *ATM backbone*:

$$\frac{\partial \sigma_B(0, S; T, S)}{\partial S} \Big|_{S=S(0)}.$$

This quantity specifies how the at-the-money volatility $\sigma_B(0, S(0); T, S(0))$ changes with the underlying $S(0)$. Simple calculus yields

$$\begin{aligned} \frac{\partial \sigma_B(0, S; T, S)}{\partial S} \Big|_{S=S(0)} &= \frac{\partial \sigma_B(0, S(0); T, K)}{\partial S(0)} \Big|_{K=S(0)} \\ &\quad + \frac{\partial \sigma_B(0, S(0); T, K)}{\partial K} \Big|_{K=S(0)} \end{aligned}$$

and we obtain, for the model (16.5) assuming $L \approx S(0)$,

$$\frac{\partial \sigma_B(0, S; T, S)}{\partial S} \Big|_{S=S(0)} \approx -(1-m) \frac{\lambda}{S(0)}. \quad (16.7)$$

The ATM backbone in the model (16.5) is independent of the skew b . When $m = b$ (model (16.3)) the ATM backbone is twice the slope of the volatility smile (see (16.6)), and when $m = 1$ (model (16.4)), the ATM backbone is zero, i.e. at-the-money implied volatility does not change as the underlying moves. Other regimes are easy to simulate. For example, a trader may believe that the at-the-money implied volatility should “slide along” the smile, i.e. exhibit (a weaker form of) the “sticky strike” behavior. Mathematically, this is expressed as

$$\frac{\partial \sigma_B(0, S; T, S)}{\partial S} \Big|_{S=S(0)} = \frac{\partial \sigma_B(0, S(0); T, K)}{\partial K} \Big|_{K=S(0)}$$

which, using (16.6) and (16.7), gives us the following condition on the mixing parameter:

$$m = \frac{b+1}{2}.$$

16.1.3 Stochastic Volatility Swaption Grid

With backbone issues out of the way, let us now discuss a typical setup for vanilla options modeling. As mentioned before, a stochastic volatility model offers a good compromise between tractability and the ability to represent typical shapes of volatility smiles. To describe a typical setup, let us introduce a tenor structure

$$0 < T_0 < T_1 < T_2 < \dots < T_N, \quad \tau_n = T_{n+1} - T_n,$$

and a collection of forward swap rates of different expiries/tenors, as in Section 5.10, see (5.13)–(5.14). For each swap rate $S_{n,m}(t)$, we specify the following SV-style dynamics (see Chapter 8) in the corresponding annuity measure

$$dS_{n,m}(t) = \lambda_{n,m} (b_{n,m} S_{n,m}(t) + (1 - b_{n,m}) S_{n,m}(0)) \sqrt{z_{n,m}(t)} dW^{n,m}(t), \quad (16.8)$$

$$dz_{n,m}(t) = \theta (1 - z_{n,m}(t)) dt + \eta_{n,m} \sqrt{z_{n,m}(t)} dZ^{n,m}(t), \quad (16.9)$$

with $\langle dZ^{n,m}(t), dW^{n,m}(t) \rangle = 0$. Alternative local volatility parameterizations as in (16.3) or even (16.5) could of course be used, but we abstain from doing so to simplify notations. Further, we assume that the mean reversion of variance parameter θ is *global*, i.e. the same for all swaptions. This does not in any way restrict the range of available smiles for each individual swap rate as explained in Section 8.2, yet allows for a measure of consistency in term structure models (e.g., as in Section 13.2) that we, eventually, calibrate to the vanilla market. The rest of the parameters

$$\{(\lambda_{n,m}, b_{n,m}, \eta_{n,m})\}, \quad n = 0, \dots, N-1, \quad m = 1, \dots, N-n,$$

form the so-called *SV swaption grid*, with the meanings of various parameters explained in Section 8.2. Relative to a full swaption volatility cube (see footnote 2), an SV swaption grid typically requires storage of fewer parameters, as the SV model produces a fairly parsimonious (and guaranteed arbitrage-free) interpolation rule in the strike dimension, eliminating the need for outright storage of implied Black or Gaussian volatilities on a strike grid. Of course, multiple other — possibly heuristic — interpolation rules can be used instead. We return to this briefly in Section 16.1.5 below.

16.1.4 Calibrating Stochastic Volatility Model to Swaptions

SV parameters for swaptions are usually obtained by individually calibrating swaptions for each expiry/maturity grid point, with an exception of the mean reversion parameter θ . Recall (Section 8.2) that the parameter θ controls the speed at which the volatility smile flattens with time to expiry, so it is possible to choose a single θ for the whole grid, in such a way as to minimize the variability of $\eta_{n,m}$ across different n 's (expiries).

Selection of θ can be done manually by choosing a particular θ , calibrating SV parameters to each swaption grid point, and then assessing how constant $\eta_{n,m}$ for different n, m are. If not sufficiently constant, a different θ can be selected, and the procedure iterated until a sufficiently good choice is found. In general, if $\eta_{n,m}$'s increase with n , we need a smaller θ to prevent volatility smiles from flattening out too fast with expiry. Conversely, if $\eta_{n,m}$'s decrease in n , a larger θ is needed.

Apart from θ , swaption calibration is performed individually for each grid point. Let us sketch the algorithm. First, we fix a particular swaption maturity and a swap tenor, as represented by indices n, m . Suppressing these indices for the moment, suppose a collection of strikes K_1, \dots, K_J is given, along with corresponding market prices of swaptions $\widehat{V}_1, \dots, \widehat{V}_J$. Given λ, b, η let

$$V(K_j; \lambda, b, \eta), \quad j = 1, \dots, J,$$

be the model prices of swaptions in the model (16.8)–(16.9) with parameters λ, b, η . Our goal is to find λ, b, η to match as closely as possible the market prices $\widehat{V}_1, \dots, \widehat{V}_J$, where nearly always $J \geq 3$. This type of problem is most conveniently solved by non-linear optimization methods. Defining the objective function

$$\mathcal{I}_1(\lambda, b, \eta) = \sum_{j=1}^J w_j \left(V(K_j; \lambda, b, \eta) - \widehat{V}_j \right)^2, \quad (16.10)$$

where w_1, \dots, w_J are user-specified weights, we obtain the calibrated parameters by solving the problem

$$(\lambda^*, b^*, \eta^*) = \operatorname{argmin}_{\{\lambda, b, \eta\}} \mathcal{I}_1(\lambda, b, \eta)$$

with a specialized algorithm such as the Fletcher-Reeves or the Levenberg-Marquardt method (see Press et al. [1992]). As the optimization problem is solved numerically, the solution typically involves multiple calculations of option prices in the SV model, and having an efficient valuation algorithm such as the one developed in Section 8.4 is important for performance.

The weights w_1, \dots, w_J serve two purposes. One is to express the view on which swaptions should be matched more accurately: the higher the weight w_j is, the more closely the algorithm will try to match the price of the swaption with strike K_j . As we typically have more confidence in the at-the-money swaption prices, we would often set the weights higher for at-the-money strikes and lower for strikes away from the ATM. The other important purpose of the weights is to normalize the magnitude of different terms in the sum in (16.10), as different scales of different terms (i.e. some \widehat{V}_j 's are bigger than others) will influence which terms are matched closer. As we often seek to ensure a good fit in terms of implied volatilities rather than absolute option values, a commonly-used scaling involves vegas of the options in the optimization problem; each weight w_j then represents a product of a user-specified importance weight and a scaling weight equal to the inverse of the swaption vega. To simplify user interface, the vega scaling can be internalized, with (16.10) replaced by

$$\mathcal{I}_2(\lambda, b, \eta) = \sum_{j=1}^J w_j \left(\frac{V(K_j; \lambda, b, \eta) - \widehat{V}_j}{\widehat{\Upsilon}_j} \right)^2, \quad (16.11)$$

where $\widehat{\Upsilon}_j$'s are vegas of corresponding options. For efficiency reasons, the vegas should not be calculated inside the calibration loop; a common shortcut is to just use vegas obtained in the Black model. The resulting objective function is a (numerically efficient) approximation to an objective function expressed in terms of implied volatilities:

$$\mathcal{I}_3(\lambda, b, \eta) = \sum_{j=1}^J w_j (\sigma_B(K_j; \lambda, b, \eta) - \widehat{\sigma}_j)^2, \quad (16.12)$$

where $\sigma_B(K_j; \lambda, b, \eta)$ is the Black volatility implied by the model for the option with strike K_j , and $\widehat{\sigma}_j$ is its market-implied volatility. While (16.12) could be used directly, the expense of calculating implied volatilities inside the calibration loop typically makes it less attractive than (16.11).

Finally, let us remind the reader that optimization of a precision norm (e.g. either \mathcal{I}_1 , \mathcal{I}_2 , \mathcal{I}_3) must be undertaken for each pair of swaption expiries and swap tenors in the SV swaption grid, a total of $N(N + 1)/2$ separate optimization problems.

16.1.5 Some Other Interpolation Rules

Usage of the SV model for swaption calibration is particularly convenient if one ultimately needs to use swaption market data for calibration of term structure models such as the quasi-Gaussian model of Section 13.2 or the Libor market model of Section 14.2.5. See, in particular, the discussion in Section 15.2. It is, however, certainly possible to represent the strike-dependence of implied swaption volatilities by different means. A particularly popular choice is to calibrate a SABR model (Section 8.6) to the smile, using the principles outlined above. The existence of a reasonably accurate expansion for implied volatilities in this model makes optimization of the norm (16.12) particularly convenient.

To improve the fitting capability of the model, we note that it is not uncommon for practitioners to “improve” the SABR model with heuristic modifications, such as making the power c or correlation ρ a smooth, bounded function of swaption strikes. Such alterations of the original model make little sense dynamically speaking, but may still represent a valid representation of the marginal distribution of forward swap rates. In a sense, the original SABR model has been used to produce a particular parametric interpolation rule for implied volatilities, where some of the parameters happen to have a convenient intuitive interpretation. Of course, it may then be tempting to skip the entire concept of a dynamic model and simply jump straight to the specification of a smooth parametric form for implied volatilities as a function of strike. There are numerous such forms in circulation; one representative example is the SVI (“stochastic volatility inspired”) 5-parameter form proposed in Gatheral [2004]:

$$\sigma_B(0, S(0); T, K) = a + b \left(\rho(k - h) + \sqrt{(k - h)^2 + s^2} \right), \quad (16.13)$$

where $k \triangleq \ln(K/S(0))$. More details about the valid range and intuition for the parameters a, b, h, ρ, s can be found in Gatheral [2004] and Gatheral and Jacquier [2010], and shall not be repeated here. Let us just note that one drawback of parametric forms is that they can produce arbitrages, in the sense that there typically are parameter combinations that will imply negative marginal densities⁵ for swap rates. This issue must be taken into consideration, e.g. by imposing constraints on the parameter space when calibrating the parametric form against market prices.

We note that similar issues with violation of arbitrage can arise if crude interpolation schemes are used for strike interpolation in a swaption cube — for instance, linear interpolation should never be used, since the second derivative of implied volatilities will not exist everywhere. A better choice for an interpolation scheme would be a twice differentiable spline, such as those described at length in Chapter 6.

⁵These can be computed by differentiating swaption prices twice with respect to the strike, see Section 7.1.2.

16.2 Caps and Floors

While pricing caps and floors (see Section 5.8) is typically no more complicated than pricing swaptions, calibrating a model to quoted cap or floor prices is more involved than calibrating it to swaptions. This is due to the fact that market prices of individual caplets are not directly available, since only prices of full caps — i.e. collections of caplets — are quoted. For example, in the short- to medium-term market in the US, caps of maturities 1, 2, 3, 5 and 10 years are traded. With each caplet covering 3 months, a total of $10 \times 4 = 40$ different caplet maturities are involved, each requiring its own rule for strike interpolation. The scarcity of actual quotes, and the fact that these quotes represent sums of option prices, can potentially lead to overfitting unless extra constraints, either implicit or explicit, are imposed during calibration.

16.2.1 Basic Problem

Assume for a moment that market prices of both 2 year and 2 year 3 months caps are known at multiple strikes. By simple subtraction, we could then recover the prices of individual caplets fixing in two years time, and would be able to fit model parameters to prices of those caplets across strikes, in a manner identical to that employed for swaptions. In reality, there is no market in 2 year 3 months caps, but we could always attempt to obtain the required prices by interpolating between the known prices of 2 year and 3 year caps. Alas, this idea is hampered by the fact that 2 year caps are typically quoted in a range of strikes that is *different* from the strikes for a 3 year cap: for a cap with a given maturity, the quoted strikes are typically fixed offsets from the forward swap rate of the corresponding maturity, i.e. $\text{ATM} \pm 100$ basis points, $\text{ATM} \pm 200$ basis points, and so on. As a consequence, we would need to perform interpolation between 2 year and 3 year cap prices across both expiries *and strikes*. Ensuring that such an interpolation scheme is both free of arbitrage and will give rise to reasonable (i.e. smooth over time) model parameters is not an easy task; our advice is to avoid it.

A more reasonable approach to cap calibration is to employ interpolation directly in model parameters, borrowing the ideas from yield curve construction theory (see Chapter 6). Specifically, we can formalize the cap calibration problem as finding model parameter curves indexed by expiry, such that a price precision norm is minimized subject to penalties for non-smooth model parameters. Encouraging smooth model parameters across expiry makes any subsequent parameter interpolation across time (as required for seasoned trades with fixing schedules deviating from that used in calibration) both more stable and more believable. Moreover, imposing smoothness constraints promotes the stability of calibration through time, a property important for consistency of risk management.

16.2.2 Setup and Norms

To formalize our approach, we specify the tenor structure

$$0 = T_0 < T_1 < T_2 < \dots < T_N, \quad \tau_n = T_{n+1} - T_n,$$

such that $[T_n, T_{n+1}]$ is a caplet tenor (3 months in the US). For concreteness, we use an SV model to define our volatility interpolation scheme in strike space; let the SV parameters to be used for a caplet that fixes at T_n and pays at T_{n+1} be denoted (λ_n, b_n, η_n) , $n = 1, \dots, N - 1$. We denote the price of the n -th caplet with strike K in the SV model with parameters λ, b, η by $V_n(K; \lambda, b, \eta)$. Let $n_i, n_1 < \dots < n_I$, be the number of caplets in the i -th standard market cap. Furthermore, let us suppose that the i -th standard cap is available with strikes

$$K_{i,1}, \dots, K_{i,J},$$

where we for simplicity have assumed that caps of different tenors are quoted for the same *number* of strikes J (but we allow for different *values* of those strikes). Finally, let us denote by $\widehat{V}_{i,j}$, $i = 1, \dots, I$, $j = 1, \dots, J$, the market price of the i -th standard cap (with n_i caplets) at strike $K_{i,j}$.

Let us first consider the introduction of a precision norm that quantifies the amount of mispricing associated with a given set of SV parameters. For instance, we could use a standard weighted least-squares norm

$$\mathcal{I}_1 = \sum_{i=1}^I \sum_{j=1}^J w_{i,j} \left(\sum_{n=1}^{n_i} V_n(K_{i,j}; \lambda_n, b_n, \eta_n) - \widehat{V}_{i,j} \right)^2,$$

where $w_{i,j}$ is the weight associated with the i -th cap of strike K_j .

In principle, we can treat parameter triples for all n , $n = 1, \dots, N - 1$, as independent variables to be recovered in the solution of an optimization problem. However, significant performance improvements can be realized if we reduce the number of free parameters by allowing only the parameters that correspond to the *expiries of market caps* to be free inputs, while interpolating the rest. Linear interpolation seems to perform well, although more sophisticated interpolation schemes borrowed from Chapter 6 could result in further improvements. In any case, if we denote by \mathcal{X} the collection of $\{(\lambda_{n,i}, b_{n,i}, \eta_{n,i})\}$ for $i = 1, \dots, I$, then we can rewrite the objective function as

$$\mathcal{I}_2(\mathcal{X}) = \sum_{i=1}^I \sum_{j=1}^J w_{i,j} \left(\sum_{n=1}^{n_i} V_n(K_{i,j}; \lambda_n(\mathcal{X}), b_n(\mathcal{X}), \eta_n(\mathcal{X})) - \widehat{V}_{i,j} \right)^2,$$

where, as explained, $\lambda_n(\mathcal{X})$, $b_n(\mathcal{X})$, $\eta_n(\mathcal{X})$ are obtained from the elements in \mathcal{X} by suitable interpolation.

Various types of penalties for lack of smoothness are possible, with the discussion of the similar issues in LM calibration in Section 14.5.6 imminently applicable here; a reasonable choice would minimize the discrete equivalent of the integral of the square of the first-order derivative. In particular, we define a norm that, in essence, penalizes deviations of SV parameters from being constant over time,

$$\begin{aligned} \mathcal{I}_{\text{smooth}}(\mathcal{X}) &= w^\lambda \sum_{n=2}^{N-1} (\lambda_n(\mathcal{X}) - \lambda_{n-1}(\mathcal{X}))^2 \\ &+ w^b \sum_{n=2}^{N-1} (b_n(\mathcal{X}) - b_{n-1}(\mathcal{X}))^2 + w^\eta \sum_{n=2}^{N-1} (\eta_n(\mathcal{X}) - \eta_{n-1}(\mathcal{X}))^2. \end{aligned} \quad (16.14)$$

Here weights w^λ , w^b , w^η determine the relative importance of smoothing different model parameters. While the representation (16.14) is quite transparent, we can improve performance somewhat by reformulating — exactly or approximately, depending on the interpolation used — the objective function solely in terms of the components of \mathcal{X} , i.e. by imposing smoothness directly on the “free” parameters $\{(\lambda_{n,i}, b_{n,i}, \eta_{n,i})\}$ for $i = 1, \dots, I$.

16.2.3 Calibration Procedure

Having introduced precision and smoothing norms above, the SV cap calibration problem can be cast as a minimization of the following objective function,

$$\mathcal{X} = \operatorname{argmin} \{w_{\text{precision}} \mathcal{I}_2(\mathcal{X}) + w_{\text{smooth}} \mathcal{I}_{\text{smooth}}(\mathcal{X})\},$$

over the allowed domain for \mathcal{X} . The weights $w_{\text{precision}}$ and w_{smooth} determine relative importance of achieving smoothness over a good fit to market prices. This optimization problem may be solved by numerical methods (see Press et al. [1992]), just as many other calibration problems we discussed previously. To improve efficiency of the algorithm, one could here attempt to split volatility level calibration from smile slope calibration, as dividing optimization problems into smaller ones and tackling them separately often gives us better performance (see relevant discussion in Section 14.5.8). Specifically, significant gains can often be found by iterating over a split scheme where we first calibrate the volatility parameter to at-the-money cap prices only and then calibrate the other model parameters to out-of-the-money cap prices. The success of such a “relaxation” scheme lies with the relative independence of the impacts of volatility parameter and the other parameters: by successfully optimizing in two relatively orthogonal dimensions, we reach a joint minimum faster than in the full calibration. We omit straightforward details.

16.3 Terminal Swap Rate Models

The relative simplicity with which European swaptions (and caps) can be priced comes from the fact that valuation here requires only knowledge of the terminal distribution of a single swap rate, in the appropriate annuity measure. This holds true for all securities whose payoffs can be expressed as deterministic functions of the swap rate $S(T)$ in the annuity measure, as should be clear from the replication argument of Proposition 8.4.13. Unfortunately, such payoffs are relatively rare. Much more common are relatively simple payoffs that appear to depend on the rate $S(T)$ only but, in fact, require the knowledge of certain additional discount bonds, often observed on the same date. As multiple discount bonds are involved and the knowledge of the distribution of a swap rate is not sufficient for valuation, it would appear that a full term structure model is needed to price such derivatives. This, of course, is an option that is always available. However, if the dependence on additional discount bonds is sufficiently mild, we often can avoid computational cost of a full-blown term structure model through certain approximations that aim at functionally linking the values of discount bonds on date T to the “driving” rate $S(T)$, i.e. the rate that primarily determines the payoff. The basic modeling idea, which we denote the *Terminal Swap Rate* (TSR) approach, is extremely useful in handling a range of actively traded European derivatives that are not, strictly speaking, functions of a single rate, but can still be approximated accurately as such. We use the (somewhat loose, admittedly) term *approximately single-rate* for this class of securities; several common securities in the class will be presented in subsequent sections, after the TSR method has been described in detail.

16.3.1 TSR Basics

As briefly outlined in the previous section, the TSR approach treats the swap rate $S(T)$ as the single fundamental state variable for the yield curve at time T . To define the method formally, we continue with the notations of $A(t)$ being the annuity corresponding to the swap rate $S(t)$; for concreteness, we assume that (see (5.4)–(5.5))

$$A(t) \triangleq A_{0,N}(t) = \sum_{n=0}^{N-1} \tau_n P(t, T_{n+1}), \quad (16.15)$$

$$S(t) \triangleq S_{0,N}(t) = \frac{P(t, T) - P(t, T_N)}{A(t)}, \quad (16.16)$$

where

$$0 < T = T_0 < T_1 < \dots < T_N, \quad \tau_n = T_{n+1} - T_n,$$

is a tenor structure of dates. We continue denoting by Q^A the annuity measure, i.e. the measure for which $A(t)$ is the numeraire. We recall that the

market-implied distribution of $S(T)$ in \mathcal{Q}^A can be found from calibrating a vanilla model to N -period swaptions with expiry T , across multiple strikes.

The no-arbitrage valuation formula (1.15) states that the value of a derivative with an \mathcal{F}_T -measurable payoff X is given by

$$V(0) = A(0)\mathbb{E}^A \left(\frac{X}{A(T)} \right). \quad (16.17)$$

Let $\{P(T, M)\}_{M \geq T}$ be the discount bonds of various maturities, all observed at time T . A TSR model specifies a map

$$P(T, M) = \pi(S(T), M), \quad M \geq T, \quad (16.18)$$

where $\{\pi(\cdot, M)\}_{M \geq T}$ is a collection of exogenously specified maturity-indexed functions. In other words, each discount factor is assumed to be a deterministic, known function of the swap rate.

In a proper term structure model, the relationship between the market rate $S(T)$ and the discount factors $\{P(T, M)\}_{M \geq T}$ emerges from the model itself, and is ultimately derived from no-arbitrage conditions. While now we seek to impose the functional relationships (16.18) exogenously, consideration of no-arbitrage must also play a role. Indeed, a first condition to be imposed on the functions $\{\pi(\cdot, M)\}_{M \geq T}$ will be the *no-arbitrage condition*: the valuation formula (16.17), when applied to the specification (16.18), must reproduce initial discount bond prices, i.e. the following must hold for any⁶ $M \geq T$,

$$P(0, M) = A(0)\mathbb{E}^A \left(\frac{\pi(S(T), M)}{\sum_{n=0}^{N-1} \tau_n \pi(S(T), T_{n+1})} \right), \quad M \geq T. \quad (16.19)$$

The no-arbitrage condition by itself is not sufficient to obtain a workable model. Another restriction on the mapping functions is obtained by observing that the swap rate $S(T)$ itself is a function of discount factors, as evidenced by (16.16). This suggests the introduction of a *consistency condition*, i.e. the requirement that the following holds for all x ,

$$x = \frac{1 - \pi(x, T_N)}{\sum_{n=0}^{N-1} \tau_n \pi(x, T_{n+1})}. \quad (16.20)$$

The final condition that we impose on a TSR model is that the set of functions $\{\pi(\cdot, M)\}_{M \geq T}$ should be *reasonable*. While somewhat harder to quantify than the other conditions, we shall mostly impose the following restrictions:

⁶In some applications of TSR models, it may suffice that this expression holds only for a single value of M (namely the payment date of the security in question). In such cases, certain simplifications may be possible, as demonstrated in Section 16.6.4.

- For each x and $M \geq T$, $\pi(x, M)$ is between 0 and 1,

$$0 < \pi(x, M) \leq 1.$$

- For each x , $\pi(x, \cdot)$ is monotonic in M ,

$$M_1 < M_2 \implies \pi(x, M_1) \geq \pi(x, M_2).$$

- The function $\pi(x, M)$ is continuous in (x, M) .

Some of these conditions are more important than others. For example, one may choose to tolerate negative interest rates, i.e. having $\pi(x, M) > 1$ for some x, M , but not negative prices of bonds, i.e. having $\pi(x, M) < 0$ for some x, M .

The conditions listed above do not define the functions $\{\pi(\cdot, M)\}_{M \geq T}$ uniquely; however, they do, as a rule, specify the functions uniquely within a particular parametric class. A concrete model is then obtained by postulating a particular parametric class for the functions $\{\pi(\cdot, M)\}_{M \geq T}$ first, and then choosing functions within the class uniquely from the no-arbitrage and consistency conditions. Let us consider a few representative examples.

16.3.2 Linear TSR Model

The linear TSR model is obtained by specifying

$$\frac{\pi(x, M)}{\sum_{n=0}^{N-1} \tau_n \pi(x, T_{n+1})} = a(M)x + b(M), \quad M \geq T, \quad (16.21)$$

for deterministic functions $a(\cdot)$ and $b(\cdot)$. The no-arbitrage condition requires

$$P(0, M) = A(0)E^A(a(M)S(T) + b(M)),$$

implying a condition on the free coefficient $b(\cdot)$,

$$b(M) = \frac{P(0, M)}{A(0)} - a(M)S(0). \quad (16.22)$$

The consistency condition requires that

$$x = \frac{1 - \pi(x, T_N)}{\sum_{n=0}^{N-1} \tau_n \pi(x, T_{n+1})} = (a(T_0)x + b(T_0)) - (a(T_N)x + b(T_N)),$$

so that

$$b(T_0) = b(T_N), \quad (16.23)$$

$$a(T_0) = 1 + a(T_N). \quad (16.24)$$

It follows from (16.22) that if (16.23) is satisfied, then (16.24) is satisfied as well, and vice versa.

The definition (16.21) imposes new restrictions on $a(\cdot)$, $b(\cdot)$ that go beyond those considered in the previous section. In particular, the following must now hold,

$$\sum_{n=0}^{N-1} \tau_n (a(T_{n+1}) x + b(T_{n+1})) \equiv 1,$$

implying

$$\sum_{n=0}^{N-1} \tau_n a(T_{n+1}) = 0, \quad (16.25)$$

$$\sum_{n=0}^{N-1} \tau_n b(T_{n+1}) = 1. \quad (16.26)$$

It is enough to ensure that one of these two conditions is satisfied; the other will follow automatically by (16.22).

To complete the model specification, we may proceed as follows. First, choose coefficients $\{a(T_1), \dots, a(T_N)\}$, subject to the condition (16.25); a scheme for this will be discussed shortly. Then, calculate $a(T) = a(T_0)$ from (16.24), and the rest of $a(M)$'s by, for example, linear interpolation of $\{a(T), a(T_1), \dots, a(T_N)\}$. Finally, calculate all $b(M)$'s via (16.22).

The specification of a TSR model above enjoys a fair amount of numerical tractability, owing to the simple linear relationship between the market rate and annuity-discounted bonds. For that reason it is rather popular in applications. However, the linear relationship imposed by the model is not wholly realistic, as bond prices may become negative in certain states of the world. Whether this is a problem for a particular application should be decided on a case-by-case basis.

The linear TSR model (16.21) is rather flexible, as the coefficients $\{a(T_1), \dots, a(T_N)\}$ can be selected essentially independently, subject to (16.25) only. Setting these coefficients individually is not particularly convenient, however, as financial implications of various choices are not transparent. As promised above, let us therefore look for a more meaningful way of parameterizing $a(\cdot)$. For this, let us first observe that the coefficients $a(\cdot)$ essentially define the shape of the yield curve at time T for different levels of the “state variable” $S(T)$. Of course, we have seen previously that the same role is played by the mean reversion parameter in the context of a Gaussian term structure model, see Section 10.1.2. This suggests connecting the coefficients of the TSR model to a mean reversion parameter, which would not only reduce the number of parameters we need to specify, but also parametrize the model with a single parameter that has strong financial interpretation and that, in principle, could be derived from prices of traded derivatives (see Section 13.1.8).

To connect $a(\cdot)$ to mean reversion, we interpret the equality (16.21) as defining $a(M)$ via

$$a(M) = \frac{\partial}{\partial S(T)} \frac{P(T, M)}{\sum_{n=0}^{N-1} \tau_n P(T, T_{n+1})},$$

which we rewrite, in the context of a Gaussian one-factor model, as

$$\begin{aligned} a(M) &= \frac{\partial}{\partial x} \frac{P(T, M, x)}{\sum_{n=0}^{N-1} \tau_n P(T, T_{n+1}, x)} \Big|_{S(T,x)=S(0)} \\ &\quad \times \left(\frac{\partial S(T, x)}{\partial x} \Big|_{S(T,x)=S(0)} \right)^{-1}, \end{aligned}$$

where x is now the short rate state in the Gaussian model on which all discount bonds and swap rates depend. We denote by $A(T, x)$ the annuity as the function of the short rate state x ,

$$A(T, x) = \sum_{n=0}^{N-1} \tau_n P(T, T_{n+1}, x),$$

so that

$$S(T, x) = (1 - P(T, T_N, x))/A(T, x)$$

and

$$\begin{aligned} \frac{\partial}{\partial x} \frac{P(T, M, x)}{A(T, x)} &= -\frac{P(T, M, x) G(T, M)}{A(T, x)} - \frac{P(T, M, x)}{A(T, x)^2} \frac{\partial A(T, x)}{\partial x}, \\ \frac{\partial}{\partial x} S(T, x) &= \frac{P(T, T_N, x) G(T, T_N)}{A(T, x)} - \frac{S(T, x)}{A(T, x)} \frac{\partial A(T, x)}{\partial x}, \end{aligned}$$

where $G(\cdot, \cdot)$ is a function of mean reversion (see (10.18)). By using the approximation

$$P(T, t, x)|_{S(T,x)=S(0)} \approx P(0, t)$$

for all $t \geq T$, we obtain

$$a(M) = \frac{P(0, M)(\gamma - G(T, M))}{P(0, T_N)G(T, T_N) + S(0)\gamma}, \quad (16.27)$$

where

$$\begin{aligned} \gamma &\triangleq \frac{1}{A(T, x)} \frac{\partial A(T, x)}{\partial x} \Big|_{P(T,t,x)=P(0,t), \forall t \geq T} \\ &= \frac{\sum \tau_n P(0, T_{n+1}) G(T, T_{n+1})}{\sum \tau_n P(0, T_{n+1})}. \end{aligned} \quad (16.28)$$

As explained before, the coefficients $b(\cdot)$ are obtained from $a(\cdot)$ by (16.22).

With this parameterization, instead of a collection $\{a(T_1), \dots, a(T_N)\}$, only one parameter — the mean reversion κ — needs to be specified. As we shall see later in Section 16.6.8, the choice of mean reversion has a mild but non-vanishing impact on values of many approximately single-rate derivatives.

Linking $a(\cdot)$ to mean reversion leads to a more intuitive parameterization of the model, and also facilitates better risk management. This is so because the mean reversion parameter can in principle be hedged by European swaptions of the same expiry (here T) and different tenors, as we discussed in Section 13.1.8.2. For truly precise vega hedging, however, this somewhat indirect linkage of $a(\cdot)$ to swaption volatilities is less than ideal. In fact, it is not difficult to link $a(\cdot)$ to swap rate volatilities (expiring on the same date T but of different tenors) directly, using an approach similar to what we have developed in this section. We leave the details of this idea for the reader to fill.

16.3.3 Exponential TSR Model

The linear TSR model is a convenient, but by no means unique, representative of the TSR approach. The *exponential TSR model* belongs to the same class, but uses exponential functions to connect a swap rate to discount bonds; this idea originates with the exponential relationship between a discount bond and a corresponding continuously compounded spot yield.

We start the development of the exponential specification by postulating that

$$\pi(x, M) \propto \exp(-l(M)x), \quad M \geq T. \quad (16.29)$$

The intuitive meaning of the loading $l(M)$ is best understood by recalling that a continuously compounded spot yield for the period $[T, M]$, observed at T , is given by

$$y(T, M) = -\frac{1}{M-T} \ln P(T, M) = -\frac{1}{M-T} \ln \pi(S(T), M).$$

Coupled with the specification (16.29), this tells us

$$y(T, M) \propto \frac{l(M)}{M-T} S(T), \quad M \geq T,$$

so that the curve $l(M)/(M-T)$, $M \geq T$, defines the shape of the shock to the yield curve, expressed in terms of yields $\{y(T, M)\}_{M \geq T}$, for a perturbation of the market rate $S(T)$. Recycling the idea of connecting this shape to something similar in a term structure model, we use a one-factor Gaussian model and specify

$$l(M) = \frac{1 - e^{-\kappa(M-T)}}{\kappa}. \quad (16.30)$$

If more precision is desired, we can alternatively write

$$l(M) = \frac{1 - e^{-\kappa(M-T)}}{\kappa} \left(\frac{\partial S(T, x)}{\partial x} \Big|_{S(T, x)=S(0)} \right)^{-1},$$

where the partial derivative is computed in a one-factor Gaussian model, with x being the short rate state and $S(T, x)$ the swap rate as function of the short rate state in the Gaussian model (see Section 10.1.2).

Unfortunately, the expression on the right-hand side of (16.29) cannot be used directly in a TSR model, since it lacks flexibility to satisfy the consistency requirement. This, however, is easily rectified by modifying the functional form slightly and replacing $x \rightarrow \psi(x)$ in (16.29),

$$\pi(x, M) = \exp(-l(M)\psi(x) + b(M)), \quad M \geq T. \quad (16.31)$$

The maturity-dependent function $b(\cdot)$ is obtained from the no-arbitrage conditions (16.19), and the function $\psi(\cdot)$ is defined implicitly by the consistency condition: for any x , $\psi(x)$ is set to be the solution z^* of the equation

$$x = \frac{1 - \exp(-l(T_N)z + b(T_N))}{\sum_{n=0}^{N-1} \tau_n \exp(-l(T_{n+1})z + b(T_{n+1}))}.$$

This equation can easily be solved with just a couple of iterations of a numerical root search algorithm; not surprisingly, it turns out that $\psi(x) \approx x$ to high precision.

16.3.4 Swap-Yield TSR Model

Another example of a TSR model is inspired by the coupon bond yield formula (see Burghardt [2005]). The mapping functions for the *swap-yield TSR model* are defined by

$$\pi(x, M) = \left(\prod_{i=0}^{q(M)-1} \frac{1}{1 + \tau_i x} \right) \times \left(\frac{1}{1 + \tau_{q(M)} x} \right)^{(M - T_{q(M)})/\tau_{q(M)}}, \quad M \geq T, \quad (16.32)$$

where, by (14.2), the index function $q(M)$, $M \geq T$, is specified by the condition

$$M \in [T_{q(M)}, T_{q(M)+1}), \quad (16.33)$$

with the assumption that $T_{N+1} = +\infty$.

The consistency condition (16.20) is satisfied automatically as the following identity holds,

$$\frac{1 - \prod_{i=0}^{N-1} (1 + \tau_i x)^{-1}}{\sum_{n=0}^{N-1} \tau_n \prod_{i=0}^n (1 + \tau_i x)^{-1}} \equiv x. \quad (16.34)$$

The formula (16.32) essentially tells us to discount all cash flows after T at the same rate, namely a rate given by the realized swap rate $S(T)$. As mentioned, this is motivated by traditional definitions of a coupon bond yield or by the payoff of a cash-settled swaption, see Section 5.10.1. The swap-yield TSR specification is motivated by real financial constructs and it can be said to be “reasonable” in the sense of Section 16.3.1; not surprisingly the model has enjoyed widespread popularity in the financial industry. Despite this, the model is not arbitrage-free, as (16.19) is *not* satisfied. Empirically, the extent of violation of no-arbitrage conditions is fairly small, but not always negligible. Another issue with the model, at least in its basic form, is its lack of explicit control over the shape of the yield curve at time T , something we managed to introduce into the linear and exponential specifications by imposing a link between parameters of these models to a mean reversion parameter.

In defense of the swap-yield model, it should be said that, in principle, the model could be improved by few modifications. Violations of no-arbitrage conditions could be addressed by introducing a scaling parameter $b(M)$ as in (16.32). The consistency condition that would no longer hold could be satisfied by introducing a function $\psi(x)$ in the place of x in (16.32), similar to (16.31); this function would need to be calculated numerically. Even ideas about mean reversion could be applied. We invite the reader to attempt these improvements, although we generally feel that the resulting model would offer few, if any, advantages over the linear or exponential TSR models.

16.4 Libor-in-Arrears

We start our study of approximately single-rate derivatives with a closer look at Libor-in-arrears, or LIA (see Section 5.6), cash flows, probably the simplest single-rate derivatives apart from European swaptions and caps. Interestingly, it turns out that LIA cash flow valuation does not require the machinery we developed in Section 16.3, as in fact a LIA cash flow can be stated as a true single-rate product. Our discussion nevertheless shall allow us to give a convenient introduction to issues that are relevant for more complicated products to be covered in later sections.

We recall (Section 5.6) that the defining characteristic of an LIA cash flow is that it pays the Libor rate on the date when the rate fixes, rather than on the date it matures (i.e. the payment date). While most often a whole strip of such cash flows is used as a leg in a Libor-in-arrears swap, we focus our attention on a single cash flow; the valuation of a full strip follows by additivity. Let $T > 0$ be the start date, and M the end date of the period covered by a Libor rate. The forward Libor rate is given, for t such that $0 \leq t \leq T$, by

$$L(t, T, M) = \frac{P(t, T) - P(t, M)}{\tau P(t, M)}, \quad \tau = M - T;$$

we use simplified notation $L(t) = L(t, T, M)$ when there is no chance of confusion. The value, at time 0, of a Libor-in-arrears cash flow is then given by

$$V_{\text{LIA}}(0) = \beta(0) \mathbb{E}(\beta(T)^{-1} L(T)),$$

where $\beta(t)$ is the continuously compounded money market account, and the expected value is taken under the risk-neutral measure \mathbb{Q} . The standard approach to valuing payoffs that pay at time T would involve a switch to the T -forward measure, as then the expression under the expected value operator simplifies accordingly (see Section 4.2.2). Unfortunately, this is not convenient for LIA cash flows as traded caplets provide information about the distribution of the Libor rate in the M -forward measure, *not* the T -forward measure. We shall apply the M -forward measure in a moment; but for now, using the T -forward measure, we obtain

$$V_{\text{LIA}}(0) = P(0, T) \mathbb{E}^T(L(T)).$$

While the expression looks rather simple, our progress along this route is hampered by the fact that $L(t) = L(t, T, M)$ is a martingale in the M -forward measure, not the T -forward measure. Thus,

$$\mathbb{E}^T(L(T)) \neq L(0).$$

To characterize this situation, let us define the concept of a *Libor-in-arrears convexity adjustment*, defined by the difference

$$D_{\text{LIA}}(0) \triangleq \mathbb{E}^T(L(T)) - L(0).$$

This adjustment arises from the mismatch between the measure appropriate for the given payment date and the measure in which the market rate is a martingale. We shall encounter many similar examples later in the chapter; the difference of valuations under different measures is often described generically as *convexity*.

Returning to the issue of valuing an LIA cash flow, we now write the valuation formula in the M -forward measure to obtain that

$$V_{\text{LIA}}(0) = P(0, M) \mathbb{E}^M \left(\frac{1}{P(T, M)} L(T) \right).$$

Fortunately, the factor $\frac{1}{P(T, M)}$ can be rewritten in terms of the Libor rate,

$$\frac{1}{P(T, M)} = 1 + \tau L(T),$$

so that

$$V_{\text{LIA}}(0) = P(0, M) \mathbb{E}^M ((1 + \tau L(T)) L(T)). \quad (16.35)$$

The rate $L(t)$ is a martingale in the M -forward measure, i.e. it has no drift,

$$\mathbb{E}^M(L(T)) = L(0).$$

In particular,

$$V_{\text{LIA}}(0) = P(0, M) (L(0) + \tau \mathbb{E}^M(L(T)^2)). \quad (16.36)$$

The full distribution of $L(T)$ in this measure is encoded in prices of caplets on $L(T)$ with different strikes, see Section 7.1.2. So, to compute

$$\mathbb{E}^M(L(T)^2), \quad (16.37)$$

we merely need to integrate the function x^2 against the probability density of $L(T)$ in measure \mathbb{Q}^M . If one has fitted a particular vanilla caplet model to the market, the density could, in principle, be extracted from this model; in some cases (e.g. the Black or Bachelier models), the density integral can be computed in closed form. In general, however, it is preferable to establish the density directly from observed market prices of T -maturity caplets⁷ on $L(T)$, and to use the replication method of Proposition 8.4.13. Applying the proposition to the problem (16.37), we obtain

$$\begin{aligned} \mathbb{E}^M(L(T)^2) &= L(0)^2 \\ &+ 2 \int_{-\infty}^{L(0)} p(0, L(0); T, K) dK + 2 \int_{L(0)}^{\infty} c(0, L(0); T, K) dK, \end{aligned} \quad (16.38)$$

where $p(t, L; T, K)$ ($c(t, L; T, K)$) are undiscounted values of put (call) options on the rate $L(T)$ with strike K , i.e. simple undiscounted floorlets (caplets). The values of such options are available from the market, and the value of the Libor-in-arrears cash flow is computed by integrating them up.

The power of the replication method goes beyond a mere calculation of the convexity value. As should be clear from the formula above, it also provides a way to hedge the Libor-in-arrears cash flow with standard puts and calls in a *model-independent* way. In particular, to hedge a contract with the payoff

$$(1 + \tau L(T)) L(T),$$

we would

- Enter a short FRA (forward rate agreement, see Section 5.3) on $L(T)$.
- Put $\tau P(0, M) L(0)^2$ dollars into a money market account.
- Sell $2\tau \cdot (dK)$ K -strike puts for all $K \in (-\infty, L(0)]$.
- Sell $2\tau \cdot (dK)$ K -strike calls for all $K \in [L(0), \infty)$.

The hedge is static, i.e. it never needs adjustment throughout the life of the trade. And, as mentioned earlier, the hedge is model-independent, as it does not rely on any modeling assumptions.

⁷Of course establishing caplet prices may itself require some work, as only prices of full caps are quoted. See Section 16.2 for more on this.

To account for the fact that in reality one does not have an infinite number of options on the rate to construct the integrals in (16.38), the integrals can be discretized, and the following formula may be used instead,

$$\mathbb{E}^M(L(T)^2) \approx L(0)^2 + \sum_i w_i^p p(0, L(0); T, K_i) + \sum_i w_i^c c(0, L(0); T, K_i), \quad (16.39)$$

for a collection of strikes $\{K_i\}$, with weights w_i^p and w_i^c chosen so that the sums in (16.39) approximate the integrals in (16.38) at fixing time T . In particular, for a given range $x \in (-x_{\min}, x_{\max})$, the weights can be chosen to *super-replicate* the actual payoff for all values $x \in (-x_{\min}, x_{\max})$ of the realized Libor rate $x = L(T)$:

$$\begin{aligned} & \sum_i \bar{w}_i^p p(T, x; T, K_i) + \sum_i \bar{w}_i^c c(T, x; T, K_i) \\ & \geq 2 \int_{-\infty}^{L(0)} p(T, x; T, K) dK + 2 \int_{L(0)}^{\infty} c(T, x; T, K) dK = x^2 - L(0)^2. \end{aligned}$$

Similarly, we can choose the weights to *sub-replicate*:

$$\begin{aligned} & \sum_i \underline{w}_i^p p(T, x; T, K_i) + \sum_i \underline{w}_i^c c(T, x; T, K_i) \\ & \leq 2 \int_{-\infty}^{L(0)} p(T, x; T, K) dK + 2 \int_{L(0)}^{\infty} c(T, x; T, K) dK = x^2 - L(0)^2. \end{aligned}$$

The minimum value over all super-replicating (maximum over all sub-replicating) portfolios can be regarded as the upper (lower) arbitrage bound on the value of the long (short) LIA cash flow. A value of the LIA cash flow outside of these bounds is arbitrageable with options available in the market in a static, model-independent way.

16.5 Libor-with-Delay

Having considered Libor-in-arrears, we now move on to a more interesting case of Libor cash flows with an arbitrary payment delay. For this product, we can apply the lessons learned in Section 16.4, and also start using the techniques of Section 16.3 for the first time.

Continuing with the notation $L(t) = L(t, T, M)$ for the forward Libor rate covering the period $[T, M]$, we consider a cash flow that pays this rate at some arbitrary payment time T_p , $T_p \geq T$. Switching to the M -forward measure, the measure in which the market-implied distribution of $L(T)$ is known, it follows that the value of the *Libor-with-delay* cash flow is given by

$$V_{LD}(0) = P(0, M) \mathbb{E}^M \left(\frac{P(T, T_p)}{P(T, M)} L(T) \right). \quad (16.40)$$

The presence of the term $P(T, T_p)$ inside the expected value operator now generally prevents us from representing the payoff as a function of the rate $L(T)$ only, making this payoff a simple example of what we defined as an approximately single-rate derivative in Section 16.3.

Valuing Libor-with-delay cash flows presents no theoretical difficulties if one uses a full term structure model, such as the quasi-Gaussian model or a version of the Libor market model. While possible, such a brute-force approach is generally not recommended here. For instance, we note the value of a Libor-in-arrears cash flow and, by extension, of a Libor-with-delay cash flow will depend on values of options of all strikes on a given rate, suggesting that the underlying model should ideally match the entire volatility smile for the underlying Libor rate, something that will be a stretch for a full-blown term structure model. On top of this, there are obvious computational issues in employing a full term structure model for something as vanilla as Libor-with-delay cash flows.

A more sensible approach to the pricing of a Libor-with-delay cash flow utilizes the replication method from Section 16.4 above, along with a method to represent the payoff in (16.40) as a function of the single rate ($L(T)$) only. For the latter, we may use the methods in Section 16.3; we proceed to show an example.

16.5.1 Swap-Yield TSR Model

The simplest and probably most popular method for establishing a functional relationship between the payoff in (16.40) and the rate $L(T)$ is an application of the swap-yield terminal swap rate model of Section 16.3.4. While this method is not fully arbitrage-free (as mentioned), the degree to which no-arbitrage is violated is typically immaterial for Libor-with-delay cash flows. To apply the model to the problem at hand, the rate $L(T) = L(T, T, M)$ is specified to be the market rate, and $P(T, T_p)$ is linked to this rate via the relationship

$$P(T, T_p) = \left(\frac{1}{1 + \tau L(T)} \right)^{(T_p - T)/\tau}.$$

The formula (16.40) then becomes

$$V_{LD}(0) = P(0, M) E^M \left((1 + \tau L(T))^{1 - (T_p - T)/\tau} L(T) \right),$$

allowing us to apply the replication method outlined in Section 16.4 to obtain the value, and the model-independent hedge, of the cash flow. While the resulting formula would not be perfectly arbitrage-free, it will handle correctly the two special cases $T_p = T$ (the formula (16.35) is recovered) and $T_p = M$ (zero convexity adjustment for Libor paid at the end date of the Libor period). Of course, as with any approximation, one should be mindful of pushing the formula beyond its limits; in particular it should not be used for $T_p \gg M$.

16.5.2 Other Terminal Swap Rate Models

As an alternative to the approach in Section 16.5.1, we may apply any of the arbitrage-free terminal swap rate models of Sections 16.3.2 and 16.3.3. As these two TSR models are specifically designed to relate discount factors observed at a particular date T to a market rate observed on the same date, the task of linking $P(T, T_p)$ to $L(T)$ for the Libor-with-delay contract is easily accomplished. We trust the reader can see how to proceed; if not, a more general case is considered in Section 16.6.4 below.

16.5.3 Approximations Inspired by Term Structure Models

While we do not recommend using a full term structure model to value Libor-with-delay cash flows, it is possible to use such models to derive the relationship between the discount factor $P(T, T_p)$ and the Libor rate $L(T)$. For instance, Andreasen [2002] suggests the one-factor quasi-Gaussian model for the task. While a purely Gaussian model (which has been our standard for these types of approximations in earlier chapters) would give very similar results, we use the qG model here for variety. To develop the approximation, we recall (13.5),

$$P(T, M) = P(T, M, x(T), y(T)), \quad M \geq T,$$

$$P(T, M, x, y) = \frac{P(0, M)}{P(0, T)} \exp \left(-G(T, M)x - \frac{1}{2}G(T, M)^2 y \right),$$

where $x(T)$, $y(T)$ are the state variables in the model and $G(T, M)$ is a function of mean reversion, see (13.3).

Writing the forward Libor rate $L(t) = L(t, T, M)$ as $L(t, x, y)$ to emphasize its dependence on state variables, we note that $P(T, T_p, x, y)$ and $L(T, x, y)$ are monotonic in x , for any fixed value of the state variable $y(T)$. We recall from Chapter 13 that the state variable $y(T)$ is a locally deterministic auxiliary variable whose role is to keep the model arbitrage-free. As done many times in Chapter 13, for the purposes of deriving an approximation let us fix its time T value at some deterministic level $\bar{y}(T)$. Then the discount bond can be expressed in terms of the Libor rate directly,

$$P(T, T_p) = P(T, T_p, X(T, L(T)), \bar{y}(T)), \quad (16.41)$$

where $X(T, l)$ is an inverse (in x) function to $L(T, x, \bar{y}(T))$.

To derive a suitable expression for $\bar{y}(T)$, we recall (Proposition 13.1.4, equation (10.41)) that in the quasi-Gaussian model

$$\mathbb{E}(y(T)) \approx \text{Var}(x(T)).$$

Then, assuming a nearly linear relationship between the Libor rate and the state variable x , we write

$$\text{Var}(L(T)) \approx \left(\frac{\partial L(T, x, y)}{\partial x} \Big|_{x=y=0} \right)^2 \text{Var}(x(T)).$$

This expression suggests the following approximation for $\bar{y}(T)$,

$$\bar{y}(T) = \left(\frac{\partial L(T, x, y)}{\partial x} \Big|_{x=y=0} \right)^{-2} \text{Var}(L(T)),$$

where we approximate $\text{Var}(L(T))$ with the variance in the corresponding forward measure, $\text{Var}(L(T)) \approx \text{Var}^M(L(T))$, and compute the latter either in the vanilla model calibrated to the volatility smile of options on $L(T)$, or directly by the replication method. With $\bar{y}(T)$ set this way, we can solve (16.41) numerically, thereby establishing the required relationship between $P(T, T_p)$ and $L(T)$ and allowing the replication method to be applied. The model thus constructed will violate the no-arbitrage condition (16.19), but only mildly so; moreover we can fix this violation with the application of the scaling idea from Section 16.6.7.

Other term structure models can be used to link $P(T, T_p)$ to the Libor rate. For example, later in the chapter (in Section 16.6.6) we develop approximations that are inspired by Libor market models.

16.5.4 Applications to Averaging Swaps

Libor-with-delay cash flows do not, as a rule, trade individually but instead serve as building blocks for other derivatives; common among them are the so-called *averaging swaps*, i.e. swaps that are composed of *averaging cash flows*. An averaging cash flow (recall Section 5.7) pays, at time T_p , an average Libor rate \bar{L} over the period,

$$\bar{L} = \sum_{i=1}^k w_i L(t_i^f, t_i^s, t_i^e),$$

where we use the notation of Section 5.7. The value of the averaging cash flow at time 0 is given by

$$V_{\text{avg}}(0) = \beta(0) E(\beta(T_p)^{-1} \bar{L})$$

which, using linearity of the payoff and applying appropriate measure changes, is given by

$$\begin{aligned} V_{\text{avg}}(0) &= \sum_{i=1}^k w_i \beta(0) E\left(\beta(T_p)^{-1} L(t_i^f, t_i^s, t_i^e)\right) \\ &= \sum_{i=1}^k w_i P(0, t_i^e) E^{t_i^e} \left(\frac{P(t_i^f, T_p)}{P(t_i^f, t_i^e)} L(t_i^f, t_i^s, t_i^e) \right). \end{aligned}$$

Each term in the sum is the value of a Libor-with-delay cash flow, and can be evaluated by one of the methods developed above.

16.6 CMS and CMS-Linked Cash Flows

The discussion of issues around valuation of Libor-in-arrears and Libor-with-delay cash flows provides us with a useful blueprint for modeling the larger, and more important, class of CMS and CMS-linked cash flows. Using the notations (16.15)–(16.16), we recall from Section 5.11 that a CMS cash flow pays the swap rate $S(T)$ at time T_p , $T_p \geq T$, typically with⁸ $T_p \leq T_1$. More generally, a CMS-linked cash flow pays some function of the swap rate $S(T)$.

In a direct analogy to the LIA case, the market-implied distribution of $S(T)$ in the swap measure Q^A is known from market values of European swaptions; yet, a CMS cash flow is more naturally valued in the T_p -forward measure, the measure linked to a discount bond maturing on the CMS payment date. Not surprisingly, this gives rise to a convexity adjustment. More precisely, note that the value of a CMS cash flow is given by the following expectation in the annuity measure,

$$V_{\text{CMS}}(0) = A(0)E^A \left(\frac{P(T, T_p)}{A(T)} S(T) \right), \quad (16.42)$$

whereby we can define the *CMS convexity adjustment* to be

$$\begin{aligned} D_{\text{CMS}}(0) &\triangleq E^{T_p}(S(T)) - S(0) \\ &= \frac{A(0)}{P(0, T_p)} E^A \left(\frac{P(T, T_p)}{A(T)} S(T) \right) - S(0). \end{aligned} \quad (16.43)$$

At a high level, our discussion of Libor-with-delay valuation issues (Section 16.5) readily extends to the case of CMS cash flows. As always, the expected value in (16.42) can, in principle, be computed with the help of a term structure model, but this is typically too inaccurate and too slow⁹, so the replication method is typically a better choice. Of course, in order to apply it, the payoff of (16.42) and, in particular, the multiplier $P(T, T_p)/A(T)$, needs to be represented as a function of $S(T)$ only. We shall discuss this in a moment, but first we briefly review the replication method for CMS cash flows. We have already discussed replication for Libor-linked cash flows, but we find it worthwhile to examine the method again in a CMS-specific setting.

⁸Recall that T_1 is the first payment date of the swap underlying the swap rate S .

⁹It is nevertheless often useful to be able to calculate CMS convexity adjustments in a given term structure model in closed form, e.g. for assessing the loss of precision of the model when pricing exotics linked to CMS rates. We return to this task later in the chapter.

16.6.1 The Replication Method for CMS

In close analogy to the LIA case, when the replication method of Proposition 8.4.13 is applied to the CMS payoff in (16.42), it decomposes the CMS payoff into a portfolio of standard European options on the swap rate, i.e. swaptions; from this representation the market value of the payoff may be obtained by simply summing up swaption values. These swaption values can be taken directly from the market or, if our goal is to compute a model price, from a given model.

We shall discuss ways of linking the term $P(T, T_p)/A(T)$ to the swap rate $S(T)$ momentarily; for now let us simply assume that an *annuity mapping function* $\alpha(s)$ has been found such that

$$E^A \left(\frac{P(T, T_p)}{A(T)} S(T) \right) = E^A (\alpha(S(T)) S(T)). \quad (16.44)$$

Proposition 8.4.13 then stipulates that

$$\begin{aligned} E^A (\alpha(S(T)) S(T)) &= S(0)\alpha(S(0)) \\ &+ \int_{-\infty}^{S(0)} w(K)p(0, S(0); T, K) dK + \int_{S(0)}^{\infty} w(K)c(0, S(0); T, K) dK, \end{aligned} \quad (16.45)$$

where the hedge weights $w(s)$ are given by

$$w(s) = \frac{d^2}{ds^2}(\alpha(s)s),$$

and $p(t, S; T, K)$ ($c(t, S; T, K)$) are put (call) options on the rate $S(T)$ with strike K , forward S and fixing at T , as observed at t . Combining (16.42), (16.44), and (16.45), we obtain

$$\begin{aligned} V_{\text{CMS}}(0) &= A(0)S(0)\alpha(S(0)) \\ &+ \int_{-\infty}^{S(0)} w(K)V_{\text{rec}}(0, K) dK + \int_{S(0)}^{\infty} w(K)V_{\text{pay}}(0, K) dK, \end{aligned} \quad (16.46)$$

where $V_{\text{rec}}(0, K)$ ($V_{\text{pay}}(0, K)$) are the values, at time 0, of receiver (payer) European swaptions, respectively:

$$\begin{aligned} V_{\text{rec}}(0, K) &= A(0)E^A ((K - S(T))^+), \\ V_{\text{pay}}(0, K) &= A(0)E^A ((S(T) - K)^+). \end{aligned}$$

As mentioned above, the swaption values can either be computed in a model of choice or directly observed in the market. We emphasize again that not

only does the replication method calculate the value of a CMS cash flow consistently with the market in swaptions for all strikes, but it also provides a static, model-independent (up to the choice of $\alpha(s)$) hedging portfolio of payer and receiver swaptions.

In the basic expression (16.46) we can impose various restrictions on swaption hedge positions to, say, incorporate liquidity constraints into the price of the CMS cash flow. For example, swaptions of very low or very high strikes may not be easily tradeable. Then, supposing that the lowest available strike is K_{\min} , and the highest one is K_{\max} , one can choose to pay no more than

$$A(0)S(0)\alpha(S(0)) + \int_{K_{\min}}^{S(0)} w(K)V_{\text{rec}}(0, K) dK + \int_{S(0)}^{K_{\max}} w(K)V_{\text{pay}}(0, K) dK,$$

on the grounds that this is the value that can be “locked in” by hedging with available vanillas. Adjustment for the fact that only a finite number of strikes are traded can proceed along the lines of the discussion in Section 16.4.

The replication approach extends virtually unchanged to cash flows that pay an arbitrary — but reasonably smooth — function of the swap rate, say $g(S(T))$, paid at time $T_p \geq T$. Notable examples of such cash flows include CMS caplets and floorlets (see Section 5.11),

$$g(s) = g_{\text{caplet}}(s) = (s - K)^+, \quad g(s) = g_{\text{floorlet}}(s) = (K - s)^+. \quad (16.47)$$

The value of a CMS-linked cash flow is then, naturally, equal to

$$V_{g^{\text{CMS}}}(0) = A(0)E^A \left(\frac{P(T, T_p)}{A(T)} g(S(T)) \right) = A(0)E^A (\alpha(S(T))g(S(T))), \quad (16.48)$$

and the replication method, as given by (16.46), applies with the weights calculated by

$$w(s) = \frac{d^2}{ds^2} (\alpha(s)g(s)). \quad (16.49)$$

For many payoffs of interest the second derivative here will not be defined in a conventional sense at all points, and may, in particular, contain Dirac delta functions. For example, for the CMS caplets and floorlets with strike K , the second derivative in (16.49) would include a delta function centered at K . This is, however, not a cause for concern, as delta functions are easy to handle in the integrals in (16.46): a delta function centered at some point s_0 would just contribute a term $V_{\text{rec}}(0, s_0)$ (or $V_{\text{pay}}(0, s_0)$, depending on the relationship between s_0 and K) to the integrals in the replication method.

We observe in passing that the replication method requires calculation of values for a collection of swaptions of different strikes. For some models, such calculations can be optimized, see e.g. the discussion of Section 8.4.5.

16.6.2 Annuity Mapping Function as a Conditional Expected Value

The previous section introduced the annuity mapping function $\alpha(s)$ as a critical ingredient of CMS valuation, but stopped short of developing a method to determine it. From examples presented earlier in this chapter, the reader might reasonably expect that terminal swap rate models and/or approximations inspired by term structure models could be used for that purpose. This is indeed the case, as we shall show momentarily. First, however, we find it useful to step back a little to determine the actual theoretical meaning of annuity mapping functions; this analysis is illuminating in its own right and is also helpful in developing a systematic approach to finding good approximations.

Let us start with the main valuation formula (16.42). We obtain

$$\begin{aligned} V_{\text{CMS}}(0) &= A(0)E^A \left(\frac{P(T, T_p)}{A(T)} S(T) \right) \\ &= A(0)E^A \left(E^A \left(\frac{P(T, T_p)}{A(T)} S(T) \middle| S(T) \right) \right) \\ &= A(0)E^A \left(S(T)E^A \left(\frac{P(T, T_p)}{A(T)} \middle| S(T) \right) \right). \end{aligned}$$

Now, if we compare this formula to (16.44), we obtain the following useful result.

Proposition 16.6.1. *The annuity mapping function $\alpha(s)$ in (16.44) or, more generally, (16.48) may be written as the conditional expectation*

$$\alpha(s) = E^A \left(\frac{P(T, T_p)}{A(T)} \middle| S(T) = s \right). \quad (16.50)$$

This result is model-independent.

The proposition clarifies the role of various methods of linking discount bond values to rates that we introduced previously in order to value approximately single-rate derivatives. These methods, in fact, can be seen as approximations to the true annuity mapping function defined by the conditional expected value in (16.50). We shall return to this interpretation and explore it in more detail as we discuss various methods individually below. For now, we note that the problem of calculating the conditional expected value in (16.50) could be attacked directly, as we demonstrate later in Section 16.6.6 for the LM model, or by projection methods. To elaborate briefly, the expected value of random variable X conditional on some other random variable Y can be interpreted as a projection of X on the space of all (suitably regular) functions of Y . Let us denote such a space by \mathcal{B} ; then

$$E(X|Y) = f^*(Y), \quad \text{where } f^* = \operatorname{argmin} \left\{ E \left((X - f(Y))^2 \right), \quad f \in \mathcal{B} \right\}.$$

Following Antonov and Misirpashaev [2009a], we can then obtain a tractable approximation to the true value of the conditional expected value by restricting the subspace of functions of Y to project on. For a given subspace $\tilde{\mathcal{B}} \subset \mathcal{B}$, an approximation is then defined as the closest, in the least-squares sense, element of the subspace to X ,

$$E(X|Y) \approx f^*(Y), \quad \text{where } f^* = \operatorname{argmin} \left\{ E \left((X - f(Y))^2 \right), \quad f \in \tilde{\mathcal{B}} \right\}.$$

If the subspace $\tilde{\mathcal{B}}$ is defined by a parametric functional form,

$$\tilde{\mathcal{B}} = \{f(y; \theta), \theta \in \Theta\}$$

for some parametric set $\Theta \subset \mathbb{R}^d$, then the necessary condition for $f^*(y) \triangleq f(y; \theta^*)$ to be optimal is given by the equations

$$\frac{\partial}{\partial \theta_i} E \left((X - f(Y; \theta))^2 \right) = 0, \quad i = 1, \dots, d.$$

For later use, let us formalize this result as a proposition.

Proposition 16.6.2. *Given two random variables X and Y and a parametric set of functions $\{f(y; \theta)\}$, $\theta \in \Theta \subset \mathbb{R}^d$, an approximation to $E(X|Y)$ is given by*

$$E(X|Y) \approx f(Y; \theta^*),$$

where θ^* is a solution to the set of equations

$$E \left(X \frac{\partial f}{\partial \theta_i} (Y; \theta) \right) = E \left(f(Y; \theta) \frac{\partial f}{\partial \theta_i} (Y; \theta) \right), \quad i = 1, \dots, d. \quad (16.51)$$

16.6.3 Swap-Yield TSR Model

As our first concrete model, let us consider the swap-yield model of Section 16.3.4, a model that has long been a de-facto standard for linking the annuity to the swap rate. Recalling the index function $q(M)$ defined by (16.33), we link discount factors $P(T, M)$ to the swap rate by the formula

$$P(T, M) = \left(\prod_{i=0}^{q(M)-1} \frac{1}{1 + \tau_i S(T)} \right) \times \left(\frac{1}{1 + \tau_{q(M)} S(T)} \right)^{(M - T_{q(M)})/\tau_{q(M)}},$$

with $M \geq T$. As (16.34) holds, we have

$$\begin{aligned}
A(T) &= \sum_{n=0}^{N-1} \tau_n P(T, T_{n+1}) \\
&= \sum_{n=0}^{N-1} \tau_n \prod_{i=0}^n \frac{1}{1 + \tau_i S(T)} = \frac{1}{S(T)} \left(1 - \prod_{i=0}^{N-1} \frac{1}{1 + \tau_i S(T)} \right).
\end{aligned}$$

Also, assuming $T_p \in [T_0, T_1]$ (with obvious modifications for the general case),

$$P(T, T_p) = \left(\frac{1}{1 + \tau_0 S(T)} \right)^{(T_p - T)/\tau_0}.$$

Then

$$\alpha(s) = s \frac{\left(\frac{1}{1 + \tau_0 s} \right)^{(T_p - T)/\tau_0}}{1 - \prod_{i=0}^{N-1} \frac{1}{1 + \tau_i s}}$$

defines the function $\alpha(s)$ to be used in (16.44). We note, as before, that the model will violate basic arbitrage restrictions, in the sense that

$$E^A \left(\frac{P(T, T_p)}{A(T)} \right) \neq \frac{P(0, T_p)}{A(0)}.$$

A method to correct for this is shown in Section 16.6.7.

16.6.4 Linear and Other TSR Models

As all terminal swap rate models are specifically designed to relate discount bonds of various maturities on a particular date T to a market rate $S(T)$, it is easy to extend the discussion in Section 16.6.3 to the general TSR model class. In the TSR class, the linear TSR model (see Section 16.3.2) is arguably the simplest and probably the most popular, so let us use this model as our second concrete example. Applied to the problem of CMS cash flow valuation, the model postulates a linear relationship between the inverse annuity and the swap rate,

$$\alpha(s) = \alpha_1 s + \alpha_2 \quad (16.52)$$

(in the notation of Section 16.3.2 we have $\alpha_1 = a(T_p)$, $\alpha_2 = b(T_p)$). The parameter α_1 can be considered an exogenous input, and α_2 determined by the no-arbitrage requirement that

$$\frac{P(0, T_p)}{A(0)} = E^A \left(\frac{P(T, T_p)}{A(T)} \right) = E^A (\alpha_1 S(T) + \alpha_2) = \alpha_1 S(0) + \alpha_2,$$

implying

$$\alpha_2 = \frac{P(0, T_p)}{A(0)} - \alpha_1 S(0). \quad (16.53)$$

With this specification,

$$\begin{aligned}
 V_{\text{CMS}}(0) &= A(0)E^A((\alpha_1 S(T) + \alpha_2)S(T)) \\
 &= \alpha_2 A(0)S(0) + \alpha_1 A(0)E^A(S(T))^2 \\
 &= P(0, T_p)S(0) - \alpha_1 A(0)S(0)^2 + \alpha_1 A(0)E^A(S(T))^2 \\
 &= P(0, T_p)S(0) + \alpha_1 A(0)\text{Var}^A(S(T)),
 \end{aligned} \tag{16.54}$$

and the convexity adjustment is then given simply by

$$D_{\text{CMS}}(0) = \alpha_1 \frac{A(0)}{P(0, T_p)} \text{Var}^A(S(T)). \tag{16.55}$$

The variance of $S(T)$ is computed either directly in a model of choice, or by integrating the function $(s - S(0))^2$ against the market-implied probability density of the swap rate obtained from swaption prices or, equivalently¹⁰, by the replication method. The elegant formula (16.55), reminiscent of the Libor-in-arrears formula (16.36), is commonly used in practice, despite the fact that discount factors can become negative in certain states of the simulated world under the specification (16.52). The parameter α_1 can be linked to mean reversion as discussed in Section 16.3.2; we touch on this in more detail in Section 16.6.8. More crudely, we can estimate α_1 from a boundary argument, where we simply observe that for scenarios where the time T yield curve is very low (and $S(T)$ therefore is close to zero), we must have

$$\frac{P(T, T_p)}{A(T)} \approx \frac{1}{\sum_{n=0}^{N-1} \tau_n}.$$

As we may write $\alpha_2 = E^A(\alpha_1 S(T) + \alpha_2 | S(T) = 0)$, this suggests setting

$$\alpha_2 = \frac{1}{\sum_{n=0}^{N-1} \tau_n}, \quad \alpha_1 = \frac{1}{S(0)} \left(\frac{P(0, T_p)}{A(0)} - \alpha_2 \right), \tag{16.56}$$

where the equation for α_1 follows from (16.53). While not exactly state-of-the-art, this simplified approach often yields decent precision. See Figure 16.1 on p. 736 for some representative test results.

To wrap up the discussion of the applications of linear TSR models to CMS pricing, let us briefly touch upon its relationship to the result of Proposition 16.6.1. Setting $X = P(T, T_p)/A(T)$, $Y = S(T)$, and $f(y; \theta) = \theta_1 + \theta_2 x$, $(\theta_1, \theta_2)^\top \in \mathbb{R}^2$, we obtain from Proposition 16.6.2 that the coefficients of the best linear approximation to the annuity mapping function $\alpha(s)$,

$$\alpha(s) = \theta_1^* s + \theta_2^*,$$

are given by the solution to the equations

¹⁰In theory; numerical implementation could lead to slight differences.

$$\begin{aligned}\mathbb{E}^A \left(\frac{P(T, T_p)}{A(T)} \right) &= \mathbb{E}^A (\theta_1 S(T) + \theta_2), \\ \mathbb{E}^A \left(\frac{P(T, T_p)}{A(T)} S(T) \right) &= \mathbb{E}^A ((\theta_1 S(T) + \theta_2) S(T)).\end{aligned}$$

Solving the equations, we obtain the optimal coefficients

$$\theta_1^* = \frac{P(0, T_p)}{A(0)} \frac{D_{\text{CMS}}(0)}{\text{Var}^A(S(T))}, \quad \theta_2^* = \frac{P(0, T_p)}{A(0)} - \theta_1^* S(0). \quad (16.57)$$

The same result could have been obtained by backing out α_1 and α_2 from (16.53) and (16.55), a fact that is not surprising given that both calculations started with a linear approximation to $\alpha(s)$.

The result (16.57), even if somewhat trivially obtainable from (16.53) and (16.55), emphasizes the point that a known magnitude of the CMS adjustment will often uniquely identify the annuity mapping function within a parametric class. Importantly, this can be used to calibrate the annuity mapping function (within a given parametric class) to liquidly traded CMS swaps, the market values of which reveal the size of the convexity adjustment. The calibrated annuity mapping function can then be used to value more complicated, and less liquid, CMS-linked derivatives such as CMS caps or CMS range accruals.

Let us finally note that while the linear TSR model gives us the most amount of analytic tractability, the ideas behind Proposition 16.6.2 could be applied to other types of TSR models as well. For example, if we were to choose functions for discount bonds from the exponential class and apply Proposition 16.6.2, we would obtain a model that is quite similar to the exponential TSR model.

16.6.5 The Quasi-Gaussian Model

As was the case for Libor-with-delay cash flows (see Section 16.5.3), the quasi-Gaussian (qG) model can be used as a source of inspiration for the functional relationship between the annuity and the swap rate. We recall the bond reconstruction formula in the quasi-Gaussian model (see (13.5))

$$\begin{aligned}P(T, M) &= P(T, M, x(T), y(T)), \quad M \geq T, \\ P(T, M, x, y) &= \frac{P(0, M)}{P(0, T)} \exp \left(-G(T, M)x - \frac{1}{2}G(T, M)^2 y \right),\end{aligned}$$

with $x(T)$, $y(T)$ being the state variables of the model and $G(T, M)$ a deterministic function of mean reversion, and define $A(T, x, y)$, $S(T, x, y)$ accordingly. Motivated by Section 16.5.3, we set

$$\bar{y}(T) = \left(\frac{\partial S(T, x, y)}{\partial x} \Big|_{x=y=0} \right)^{-2} \text{Var}^A(S(T)),$$

where $\text{Var}^A(S(T))$, the variance of the swap rate $S(T)$ in the annuity measure, is computed consistently with the model used in the replication method, and define $X(T, s)$ to be the inverse function, in x , of $S(T, x, \bar{y}(T))$. Then we can define the mapping function $\alpha(s)$ by

$$\alpha(s) = \frac{P(T, T_p, X(T, s), \bar{y}(T))}{A(T, X(T, s), \bar{y}(T))}, \quad (16.58)$$

and calculate $V_{\text{CMS}}(0)$ via (16.46). For calculating market-consistent CMS values, the values of swaptions in the replication algorithm should be either taken directly from the market or calculated using a market-calibrated vanilla model. If, on the other hand, our objective is to calculate an analytical approximation to the value of a CMS cash flow in the quasi-Gaussian model (a value that could be used to assess and adjust the valuation of more exotic payoffs linked to CMS rates in the model), then we should value swaptions in the qG model directly, perhaps using an approximation such as Proposition 13.1.10.

16.6.6 The Libor Market Model

A Libor market model can also be used to specify the form of dependency of the annuity on the forward swap rate. While perhaps too complicated for valuing CMS cash flows *per se*, establishing the dependency explicitly would be useful for applications of Libor market models to exotic derivatives that are linked to CMS rates, such as callable CMS range accruals, CMS spread TARNs, and the like (see e.g. Sections 5.13.2 and 5.14). When valuing CMS-linked exotic derivatives, it is often desirable to confirm that the values of CMS convexity adjustments in the Libor market model agree with the “market” convexity adjustments or, at the very least, to quantify, and potentially correct for, any observed differences (see Chapter 21). Of course, one can always use Monte Carlo simulation to calculate CMS adjustments in a Libor market model, but the usual performance considerations favor an analytical or semi-analytical approach.

The subject of calculating CMS adjustments in Libor market models has received some attention in the literature — see e.g. Gatarek [2003] for a representative approach — but most published methods generally boil down to using “freezing” techniques to approximate the drift of the swap rate in the forward measure, a method that is not particularly accurate. A notable exception to this trend is the recent work by Antonov and Arneguy [2009] who calculate the expected value

$$\frac{P(0, T_p)}{A(0)} \mathbb{E}^{T_p} (S(T)) = \mathbb{E}^A \left(\left(\frac{P(T, T_p)}{A(T)} \right) S(T) \right)$$

by deriving an approximate SDE for $P(t, T_p)/A(t)$, and then obtain a linear annuity mapping function via (16.57). Test results given in the paper suggest

that the approach is reasonably accurate; however, we believe that it is important to capture the non-linearity of the annuity mapping function in LMM in order to obtain a truly precise approximation. Our preferred alternative is developed below.

First, we recall (16.50) which states that, independently of the underlying model, $\alpha(s)$ can be interpreted as the expected value of $P(T, T_p)/A(T)$ conditioned on $S(T) = s$. We proceed to derive an approximation to this conditional expected value, consistent with the Libor market model. As we did in Chapter 14, we denote the spanning Libor rates by

$$L_n(t) = L(t, T_n, T_{n+1}), \quad n = 0, \dots, N - 1.$$

Assuming $T_p = T$ for notational simplicity¹¹, we observe that the argument of the conditional expected value, the inverse numeraire $1/A(T)$, is a deterministic function of the vector of Libor rates $\mathbf{L}(T) = (L_0(T), \dots, L_{N-1}(T))^\top$,

$$\frac{1}{A(T)} = \rho(\mathbf{L}(T)), \quad \rho(x) = \left(\sum_{n=0}^{N-1} \tau_n \prod_{i=0}^n (1 + \tau_i x_i)^{-1} \right)^{-1}.$$

Approximating

$$\alpha(s) = \mathbb{E}^A (\rho(\mathbf{L}(T)) | S(T_n) = s) \approx \rho(\mathbb{E}^A (\mathbf{L}(T) | S(T) = s)),$$

we reduce the problem to that of computing $\mathbb{E}^A(\mathbf{L}(T) | S(T) = s)$ in a Libor market model, a problem we can tackle in the usual fashion, through an application of a Gaussian approximation. For concreteness, let us consider the following form of the Libor market model (see Section 14.2.5),

$$\begin{aligned} dL_n(t) &= \sqrt{z(t)} \varphi(L_n(t)) \lambda_n(t)^\top dW^{T_{n+1}}(t), \\ dz(t) &= \theta(z_0 - z(t)) dt + \eta(t) \sqrt{z(t)} dZ(t), \quad z(0) = z_0 = 1, \end{aligned} \tag{16.59}$$

with $\langle dW^{T_{n+1}}(t), dZ(t) \rangle = 0$ for $n = 0, \dots, N - 1$. Here $\lambda_n(t)$ is an m -dimensional deterministic volatility function and $W^{T_{n+1}}(t)$ is an m -dimensional $Q^{T_{n+1}}$ -Brownian motion. To compute $\mathbb{E}^A(\mathbf{L}(T) | S(T) = s)$, we use the following Gaussian approximation to the Q^A -dynamics of Libor and swap rates,

$$L_n(t) \approx \widehat{L}_n(t), \quad S(t) \approx \widehat{S}(t),$$

where

$$\begin{aligned} d\widehat{L}_n(t) &= \varphi(L_n(0)) \lambda_n(t)^\top dW^A(t), \quad \widehat{L}_n(0) = L_n(0), \quad n = 0, \dots, N - 1, \\ d\widehat{S}(t) &= \varphi(S(0)) \left(\sum_{i=0}^{N-1} w_i \lambda_i(t)^\top \right) dW^A(t), \quad \widehat{S}(0) = S(0), \end{aligned}$$

¹¹For $T \neq T_p \leq T_1$ the function ρ to be defined momentarily would be multiplied by a term that expresses $P(T, T_p)$ as an (approximate) function of $x_0 = L_0(T)$ as in, e.g., Section 16.5.1.

with

$$w_i = \frac{\varphi(L_i(0))}{\varphi(S(0))} \frac{\partial S(0)}{\partial L_i(0)}, \quad i = 0, \dots, N-1, \quad (16.60)$$

(see Section 14.4.2 for details on approximating swap rate dynamics in Libor market models). The required expected value is then computed by

$$\begin{aligned} \mathbb{E}^A (L_n(T) | S(T) = s) \\ \approx \mathbb{E}^A \left(\widehat{L}_n(T) \middle| \widehat{S}(T) = s \right) = L_n(0) \left(1 + c_n \frac{s - S(0)}{S(0)} \right), \end{aligned}$$

where

$$c_n = \frac{\varphi(L_n(0))S(0) \int_0^T \lambda_n(t)^\top \left(\sum_{i=0}^{N-1} w_i \lambda_i(t) \right) dt}{\varphi(S(0))L_n(0) \int_0^T \left\| \sum_{i=0}^{N-1} w_i \lambda_i(t) \right\|^2 dt}. \quad (16.61)$$

Putting all steps together, we obtain the following proposition.

Proposition 16.6.3. *The mapping function $\alpha(s)$ defined by (16.50) in the Libor market model (16.59) is approximately given by*

$$\alpha(s) = \mathbb{E}^A \left(\frac{1}{A(T)} \middle| S(T) = s \right) \approx \left(\sum_{n=0}^{N-1} \tau_n \prod_{i=0}^n (1 + \tau_i l_i(s))^{-1} \right)^{-1},$$

where

$$l_n(s) = L_n(0) \left(1 + c_n \frac{s - S(0)}{S(0)} \right),$$

for $n = 0, \dots, N-1$, with coefficients c_n given by (16.61) and weights w_n given by (16.60).

16.6.7 Correcting Non-Arbitrage-Free Methods

Several of the annuity mapping methods developed in previous sections (e.g. in Sections 16.6.3, 16.6.5 and 16.6.6) are not arbitrage-free by construction. Others, such as the linear TSR model, may theoretically be arbitrage-free, but numerical methods may induce slight errors. In this section we introduce a simple adjustment to all methods that will remedy the main arbitrage issues, be they theoretical or numerical.

We recall that the principal valuation formula for CMS cash flows specifies

$$V_{\text{CMS}}(0) = A(0) \mathbb{E}^A \left(\frac{P(T, T_p)}{A(T)} S(T) \right) = A(0) \mathbb{E}^A (\alpha(S(T)) S(T)),$$

where $\alpha(s)$ is obtained by one of the methods discussed above. The quantity $P(T, T_p)/A(T)$, being a ratio of a tradeable asset and the numeraire, is a

martingale in the annuity measure. Hence, in any arbitrage-free model the following should hold,

$$\mathbb{E}^A \left(\frac{P(T, T_p)}{A(T)} \right) = \frac{P(0, T_p)}{A(0)}.$$

That is, we should have

$$\mathbb{E}^A (\alpha(S(T))) = \frac{P(0, T_p)}{A(0)}. \quad (16.62)$$

If, however, the function $\alpha(s)$ is obtained by one of the methods that does not satisfy the no-arbitrage condition, we would see that

$$\bar{\alpha} \triangleq \mathbb{E}^A (\alpha(S(T))) \neq \frac{P(0, T_p)}{A(0)}. \quad (16.63)$$

For purposes of CMS product valuations, the inequality (16.63) is, by far, the most important manifestation of the arbitrage in the model. Pragmatically, we can compensate by rescaling the original function $\alpha(s)$ to force (16.62) to be satisfied. In particular, defining

$$\tilde{\alpha}(s) = \frac{P(0, T_p)}{A(0)} \frac{\alpha(s)}{\bar{\alpha}}, \quad (16.64)$$

we obtain the “improved” CMS valuation formula,

$$V_{\text{CMS}}(0) = A(0) \mathbb{E}^A (S(T) \tilde{\alpha}(S(T))) = P(0, T_p) \frac{\mathbb{E}^A (S(T) \alpha(S(T)))}{\mathbb{E}^A (\alpha(S(T)))}.$$

In fact, the correction (16.64) is useful even for arbitrage-free models; while the no-arbitrage property (16.62) holds in theory, in practice it can be violated in the numerical scheme used.

Apart from the fundamental “test” (16.62) that any annuity mapping method must pass, there are other checks that are useful to keep in mind while looking at any particular method for CMS product valuation. One such test is obtained from the trivial identity

$$\sum_{n=0}^{N-1} \frac{\tau_n P(T, T_{n+1})}{A(T)} = 1,$$

which implies the following relationship between annuity mapping functions $\alpha(s, T_n)$ that correspond to different payment dates T_n , $n = 1, \dots, N$ (note how we enriched the notation for the annuity mapping function with the payment date for the moment),

$$\sum_{n=0}^{N-1} \tau_n \mathbb{E}^A (\alpha(S(T), T_{n+1}) S(T)) = S(0).$$

This identity, in effect, states that the sum of CMS convexity adjustments with payment dates running over all tenor dates of the swap rate should be equal to zero.

Another useful check is obtained if we recall that

$$\frac{P(T, T_0)}{A(T)} - \frac{P(T, T_N)}{A(T)} = S(T)$$

(here of course $T_0 = T$, but we write it as such to highlight the symmetry in the expression). Multiplying both sides by $S(T)$, applying the expected value operator, and using the extended notation $\alpha(s, M)$ for the M -payment-date annuity mapping function, we obtain another identity that should be satisfied,

$$\mathbb{E}^A(\alpha(S(T), T_0)S(T)) - \mathbb{E}^A(\alpha(S(T), T_N)S(T)) = \mathbb{E}^A(S(T)^2).$$

The right-hand side here can be obtained from the Q^A -distribution of the swap rate $S(T)$ by replication and, as such, is independent of any annuity mapping function. Therefore, for any annuity mapping method, this identity — i.e. the requirement that the difference of CMS payments paid on the swap rate fixing date and on the last payment date be annuity-mapping-independent — represents another constraint that should be satisfied by the method.

16.6.8 Impact of Annuity Mapping Function and Mean Reversion

The importance of capturing volatility smile in CMS valuation, typically through the replication method (16.46), is widely acknowledged. On the other hand, the impact of other components entering into CMS valuation, especially the annuity mapping function $\alpha(s)$, is sometimes overlooked. One does not need to look any further than at the formula (16.54) for the CMS value in a linear TSR model to understand potential issues: if the parameter α_1 in (16.54) is allowed to vary freely, the CMS convexity adjustment can be made arbitrarily small or large.

Of course not all values of α_1 are compatible with financial reality, but choosing a reasonable range for α_1 is not entirely trivial. Relating the parameter to mean reversion as we did in Section 16.3.2 is useful, since we understand the role of mean reversion and its impact on model dynamics reasonably well. Moreover, mean reversion can be directly linked to market prices of traded securities, as shown in Section 13.1.8. It turns out that CMS convexity adjustment can vary by 10%-20% (in relative terms) when using different but reasonable levels of mean reversions.

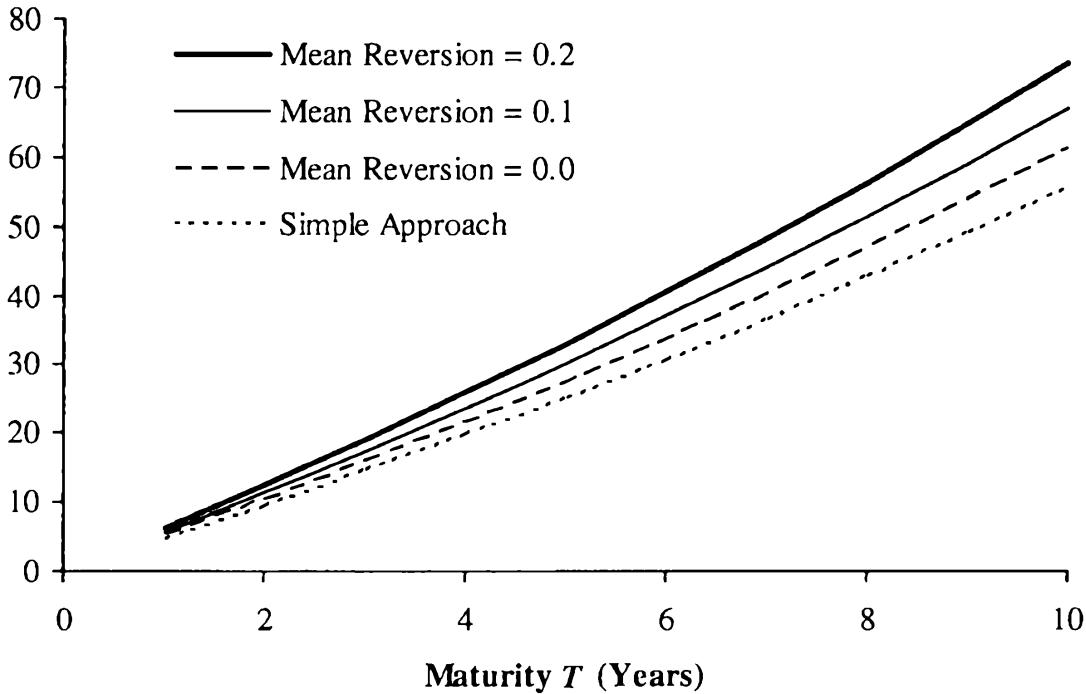
To demonstrate the effect of mean reversion on the CMS convexity adjustment, consider the concrete problem of estimating the time 0 forward rate for a 10 year swap rate with semi-annual fixings. We use the linear

TSR method and assume that interest rates are flat at 5% (continuously compounded), and that the par swap rate is log-normally distributed with a constant volatility of $\sigma_S = 17\%$ for all fixing dates (which is hardly consistent with the presence of mean reversion, but good enough for a numerical example). Under this assumption,

$$\text{Var}^A(S(T)) = S(0)^2 \left(e^{\sigma_S^2 T} - 1 \right),$$

allowing us to compute the convexity adjustment (16.55) in closed form for a given maturity T . We estimate the coefficient α_1 in (16.55) in two ways: by the simplified approach (16.56), or from the more elaborate formula (16.27) that takes mean reversion κ as a parameter. Results are shown in Figure 16.1, at multiple values of T and κ .

Fig. 16.1. CMS Convexity Adjustment (Basis Points)



Notes: CMS convexity adjustment (16.55) in basis points for the linear TSR model, as computed by formula (16.27) (at three different mean reversion levels) and by the simplified formula (16.56). The rate and volatility settings are described in the text.

Notice that the convexity adjustment increases with mean reversion, a consequence of the fact that the volatility of the annuity factor effectively increases when the mean reversion goes up. Consistent with this, the slope parameter α_1 increases in mean reversion. The simple approach in formula (16.56) results in lower convexity adjustments than the mean reversion based approach.

16.6.9 CDF and PDF of CMS Rate in Forward Measure

The replication method is very useful for pricing CMS-linked cash flows, but it is not always convenient to apply. In particular, if the payoff is discontinuous, then the calculation of weights in, e.g., (16.49) will require special care. Let us attempt to develop a suitable alternative. We start by noting that the problem of pricing a cash flow that pays $g(S(T))$ at time T_p , $T_p \geq T$, (as in Section 16.6.1) can be seen as a problem of determining a density of $S(T)$ in the T_p -forward measure, $\psi^{T_p}(s)$, as we can always write

$$\mathbb{E}^{T_p}(g(S(T))) = \int_{-\infty}^{\infty} g(s) \psi^{T_p}(s) ds. \quad (16.65)$$

This density is not directly available; however, the cumulative distribution function (CDF) $\Psi^A(\cdot)$ and the probability density function (PDF) $\psi^A(\cdot)$ of a swap rate in the *annuity* measure are available in either closed form for a particular vanilla model calibrated to market, or non-parametrically from the market prices of swaptions of all strikes (see Section 7.1.2) via

$$\Psi^A(K) = 1 + \frac{\partial}{\partial K} c(K), \quad (16.66)$$

$$\psi^A(K) = \frac{\partial^2}{\partial K^2} c(K), \quad (16.67)$$

$$c(K) = \mathbb{E}^A((S(T) - K)^+). \quad (16.68)$$

The following proposition allows us to obtain the CDF and PDF of the swap rate in the forward measure from its distributional characteristics in the annuity measure.

Proposition 16.6.4. *Given an annuity mapping function $\alpha(s)$ defined by (16.50), the PDF $\psi^{T_p}(s)$ and the CDF $\Psi^{T_p}(s)$ of the swap rate in the T_p -forward measure are linked to the PDF $\psi^A(s)$ and the CDF $\Psi^A(s)$ of the swap rate in the annuity measure by*

$$\psi^{T_p}(s) = \frac{A(0)}{P(0, T_p)} \alpha(s) \psi^A(s), \quad (16.69)$$

$$\Psi^{T_p}(s) = \frac{A(0)}{P(0, T_p)} \int_{-\infty}^s \alpha(u) \psi^A(u) du. \quad (16.70)$$

Proof. Proceeding somewhat informally, we observe that the value of the density $\psi^{T_p}(K)$ at point K is equal to the (undiscounted) value of the security with the delta-function payoff, $\delta(S(T) - K)$,

$$\psi^{T_p}(K) = \mathbb{E}^{T_p}(\delta(S(T) - K)).$$

By switching to the annuity measure, using the law of iterated conditional expectations, and the definition (16.50) of $\alpha(s)$, we obtain

$$\begin{aligned}
\psi^{T_p}(K) &= \frac{A(0)}{P(0, T_p)} E^A \left(\frac{P(T, T_p)}{A(T)} \delta(S(T) - K) \right) \\
&= \frac{A(0)}{P(0, T_p)} E^A (\alpha(S(T)) \delta(S(T) - K)) \\
&= \frac{A(0)}{P(0, T_p)} \alpha(K) E^A (\delta(S(T) - K)) \\
&= \frac{A(0)}{P(0, T_p)} \alpha(K) \psi^A(K).
\end{aligned}$$

The statement (16.70) follows trivially. \square

In practice, we would use one of the approximations to $\alpha(s)$ as derived in, for example, Sections 16.6.3, 16.6.4, 16.6.5, or 16.6.6. The density integration method (16.65) with the density $\psi^{T_p}(s)$ given by (16.69) is theoretically equivalent to the replication method of Section 16.6.1, but can, as hinted at earlier, have better numerical properties for non-smooth or discontinuous payoffs such as digital options or range accruals on a CMS rate. Indeed, unlike (16.49), the density integration method does not involve payoff differentiation. Another important application of the method arises when valuing cash flows linked to *multiple* CMS rates as we discuss in Chapter 17.

The expression (16.70) for the CDF has a particularly simple form when the function $\alpha(s)$ is linear as in, for example, the linear TSR model (Section 16.3.2).

Corollary 16.6.5. *In the linear TSR model (16.52), the CDF $\Psi^{T_p}(s)$ of the swap rate in the T_p -forward measure is given by*

$$\Psi^{T_p}(s) = \frac{A(0)}{P(0, T_p)} (\alpha_1 (S(0) - s - c(s)) + \alpha(s) \Psi^A(s)), \quad (16.71)$$

where the CDF in the annuity measure $\Psi^A(s)$ is given by (16.66), and $c(s)$ is the option price with strike s in (16.68).

Proof. We have,

$$\begin{aligned}
\frac{P(0, T_p)}{A(0)} \Psi^{T_p}(s) &= \int_{-\infty}^s (\alpha_1 u + \alpha_2) \psi^A(u) du \\
&= \alpha_1 \left(\int_{-\infty}^{\infty} u \psi^A(u) du - \int_s^{\infty} (u - s) \psi^A(u) du - s \int_s^{\infty} \psi^A(u) du \right) \\
&\quad + \alpha_2 \int_{-\infty}^s \psi^A(u) du.
\end{aligned}$$

We note that

$$\begin{aligned}
\int_{-\infty}^{\infty} u \psi^A(u) du &= S(0), & \int_s^{\infty} (u - s) \psi^A(u) du &= c(s), \\
\int_s^{\infty} \psi^A(u) du &= 1 - \Psi^A(s), & \int_{-\infty}^s \psi^A(u) du &= \Psi^A(s),
\end{aligned}$$

and the result follows. \square

Corollary 16.6.6. *In the linear TSR model (16.52), the PDF $\psi^{T_p}(s)$ of the swap rate in the T_p -forward measure is given by*

$$\psi^{T_p}(s) = \frac{A(0)}{P(0, T_p)} (\alpha_1 s + \alpha_2) \psi^A(s), \quad (16.72)$$

where its PDF in the annuity measure $\psi^A(s)$ is given by (16.67).

Proof. Either directly from (16.69) or by differentiating (16.71) and using (16.66). \square

At this point the reader may wonder if the PDF $\psi^{T_p}(s)$ is directly linked to prices of traded derivatives. To answer this, we should recall the definition of a CMS caplet payoff in (16.47). Clearly,

$$V_{\text{cmscaplet}}(0, K) = P(0, T_p) E^{T_p} \left((S(T) - K)^+ \right), \quad (16.73)$$

and we therefore have the following result (compare with (16.67)):

Lemma 16.6.7. *The market-implied PDF $\psi^{T_p}(s)$ of the swap rate in the T_p -forward measure can be directly obtained from values of CMS caplets by*

$$\psi^{T_p}(K) = \frac{1}{P(0, T_p)} \frac{\partial^2}{\partial K^2} V_{\text{cmscaplet}}(0, K). \quad (16.74)$$

16.6.10 SV Model for CMS Rate

A forward measure associated with the payment date of a cash flow is often the most convenient measure to use when dealing with vanilla derivatives linked to *multiple* market rates, as we shall find out in the next chapter. So, it would be useful to have PDFs and CDFs of market rates in the forward measure to be of tractable form, i.e. come from some common parameterizations such as the SV model of Chapter 8. Alas, this is not the case if we start with the SV model in the *annuity* measure for the swap rate, which is of course a common procedure. In this section we discuss these issues and proceed to derive useful approximations.

The formulas in Proposition 16.6.4 and Corollary 16.6.6 specify how PDFs of swap rates change under a measure change, from the annuity measure to the T_p -forward measure. These transformations are independent of the actual models (PDFs) used and are, of course, exact (to the extent that $\alpha(s)$ represents a true conditional expected value). As we indicated above it is useful to have an approximation to the PDF $\psi^{T_p}(s)$ that is from the same family as the PDF in the annuity measure $\psi^A(s)$. To elaborate, assume the swap rate follows

$$dS(t) = \lambda(bS(t) + (1 - b)S(0)) \sqrt{z(t)} dW^A(t), \quad (16.75)$$

$$dz(t) = \theta(1 - z(t)) dt + \eta\sqrt{z(t)} dZ^A(t), \quad z(0) = 1, \quad (16.76)$$

where $\langle dZ^A(t), dW^A(t) \rangle = 0$ and $Z^A(t)$, $W^A(t)$ are Brownian motions in Q^A , the annuity measure. This SDE system defines the distribution, and in particular the PDF $\psi^A(\cdot)$ of $S(T)$, in measure Q^A . Let us now define an adjusted process $\tilde{S}(t)$ given by the following SV dynamics in the T_p -forward measure,

$$\begin{aligned} d\tilde{S}(t) &= \tilde{\lambda} \left(\tilde{b}\tilde{S}(t) + (1 - \tilde{b})\tilde{S}(0) \right) \sqrt{\tilde{z}(t)} dW^{T_p}(t), \\ d\tilde{z}(t) &= \theta(1 - \tilde{z}(t)) dt + \tilde{\eta}\sqrt{\tilde{z}(t)} dZ^{T_p}(t), \quad \tilde{z}(0) = 1, \end{aligned} \quad (16.77)$$

where $Z^{T_p}(t)$, $W^{T_p}(t)$ are Q^{T_p} -Brownian motions satisfying $\langle dZ^{T_p}(t), dW^{T_p}(t) \rangle = 0$, and where we align the mean of $\tilde{S}(T)$ to equal the CMS-adjusted value of $S(T)$, i.e.

$$\tilde{S}(0) = E^{T_p}(S(T)).$$

Finally, we define $\psi^{T_p}(\cdot)$ to be the PDF of $\tilde{S}(T)$, and aim to set the adjusted parameters $\tilde{\lambda}$, \tilde{b} , $\tilde{\eta}$ such that the distribution of $\tilde{S}(T)$ is as close as possible to the distribution of $S(T)$ in the T_p -forward measure.

As measure transformations affect the drift of an SDE, it is clear that the equality (in distribution) of $S(T)$ and $\tilde{S}(T)$ under Q^{T_p} cannot, in general, be achieved exactly, as we here attempt to represent the measure transform as a parameter change solely affecting the diffusion term in the SDE for $S(t)$. Still, as we said in the beginning of the section, such a representation, even if approximate, is often useful for the multi-rate derivatives that we consider in Chapter 17.

The calculation of $(\tilde{S}(0), \tilde{\lambda}, \tilde{b}, \tilde{\eta})$ from $(S(0), \lambda, b, \eta)$ is not trivial and best done numerically. The convexity-adjusted CMS rate $\tilde{S}(0)$ is calculated by the replication method of Section 16.6.1. As we can rewrite (16.73) in the form

$$V_{\text{cmscaplet}}(0, K) = P(0, T_p) E^{T_p} \left((\tilde{S}(T) - K)^+ \right),$$

we note that CMS caplets are nothing more than European call options on $\tilde{S}(T)$. Hence, we can obtain $(\tilde{\lambda}, \tilde{b}, \tilde{\eta})$ by direct calibration of the SV model (16.77) (which is in T_p -forward measure) to prices of CMS caplets *that are computed in the original SV model* for $S(T)$ (with dynamics specified in the annuity measure). These CMS caplet prices in the SV model for $S(t)$ are best obtained by the replication algorithm (16.46) with weights (16.49). As we need CMS caplet prices across a range of strikes for $(\tilde{\lambda}, \tilde{b}, \tilde{\eta})$ -calibration, we can reuse much of the calculations in (16.46), as only the weights (but not the swaption prices used in replication) change in (16.46) for CMS caplets of different strikes.

At this point one may question whether the SV parameters (λ, b, η) perhaps do not change when we switch from the annuity to the forward measure, i.e. that all measure-related changes can be embedded in the change of the forward value, $S(0) \rightarrow \tilde{S}(0) = E^{T_p}(S(T))$. The answer is, of course, a clear no: the relationship between the densities as given by, for example, (16.71) is *not* just a shift of the mean of $S(T)$. In fact we see that the measure change affects the whole distribution, in particular re-distributing the probability mass from the region of lower values of the swap rate to the region of the higher values of the rate (as $\alpha_1 > 0$ typically). Hence, we would expect, at the very least, a change in the skew parameter b ; other parameters will also be affected. This highlights an important point: CMS caps/floors *should not* be valued by simply convexity-adjusting the forward swap rate and then otherwise using the same model with the same parameters as for European swaptions. Despite the fact that it often produces sizable errors, this type of computation nevertheless appears quite common in practice.

16.6.11 Dynamics of CMS Rate in Forward Measure

While in the previous section we absorbed the measure change into the SV diffusion parameters, in reality measure changes affect only (instantaneous) drifts of the SDEs defining the dynamics. Let us explore how this will work for the SV model. While not a particularly useful consideration for single-rate derivatives, this becomes important when we consider Monte Carlo pricing of derivatives linked to multiple rates; see Section 17.8.1 for such applications.

We continue looking at a single swap rate $S(t)$ associated with the annuity $A(t)$, and assume that the following stochastic volatility model is specified in the annuity measure,

$$dS(t) = \lambda \varphi(S(t)) \sqrt{z(t)} dW^A(t), \quad (16.78)$$

$$dz(t) = \theta(1 - z(t)) dt + \eta \psi(z(t)) dZ^A(t), \quad z(0) = 1, \quad (16.79)$$

where $\langle dZ^A(t), dW^A(t) \rangle = 0$ and $Z^A(t)$, $W^A(t)$ are Brownian motions in Q^A , the annuity measure.

We are interested in bringing the two-dimensional SDE (16.78)–(16.79) into the T_p -forward measure, where T_p is the payment time of the CMS contract. However, the dynamics (16.78)–(16.79) in the annuity measure *do not* uniquely define a T_p -forward measure — for that we would need a full term structure model (or, at least, an additional specification for the density process $P(t, T_p)/A(t)$). On the other hand, if we assume that we have knowledge of the conditional expectation (16.50) we can nevertheless construct a probability measure that would resemble a forward measure for European-type payoffs fixing at T and paying at T_p (but not for any other payoffs such as e.g. European derivatives fixing at time other than T). More precisely, we are interested in constructing a measure \tilde{Q}^{T_p} such that, for any function $f(\cdot)$,

$$\mathbb{E}^A \left(\frac{P(T, T_p)}{A(T)} f(S(T)) \right) = \frac{P(0, T_p)}{A(0)} \tilde{\mathbb{E}}^{T_p} (f(S(T))),$$

where $\tilde{\mathbb{E}}^{T_p}$ denotes expectation in measure $\tilde{\mathbb{Q}}^{T_p}$.

Before stating our result, we recall the definition of the function $\alpha(s)$ from (16.50). It is more convenient to deal with a rescaled function, so let us define

$$\hat{\alpha}(s) = \frac{A(0)}{P(0, T_p)} \alpha(s) = \frac{A(0)}{P(0, T_p)} \mathbb{E}^A \left(\frac{P(T, T_p)}{A(T)} \middle| S(T) = s \right).$$

Proposition 16.6.8. Define the measure $\tilde{\mathbb{Q}}^{T_p}$ by the condition that the process $(z(t), S(t))$ satisfies the following SDE in $\tilde{\mathbb{Q}}^{T_p}$,

$$\begin{aligned} dS(t) &= \lambda \varphi(S(t)) \sqrt{z(t)} (dW^{T_p}(t) + v^S(t) dt), \\ dz(t) &= \theta(1 - z(t)) dt + \eta \psi(z(t)) (dZ^{T_p}(t) + v^z(t) dt), \quad z(0) = 1, \end{aligned} \tag{16.80}$$

where $Z^{T_p}(t)$ and $W^{T_p}(t)$ are (uncorrelated) driftless Brownian motions in $\tilde{\mathbb{Q}}^{T_p}$, the drift adjustments are given by

$$\begin{aligned} v^z(t) &= \eta \psi(z(t)) \frac{\partial}{\partial z} \ln(\Lambda(t, z(t), S(t))), \\ v^S(t) &= \lambda \varphi(S(t)) \sqrt{z(t)} \frac{\partial}{\partial S} \ln(\Lambda(t, z(t), S(t))), \end{aligned} \tag{16.81}$$

and the function $\Lambda(t, z, s)$ satisfies the following PDE,

$$\begin{aligned} \frac{\partial}{\partial t} \Lambda(t, z, s) + \theta(1 - z) \frac{\partial}{\partial z} \Lambda(t, z, s) + \frac{\eta^2}{2} \psi(z)^2 \frac{\partial^2}{\partial z^2} \Lambda(t, z, s) \\ + \frac{\lambda^2}{2} \varphi(s)^2 z \frac{\partial^2}{\partial S^2} \Lambda(t, z, s) = 0, \quad t \in [0, T], \end{aligned} \tag{16.82}$$

with the terminal condition at $t = T$

$$\Lambda(T, z, s) = \hat{\alpha}(s). \tag{16.83}$$

Then for any function $f(\cdot)$ we have

$$\frac{A(0)}{P(0, T_p)} \mathbb{E}^A \left(\frac{P(T, T_p)}{A(T)} f(S(T)) \right) = \tilde{\mathbb{E}}^{T_p} (f(S(T))). \tag{16.84}$$

Proof. Clearly, for any function $f(\cdot)$,

$$\frac{A(0)}{P(0, T_p)} \mathbb{E}^A \left(\frac{P(T, T_p)}{A(T)} f(S(T)) \right) = \mathbb{E}^A (\hat{\alpha}(S(T)) f(S(T))).$$

The condition (16.84) would therefore be satisfied if we locate a measure that satisfies

$$\mathbb{E}^A(\widehat{\alpha}(S(T))f(S(T))) = \tilde{\mathbb{E}}^{T_p}(f(S(T)))$$

for any $f(\cdot)$. Recalling the results in Section 1.4, we must find a density process, i.e. a positive \mathbb{Q}^A -martingale, that equals $\widehat{\alpha}(S(T))$ at time T . Such a martingale is easy to construct,

$$\Lambda(t) = \mathbb{E}_t^A(\widehat{\alpha}(S(T))),$$

and this allows us to specify the measure $\tilde{\mathbb{Q}}^{T_p}$ as the measure for which $\Lambda(t)$ is the Radon-Nikodym derivative (with respect to \mathbb{Q}^A).

By Girsanov's theorem, moving from $\tilde{\mathbb{Q}}^A$ to $\tilde{\mathbb{Q}}^{T_p}$ is associated with certain changes in the drift terms of (16.78)–(16.79). Specifically, we recall from Section 1.5 that

$$\begin{aligned} dZ^A(t) &= dZ^{T_p}(t) + \nu^z(t) dt, \\ dW^A(t) &= dW^{T_p}(t) + \nu^S(t) dt, \end{aligned}$$

where

$$d\Lambda(t)/\Lambda(t) = \nu^z(t) dZ^A(t) + \nu^S(t) dW^A(t).$$

By Ito's lemma and the fact that $\Lambda(t)$ is a \mathbb{Q}^A -martingale,

$$\begin{aligned} d\Lambda(t, z(t), S(t)) &= \Lambda_z(t, z(t), S(t)) \eta \psi(z(t)) dZ^A(t) \\ &\quad + \Lambda_S(t, z(t), S(t)) \lambda \varphi(S(t)) \sqrt{z(t)} dW^A(t). \end{aligned}$$

Then the expressions for $\nu^z(t)$, $\nu^S(t)$ follow by matching the dZ^A , dW^A terms in the last two equations.

Finally, to find the expression for $\Lambda(t)$ we recall that since the process $(S(t), z(t))$ is Markovian, we have, with slight abuse of notations,

$$\begin{aligned} \Lambda(t) &= \Lambda(t, z(t), S(t)), \\ \Lambda(t, z, s) &= \mathbb{E}_t^A(\widehat{\alpha}(S(T))|z(t) = z, S(t) = s). \end{aligned}$$

If follows from the Feynman-Kac theorem that the function $\Lambda(t, z, s)$ satisfies the PDE (16.82)–(16.83). \square

Proposition 16.6.8 establishes a numerical scheme for simulating $(z(t), S(t))$ in $\tilde{\mathbb{Q}}^{T_p}$, for purposes of pricing European-style derivatives fixing at a given time T and paying at T_p . In general, we would determine the function $\Lambda(t, z, s)$ by numerically solving the PDE (16.82)–(16.83) on a grid of (t, z, s) , and then perform the Monte Carlo simulation for (16.80), with the drift adjustments $\nu^z(t)$, $\nu^S(t)$ computed for each path using (16.81). In some important cases, however, no finite difference scheme is required, as shown in the following corollary for the often-used case where the function $\alpha(s)$ is linear.

Corollary 16.6.9. Assume that

$$\widehat{\alpha}(s) = \widehat{\alpha}_1 s + \widehat{\alpha}_2.$$

Then

$$v^z(t) = 0, \quad v^S(t) = \lambda \varphi(S(t)) \sqrt{z(t)} \frac{\widehat{\alpha}_1}{\widehat{\alpha}_1 S(t) + \widehat{\alpha}_2}. \quad (16.85)$$

Proof. The swap rate $S(t)$ is a \mathbb{Q}^A -martingale, hence

$$\Lambda(t, z, s) = \mathbb{E}^A(\widehat{\alpha}(S(T)) | z(t) = t, S(t) = t) = \widehat{\alpha}_1 s + \widehat{\alpha}_2.$$

□

16.6.12 Cash-Settled Swaptions

After the mainly theoretical considerations of Section 16.6.11, let us return to applications and consider the important topic of pricing of cash-settled European swaptions. As explained in Section 5.10.1, cash-settled swaptions are the most common type of vanilla options in European markets, especially for derivatives quoted in EUR and GBP. Cash-settled swaptions are closely linked to the swap-settled European swaptions that are standard in the US, but rather than exercising into a physical swap contract, a cash-settled swaption instead uses a particular “swap-like” formula to determine a cash amount to be paid upon option exercise. As it turns out, the replication methods we developed earlier in this chapter allow us to link a value of a cash-settled swaption to those of swap-settled swaptions across a range of strikes. While the two kinds of swaptions are rarely traded in the same market, this connection is nevertheless important as it allows us to continue treating *swap-settled* swaptions as the fundamental type of vanilla options to which we calibrate all our models, irrespective of market conventions. That is, we would maintain a swaption grid of vanilla model parameters (see Section 16.1.3) that represents values of swap-settled swaptions, even if these are not directly traded. The actual model parameter values for each grid point would be calculated by calibration to the values of the most prevalent type of swaptions in the market; for the case of cash-settled swaptions, this would require usage of the valuation formula (16.86) developed below.

As we recall, the payoff of a cash-settled swaption is given by a deterministic function $g(\cdot)$ applied to the swap rate $S(T)$ and paid at T , where the function $g(\cdot)$ is given by

$$g(s) = \left(\sum_{n=0}^{N-1} \tau_n \prod_{i=0}^n (1 + \tau_i s)^{-1} \right) (s - K)^+$$

for a payer swaption. Given the annuity mapping function $\alpha(s)$, the value is then given by

$$V_{\text{CSS}}(0) = A(0)E^A(\alpha(S(T))g(S(T))), \quad (16.86)$$

and is easily calculated by the replication method applied to prices for swap-settled swaptions.

In application of (16.86), we would need to fix a choice for the annuity mapping function α ; for this, we typically recommend using the linear TSR model of Section 16.3.2. Interestingly, if we were alternatively to use the swap-yield model of Section 16.6.3, we would get

$$\alpha(s)g(s) = (s - K)^+,$$

as the swap-yield annuity mapping function exactly cancels the annuity discounting term in the payoff. We would therefore get

$$V_{\text{CSS}}(0) = A(0)E^A((S(T) - K)^+),$$

which is the value of a swap-settled swaption. In reality the values of swap- and cash-settled swaptions should, of course, be different; the inability of the swap-yield model to distinguish them is a symptom of the fact that the swap-yield model is not a truly arbitrage-free model.

As cash-settled swaptions differ from swap-settled swaptions, they also do not obey the “standard” call-put parity,

$$\begin{aligned} & V_{\text{CSS,pay}}(0) - V_{\text{CSS,rec}}(0) \\ &= A(0)E^A\left(\alpha(S(T))\left(\sum_{n=0}^{N-1} \tau_n \prod_{i=0}^n (1 + \tau_i S(T))^{-1}\right)(S(T) - K)\right) \\ &\neq A(0)(S(0) - K). \end{aligned}$$

Instead, a combined long-short position in a cash-settled payer swaption and a cash-settled receiver position is equivalent to a “cash-settled swap”, i.e. a (typically non-traded) derivative with the payoff

$$\left(\sum_{n=0}^{N-1} \tau_n \prod_{i=0}^n (1 + \tau_i S(T))^{-1}\right)(S(T) - K).$$

16.7 Quanto CMS

All securities discussed so far in this chapter produce payments in the same currency as the currency of the underlying rates used to calculate the payoff. It is, however, not uncommon to use a different currency for payment, a modification that leads to the creation of so-called “quanto” cash flows, see Section 4.3. While we generally limit the scope of this book to single-currency

derivatives only, quanto extensions of CMS-linked derivatives are sufficiently common to warrant a discussion of their valuation.

We have already arrived at some preliminary results on multi-currency markets in Section 4.3, a section that the reader is advised to review before proceeding; the notations of this section will be adopted in what follows.

16.7.1 Overview

Let the swap rate $S(T)$ of Section 16.6 be computed from a *domestic* currency yield curve. A quanto CMS cash flow pays $g(S(T))$ at time T_p in some other *foreign* currency; the value of the cash flow is therefore equal to

$$V_{\text{QuantoCMS}}(0) = \beta_f(0) E^f (\beta_f(T_p)^{-1} g(S(T)))$$

in foreign currency units, where $\beta_f(t)$ is the foreign money market account and E^f is the expected value operator for the foreign risk-neutral measure. Since $S(t)$ is defined by the domestic interest rate curve, its distribution in the domestic (annuity) measure is available from the swaption market. By Lemma 4.3.1, the density process relating the foreign and domestic risk-neutral measures is given by

$$E_t^d \left(\frac{dQ^f}{dQ^d} \right) = \frac{\beta_f(t) X(t)}{\beta_d(t) X(0)}, \quad t \geq 0,$$

where $X(t)$ is the spot FX rate measured in domestic currency per foreign currency units. The value of the contract in the domestic risk-neutral measure may therefore be written (in foreign currency units, naturally) as

$$V_{\text{QuantoCMS}}(0) = \frac{\beta_d(0)}{X(0)} E^d (\beta_d(T_p)^{-1} g(S(T)) X(T_p)). \quad (16.87)$$

Of course, the same formula can be derived by observing that since $g(S(T))$ is paid in a foreign currency, we can convert the proceeds into the domestic currency at time T_p to create a domestic asset with the payoff $g(S(T))X(T_p)$.

Conditioning in (16.87) on \mathcal{F}_T , we get

$$\begin{aligned} V_{\text{QuantoCMS}}(0) &= \frac{\beta_d(0)}{X(0)} E^d (\beta_d(T)^{-1} g(S(T)) [\beta_d(T) E_T^d (\beta_d(T_p)^{-1} X(T_p))]) \\ &= \frac{\beta_d(0)}{X(0)} E^d (\beta_d(T)^{-1} g(S(T)) [P_d(T, T_p) X_{T_p}(T)]), \end{aligned}$$

where we have used the notation from Section 4.3,

$$X_{T_p}(T) = \frac{P_f(T, T_p)}{P_d(T, T_p)} X(T),$$

to denote the T_p -forward FX rate seen at time T .

As we have seen previously for CMS-linked cash flows, the (domestic) annuity measure provides the most direct market information about the distribution of the underlying swap rate. Switching to this measure, we have

$$V_{\text{QuantoCMS}}(0) = \frac{A(0)}{X(0)} \mathbb{E}^{A,d} \left(g(S(T)) \frac{P_d(T, T_p)}{A(T)} X_{T_p}(T) \right).$$

Drawing on the results of Section 16.6, we see that if the payment currency of the cash flow were domestic, then the appropriate valuation formula instead would read

$$V(0) = A(0) \mathbb{E}^{A,d} \left(g(S(T)) \frac{P_d(T, T_p)}{A(T)} \right).$$

The *quanto adjustment* is defined to be the ratio,

$$D_{\text{Quanto}}(0) = \frac{\mathbb{E}^{A,d} \left(g(S(T)) \frac{P_d(T, T_p)}{A(T)} X_{T_p}(T) \right)}{X(0) \mathbb{E}^{A,d} \left(g(S(T)) \frac{P_d(T, T_p)}{A(T)} \right)}.$$

For quanto CMS valuation or, equivalently, calculation of quanto adjustments, it is natural to search for a suitable extension of the methods developed previously for (single-currency) CMS-linked cash flows. In particular, as quanto CMS structures are quite vanilla-like, we would ideally like to avoid the usage of full term structure (multi-currency) interest rate models.

By the arguments similar to those of Section 16.6.2, we have

$$V_{\text{QuantoCMS}}(0) = \frac{A(0)}{X(0)} \mathbb{E}^{A,d} (g(S(T)) v(S(T))), \quad (16.88)$$

where

$$v(s) \triangleq \mathbb{E}^{A,d} \left(\frac{P_d(T, T_p)}{A(T)} X_{T_p}(T) \middle| S(T) = s \right).$$

Let us recall the definition (16.50) of $\alpha(s)$ and also define

$$\chi(s) \triangleq \mathbb{E}^{A,d} (X_{T_p}(T) | S(T) = s). \quad (16.89)$$

Then, approximately¹²,

$$v(s) \approx \alpha(s) \chi(s),$$

so that (also approximately)

$$V_{\text{QuantoCMS}}(0) = \frac{A(0)}{X(0)} \mathbb{E}^{A,d} (g(S(T)) \alpha(S(T)) \chi(S(T))). \quad (16.90)$$

Once the value is represented in the form (16.90) it can be computed by the replication method, as in Section 16.6.1. To complete the valuation, it only remains to determine the function $\chi(s)$ in (16.89).

¹²Here we essentially assume that the slope of the yield curve is independent of the “pure” FX component of the forward FX rate.

16.7.2 Modeling the Joint Distribution of Swap Rate and Forward Exchange Rate

To compute the function $\chi(s)$ in the previous section, a joint distribution of the swap rate $S(T)$ and the forward FX rate $X_{T_p}(T)$ needs to be specified. The marginal one-dimensional distribution of $S(T)$ in $Q^{A,d}$ is given by the swaption model; we denote the cumulative distribution function (CDF) of $S(T)$ by $\Psi^A(s)$, see (16.66). The payoff of the quanto CMS cash flow in (16.88) depends on $X_{T_p}(T)$ linearly, indicating that the value of the derivative has rather limited dependence on the particular form of the distribution of the FX rate, a fact well-supported by numerical tests¹³. In the interest of analytical tractability, we simply model $X_{T_p}(T)$ as being log-normal in the domestic annuity measure $Q^{A,d}$, i.e. we assume that there exists a standard Gaussian random variable ξ_1 , a volatility σ_X , and a scaling constant m_X such that

$$X_{T_p}(T) = X(0)e^{\sigma_X \sqrt{T} \xi_1 + m_X T}. \quad (16.91)$$

The volatility σ_X is obtained by calibrating (16.91) to T -expiry ATM options on the FX rate¹⁴, whereas the choice of the constant m_X is clarified below.

With marginal distributions of $S(T)$ and $X_{T_p}(T)$ specified, we impose the dependence structure with a simple application of the so-called *copula method*. Chapter 17 contains a full review of copula methods and their applications to multi-rate vanilla derivatives, but for our needs here it suffices to note that if ξ_2 is a standard Gaussian random variable, then clearly

$$S(T) \stackrel{d}{=} (\Psi^A)^{-1}(\Phi(\xi_2)),$$

where $\Phi(\cdot)$ is the standard Gaussian CDF, and the equality is in terms of distribution. The dependence structure between $S(T)$ and $X_{T_p}(T)$ may now be imposed by correlating the two standard Gaussian random variables, ξ_1 and ξ_2 , with a correlation ρ_{XS} , leading to the following specification of the joint distribution

$$\begin{aligned} X_{T_p}(T) &= X(0)e^{\sigma_X \sqrt{T} \xi_1 + m_X T}, & S(T) &= (\Psi^A)^{-1}(\Phi(\xi_2)), \\ \text{Corr}(\xi_1, \xi_2) &= \rho_{XS}. \end{aligned}$$

The function $\chi(s)$ in (16.89) can now easily be computed,

¹³For very long-dated quanto contracts, the FX volatility smile may start to matter. The method we develop here can be extended to incorporate the FX smile, by techniques discussed in Chapter 17.

¹⁴This is not exact as the T -expiry FX option is written on $X_T(T)$, not $X_{T_p}(T)$. The difference is rarely material as $T_p - T$ is often small, but more elaborate schemes are not difficult to design, if desired.

$$\begin{aligned}
\chi(s) &= \mathbb{E}^{A,d}(X_{T_p}(T) \mid S(T) = s) \\
&= X(0)e^{m_X T} \mathbb{E}^{A,d}\left(e^{\sigma_X \sqrt{T} \xi_1} \mid \xi_2 = \Phi^{-1}(\Psi^A(s))\right) \\
&= X(0)e^{m_X T} \tilde{\chi}(s), \\
\tilde{\chi}(s) &= \exp\left(\rho_{XS}\sigma_X \sqrt{T}\Phi^{-1}(\Psi^A(s)) + \frac{\sigma_X^2 T}{2}(1 - \rho_{XS}^2)\right).
\end{aligned} \tag{16.92}$$

16.7.3 Normalizing Constant and Final Formula

To complete the development of a pricing formula for quanto CMS cash flows, we now only need to establish the constant m_X in (16.91) (and (16.92)). For this, we note that $X_{T_p}(\cdot)$ is a martingale in the domestic T_p -forward measure $Q^{T_p,d}$ (see Section 4.3), and in particular,

$$X_{T_p}(0) = \mathbb{E}^{T_p,d}(X_{T_p}(T)).$$

Changing to the domestic annuity measure $Q^{A,d}$, the following holds,

$$X_{T_p}(0) = \frac{A(0)}{P_d(0, T_p)} \mathbb{E}^{A,d}\left(\frac{P_d(T, T_p)}{A(T)} X_{T_p}(T)\right).$$

Recalling the definition of $\alpha(s)$ and $\chi(s)$ we finally write

$$\begin{aligned}
X_{T_p}(0) &= \frac{A(0)}{P_d(0, T_p)} \mathbb{E}^{A,d}(\alpha(S(T))\chi(S(T))) \\
&= X(0)e^{m_X T} \frac{A(0)}{P_d(0, T_p)} \mathbb{E}^{A,d}(\alpha(S(T))\tilde{\chi}(S(T))),
\end{aligned}$$

and hence

$$e^{-m_X T} = \frac{X(0)}{X_{T_p}(0)} \frac{A(0)}{P_d(0, T_p)} \mathbb{E}^{A,d}(\alpha(S(T))\tilde{\chi}(S(T))).$$

Combining all previous results, we have arrived at the following pricing formula.

Proposition 16.7.1. *Let the forward FX rate $X_{T_p}(T)$ be log-normal with volatility σ_X , and let the co-dependence between $X_{T_p}(T)$ and the swap rate $S(T)$ be characterized by a Gaussian copula with correlation ρ_{XS} . The value of a quanto CMS cash flow $g(S(T))$ paid in a foreign currency at time $T_p \geq T$ is then approximated by*

$$V_{\text{QuantoCMS}}(0) \approx P_f(0, T_p) \frac{\mathbb{E}^{A,d}(g(S(T))\alpha(S(T))\tilde{\chi}(S(T)))}{\mathbb{E}^{A,d}(\alpha(S(T))\tilde{\chi}(S(T)))}, \tag{16.93}$$

where the annuity mapping function $\alpha(\cdot)$ is defined by (16.50), and

$$\tilde{\chi}(s) = \exp\left(\rho_{XS}\sigma_X \sqrt{T}\Phi^{-1}(\Psi^A(s)) + \frac{\sigma_X^2 T}{2}(1 - \rho_{XS}^2)\right),$$

with $\Psi^A(s) \triangleq \mathbb{P}^{A,d}(S(T) < s)$.

Remark 16.7.2. The expected values in the denominator and the numerator of the right-hand side of (16.93) can be computed by the replication method from Section 16.6.1.

16.8 Eurodollar Futures

Eurodollar (ED) futures are exchange-traded futures contracts on Libor rates (see Section 5.4) and serve as fundamental inputs in the construction of the yield curve, as explained in Chapter 6. As we explained in Section 4.1.2, daily mark-to-market causes the value of the futures contract on a Libor rate to differ from the value of a forward contract on a Libor rate. Only the latter is an input into an interest rate curve construction, while only the former is liquidly quoted in the market. The difference between the two is called the *ED future convexity adjustment*, a quantity that we shall analyze and quantify in the following sections.

16.8.1 Fundamental Results on Futures

As in Section 4.1.2, we let $F(t, T, M)$ denote the futures rate at time t covering the period $[T, M]$, with $0 \leq t \leq T < M$. The forward (Libor) rate for the same period is, as always, denoted by $L(t, T, M)$. The next lemma establishes the relationship between the two.

Lemma 16.8.1. *Let Q be the risk-neutral measure and Q^M be the M -forward measure, with \mathbb{E} and \mathbb{E}^M being the corresponding expected value operators. Then*

$$L(t, T, M) = \mathbb{E}_t^M(L(T, T, M)), \quad F(t, T, M) = \mathbb{E}_t(L(T, T, M)).$$

Proof. The first result is from Lemma 4.2.3 and the second is from Lemma 4.2.2. \square

Lemma 16.8.1 holds for futures contracts that are marked-to-market continuously. For calculation purposes it is often more convenient to assume discrete mark-to-market; it has been established (see Hunt and Kennedy [2000]) that prices with monthly or even quarterly resettlement frequency differ little from the prices with continuous (or daily) resettlement. To work out a version of Lemma 16.8.1 that covers discrete mark-to-market, let us introduce a standard tenor structure

$$0 = T_0 < T_1 < \dots < T_N, \quad \tau_n = T_{n+1} - T_n,$$

and define Libor rates as before,

$$L_n(t) = L(t, T_n, T_{n+1}) = \frac{P(t, T_n) - P(t, T_{n+1})}{\tau_n P(t, T_{n+1})}, \quad n = 0, \dots, N - 1.$$

The discretely compounded money market account $B(t)$ is defined by (4.24), and the corresponding measure, the spot Libor measure, is defined in Section 4.2.3; we denote it by Q^B and the corresponding expected value operator by E^B . We abbreviate the notation for the expected value in the T_n -forward measure to $E^n \triangleq E^{T_n}$ (and the same for variance) and, finally, in line with the definition of spanning Libor rates, define spanning futures rates

$$F_n(t) = F(t, T_n, T_{n+1}), \quad n = 0, \dots, N - 1.$$

We are now ready for the valuation formula of discretely-resettled futures rates.

Proposition 16.8.2. *The futures rate that is marked-to-market only on the dates $T_0 < T_1 < \dots < T_n = T < M$ is given by*

$$F(t, T, M) = E_t^B(L(T, T, M)).$$

In particular,

$$E^B(L_n(T_n)) = F_n(0), \quad n = 0, \dots, N - 1.$$

Proof. At time $T_n = T$, the cash flow associated with the futures contract is

$$F(T_n, T, M) - F(T_{n-1}, T, M) = L(T, T, M) - F(T_{n-1}, T, M).$$

At time T_{n-1} , the present value of this cash flow is

$$V_{\text{fut}}(T_{n-1}) = B(T_{n-1}) E_{T_{n-1}}^B \left(\frac{L(T, T, M) - F(T_{n-1}, T, M)}{B(T_n)} \right).$$

By the definition of the rolling spot numeraire $B(t)$, the quantity

$$B(T_{n-1})/B(T_n) = P(T_{n-1}, T_n)$$

is non-random at time T_{n-1} , and so is $F(T_{n-1}, T_n, M)$, whereby

$$V_{\text{fut}}(T_{n-1}) = P(T_{n-1}, T_n) \left(E_{T_{n-1}}^B(L(T, T, M)) - F(T_{n-1}, T, M) \right).$$

As futures contracts are always entered into at a price of 0, it follows that $V_{\text{fut}}(T_{n-1}) = 0$ by definition, and therefore

$$F(T_{n-1}, T, M) = E_{T_{n-1}}^B(L(T, T, M)).$$

At time T_{n-2} , we may write

$$\begin{aligned} V_{\text{fut}}(T_{n-2}) &= B(T_{n-2}) E_{T_{n-2}}^B \left(\frac{F(T_{n-1}, T, M) - F(T_{n-2}, T, M)}{B(T_{n-1})} \right) \\ &= B(T_{n-2}) E_{T_{n-2}}^B \left(\frac{E_{T_{n-1}}^B(L(T, T, M)) - F(T_{n-2}, T, M)}{B(T_{n-1})} \right) \\ &= P(T_{n-2}, T_{n-1}) \left(E_{T_{n-2}}^B(L(T, T, M)) - F(T_{n-2}, T, M) \right), \end{aligned}$$

whereby (as $V_{\text{fut}}(T_{n-2}) = 0$)

$$F(T_{n-2}, T, M) = \mathbb{E}_{T_{n-2}}^B(L(T, T, M)).$$

Proceeding inductively, the result follows. \square

The ED future convexity adjustment is given by the difference

$$F(t, T, M) - L(t, T, M) = \mathbb{E}_t^B(L(T, T, M)) - \mathbb{E}_t^M(L(T, T, M))$$

or, for discrete settlement,

$$F(t, T, M) - L(t, T, M) = \mathbb{E}_t^B(L(T, T, M)) - \mathbb{E}_t^M(L(T, T, M)).$$

Proposition 4.5.3 derived the convexity adjustment in a general Gaussian multi-factor HJM model in closed-form. However, having demonstrated the importance of incorporating volatility smile in calculations for other types of convexity (Libor-in-arrears, CMS), we can legitimately ask whether the smile has a significant impact on ED convexity as well. This question cannot be answered within the constraints of a Gaussian model as it does not allow for smile control; we instead follow the ideas in Piterbarg and Renedo [2006] and develop a smile-enabled pricing approach in the following.

16.8.2 Motivations and Plan

Performance requirements for valuing ED futures are even more stringent than for other types of derivatives, due to their high trading volumes and, in particular, their use in yield curve construction. This rules out Monte Carlo methods, or even PDE-based schemes, necessitating the development of analytic approximations that incorporate the volatility smile yet allow for efficient numerical algorithms. In addition, we look for the formula for convexity adjustments that depends on observable market inputs in the most direct way possible, with lengthy model and curve calibrations reduced to a minimum or eliminated altogether. To achieve this, we separate model parameters into two categories: those that change often, and those that do not. The former category here covers volatility parameters, and are taken directly from the prices of options on ED contracts across different expiries and strikes. These parameters can be updated in real time as we build yield curves intra-day. The latter category of (slow-moving) parameters are essentially correlation parameters, and originate from calibrating a model with a rich volatility structure to caps and swaptions. Due to computational constraints, these parameters cannot be updated often — but they do not need to be, as they typically do not change much over time.

We use the following road map to derive the ED futures valuation formula:

1. First, an expansion technique is applied to derive a model-independent relationship that expresses a *forward rate* as a functional of a *collection of futures* rates with expiries on or before the expiry of the forward rate.

2. The variance terms that appear in the formula are separated into slow- and fast-moving parameters, as described earlier.
3. Fast-moving volatility parameters are represented in several different ways, both in model-independent fashion as functions of prices of options on ED futures across strikes¹⁵, and as closed-form expressions involving volatility smile parameters.
4. Finally, slow-moving correlation parameters are expressed in terms of the parameters of a Libor market model, properly calibrated to relevant market instruments or, more pragmatically, through the simplified formulas (16.117) or (16.118).

We emphasize that the first step of the algorithm differs from traditional methods which typically express the value of futures rates in terms of forward rates, an approach diametrically opposite of ours. We consider our approach superior, as it eliminates the need to invert equations to obtain forward rates from market-observed futures quotes (a fundamental requirement of curve building algorithms).

16.8.3 Preliminaries

As established previously, the following relations hold,

$$F_n(0) = \mathbb{E}^B(L_n(T_n)) = \mathbb{E}^B(F_n(T_n)), \quad (16.94)$$

$$L_n(0) = \mathbb{E}^{n+1}(L_n(T_n)) = \mathbb{E}^{n+1}(F_n(T_n)), \quad (16.95)$$

for all $n = 0, \dots, N - 1$. We assume that all $F_n(0)$, $n = 0, \dots, N - 1$, are known; our goal is to derive formulas that express forward rates $\{L_n(0)\}$ in terms of futures $\{F_n(0)\}$ and, potentially, other market-observed quantities. The following result is straightforward.

Lemma 16.8.3. *For each n , $n = 0, \dots, N - 1$,*

$$L_n(0) = \mathbb{E}^B \left(\left[\prod_{i=0}^n \frac{1 + \tau_i L_i(0)}{1 + \tau_i L_i(T_i)} \right] L_n(T_n) \right). \quad (16.96)$$

Proof. Follows by a measure change. \square

Lemma 16.8.3 expresses the forward rate $L_n(0)$ as an expectation of a certain payoff in the spot Libor measure (not forward measure as in (16.95)), the measure that is used in defining futures in (16.94). This result is used as a starting point for deriving convexity adjustments.

¹⁵This is similar to the replication method for computing CMS and Libor-in-arrears convexity adjustments.

16.8.4 Expansion Around the Futures Value

To express the expected value in (16.96) in terms of market-observed quantities, we derive a Taylor series expansion of (16.96) in powers of a small parameter that measures the deviation of each of $L_n(T_n)$ from its mean in the spot Libor measure, $F_n(0) = \mathbb{E}^B(L_n(T_n))$.

Fix $\epsilon > 0$, and define L_n^ϵ 's by

$$L_n^\epsilon(t) = \epsilon(L_n(t) - F_n(0)) + F_n(0).$$

Note that for any n ,

$$L_n^1(t) = L_n(t), \quad (16.97)$$

$$L_n^0(t) = F_n(0), \quad (16.98)$$

$$\frac{\partial L_n^\epsilon(t)}{\partial \epsilon} = L_n(t) - F_n(0), \quad (16.99)$$

$$L_n^\epsilon(T_n) = \epsilon(L_n(T_n) - \mathbb{E}^B(L_n(T_n))) + \mathbb{E}^B(L_n(T_n)). \quad (16.100)$$

Define

$$V(\epsilon) = \left[\prod_{i=0}^n \frac{1 + \tau_i L_i(0)}{1 + \tau_i L_i^\epsilon(T_i)} \right] L_n^\epsilon(T_n). \quad (16.101)$$

It should be clear that $V(1)$ is the value inside the expectation on the right-hand side of (16.96),

$$L_n(0) = \mathbb{E}^B(V(1)). \quad (16.102)$$

Expanding $V(\epsilon)$ into a Taylor series in ϵ yields,

$$V(\epsilon) = V(0) + \mathbb{E}^B \left(\frac{dV}{d\epsilon}(0) \right) \times \epsilon + \frac{1}{2} \mathbb{E}^B \left(\frac{d^2V}{d\epsilon^2}(0) \right) \times \epsilon^2 + O(\epsilon^3). \quad (16.103)$$

The values of the derivatives of $V(\epsilon)$ are computed in the following lemma.

Lemma 16.8.4. *For any n , $n = 0, \dots, N - 1$,*

$$V(0) = \left[\prod_{i=0}^n \frac{1 + \tau_i L_i(0)}{1 + \tau_i F_i(0)} \right] F_n(0), \quad (16.104)$$

$$\mathbb{E}^B \left(\frac{dV}{d\epsilon}(0) \right) = 0, \quad (16.105)$$

$$\mathbb{E}^B \left(\frac{d^2V}{d\epsilon^2}(0) \right) = V(0) \sum_{j,m=0}^n D_{j,m} \text{Cov}^B(L_j(T_j), L_m(T_m)), \quad (16.106)$$

where the coefficients $D_{j,m}$ are given by

$$\begin{aligned} D_{j,m} &= \left(-\frac{\tau_j}{1 + \tau_j F_j(0)} + \frac{1_{\{j=n\}}}{F_n(0)} \right) \left(-\frac{\tau_m}{1 + \tau_m F_m(0)} + \frac{1_{\{m=n\}}}{F_n(0)} \right) \\ &\quad + 1_{\{j=m\}} \left(\frac{\tau_j^2}{(1 + \tau_j F_j(0))^2} - \frac{1_{\{j=n\}}}{F_n(0)^2} \right), \quad (16.107) \end{aligned}$$

and, by definition,

$$\begin{aligned} \text{Cov}^B(L_j(T_j), L_m(T_m)) \\ = \mathbb{E}^B [(L_j(T_j) - \mathbb{E}^B(L_j(T_j))) (L_m(T_m) - \mathbb{E}^B(L_m(T_m)))] \\ = \mathbb{E}^B [(F_j(T_j) - F_j(0))(F_m(T_m) - F_m(0))]. \end{aligned}$$

Proof. It follows from (16.101) that

$$V(\epsilon) = \left(\prod_{i=0}^n (1 + \tau_i L_i(0)) \right) p(L_0^\epsilon(T_0), \dots, L_n^\epsilon(T_n)),$$

where we defined

$$p(y_0, \dots, y_n) = \left[\prod_{i=0}^n \frac{1}{1 + \tau_i y_i} \right] y_n. \quad (16.108)$$

Obviously, (16.104) follows from (16.98). Moreover, with the help of (16.99),

$$\begin{aligned} \frac{dV(\epsilon)}{d\epsilon} &= \left(\prod_{i=0}^n (1 + \tau_i L_i(0)) \right) \\ &\times \sum_{j=0}^n \frac{\partial}{\partial y_j} p(L_0^\epsilon(T_0), \dots, L_n^\epsilon(T_n)) (L_j(t) - F_j(0)), \quad (16.109) \end{aligned}$$

Since

$$\mathbb{E}^B(L_j(t) - F_j(0)) = 0,$$

the statement (16.105) is proved.

Differentiating (16.109) with respect to ϵ again, we obtain

$$\begin{aligned} \frac{d^2V(0)}{d\epsilon^2} &= \left(\prod_{i=0}^n (1 + \tau_i L_i(0)) \right) p(F_0(0), \dots, F_n(0)) \\ &\times \sum_{j,m} D_{j,m} (L_j(t) - F_j(0)) (L_m(t) - F_m(0)), \end{aligned}$$

where we have denoted

$$D_{j,m} = \frac{\partial^2}{\partial y_j \partial y_m} p(F_0(0), \dots, F_n(0)).$$

The expression (16.107) for $D_{j,m}$'s follows by calculating $\partial^2 p / \partial y_j \partial y_m$ from (16.108). Simplifying, we obtain

$$\frac{d^2 V(0)}{d\epsilon^2} = V(0) \sum_{j,m} D_{j,m} (L_j(t) - F_j(0)) (L_m(t) - F_m(0)).$$

Taking the expected value of both sides leads to (16.105). Full details of the proof are in Piterbarg and Renedo [2006]. \square

Lemma 16.8.4 expresses quantities $V(0)$, $E^B(\partial^2 V(0)/\partial\epsilon^2)$ in the series expansion (16.103) in terms of quantities that are either directly observable, such as futures rates, or computable, such as covariances of forward rates. Applying the results of Lemma 16.8.4 to the representation (16.102) and expansion (16.103), we obtain the following result.

Theorem 16.8.5. *For any n , $n = 0, \dots, N-1$, an approximation to forward rate $L_n(0)$ is obtained from the futures $\{F_i(0)\}_{i=0}^n$ and forward rates for previous periods $\{L_i(0)\}_{i=0}^{n-1}$ by solving the following equation,*

$$L_n(0) = V(0) \left(1 + \frac{1}{2} \sum_{j,m=0}^n D_{j,m} \text{Cov}^B(L_j(T_j), L_m(T_m)) \right), \quad (16.110)$$

with $V(0)$ and $D_{j,m}$ given in Lemma 16.8.4.

Remark 16.8.6. Theorem 16.8.5 specifies an algorithm for solving for forward rates $L_n(0)$ sequentially for all n , using futures prices $\{F_j(0)\}$ as inputs.

Remark 16.8.7. The expression on the right-hand side of (16.110) will be simplified in the sections that follow. In many cases the rate to be determined from the expression, $L_n(0)$, will appear on the right-hand side of (16.110) as well. In this case, (16.110) should be treated not as an identity, but as an equation on $L_n(0)$. While this may seem to complicate the problem of finding $L_n(0)$, in reality the dependence of the right-hand side of (16.110) on $L_n(0)$ is typically mild, and the equation can be solved iteratively in just a few steps.

The formula (16.110) depends on covariances between various forward rates. By the definition of the covariance,

$$\begin{aligned} & \text{Cov}^B(L_j(T_j), L_m(T_m)) \\ &= (\text{Var}^B(L_j(T_j)) \text{Var}^B(L_m(T_m)))^{1/2} \text{Corr}^B(L_j(T_j), L_m(T_m)), \end{aligned} \quad (16.111)$$

where the variances and the correlation are computed in the spot measure. We proceed to discuss how to estimate the terms on the right-hand side of (16.111) from market observations.

16.8.5 Forward Rate Variances

The variance of $L_n(T_n)$ in the formula (16.111) is to be computed in the spot measure. As an approximation, we compute the variance instead in the measure in which $L_n(t)$ is a martingale,

$$\text{Var}^B(L_n(T_n)) \approx \text{Var}^{n+1}(L_n(T_n)), \quad n = 1, \dots, N-1, \quad (16.112)$$

where by definition

$$\text{Var}^{n+1}(L_n(T_n)) \triangleq E^{n+1}(L_n(T_n) - E^{n+1}L_n(T_n))^2.$$

The error of this approximation is typically small, and allows us to rewrite the formula for computing forward rates from futures rates as

$$\begin{aligned} L_n(0) &\approx V(0) \left(1 + \frac{1}{2} \sum_{j,m=0}^n D_{j,m} \text{Var}^{j+1}(L_j(T_j))^{1/2} \right. \\ &\quad \times \left. \text{Var}^{m+1}(L_m(T_m))^{1/2} \text{Corr}^B(L_j(T_j), L_m(T_m)) \right). \end{aligned} \quad (16.113)$$

The market in options on ED futures contracts is very liquid — perhaps the most liquid market of options on interest rates. Applying the familiar replication method allows estimation of the variance of a forward rate in a model-independent way from observable prices of ED futures options (compare to (16.38)):

$$\begin{aligned} \text{Var}^{n+1}(L_n(T_n)) &= E^{n+1}((L_n(T_n) - E^{n+1}L_n(T_n))^2) \\ &= 2 \int_{-\infty}^{L_n(0)} E^{n+1}((K - L_n(T_n))^+) dK \\ &\quad + 2 \int_{L_n(0)}^{\infty} E^{n+1}((L_n(T_n) - K)^+) dK. \end{aligned} \quad (16.114)$$

In the formula (16.113), observable option prices are used directly for variance calculations and the forward rate $L_n(0)$, the rate to solve for, enters the right-hand side of the equation only as an integration limit in (16.114), so the comments of Remark 16.8.7 still apply. Equation (16.114) demonstrates explicitly that the ED convexity adjustment depends on prices of ED futures options at all strikes, i.e. on the volatility smile.

Equation (16.114) may not be easy to use in practice as only a discrete set of strikes is typically traded, and not all of them are very liquid. For these reasons, we may wish to capture the smile by a low-parametric vanilla model — or perhaps just a functional form, as in Section 16.1.5 — calibrated to liquid strikes. For instance, we could use a standard stochastic volatility model¹⁶ for the rate $L_n(t)$,

¹⁶We use σ instead of our customary λ to avoid notational conflict with LM model volatilities used later on.

$$\begin{aligned} dL_n(t) &= \sigma_n (b_n L_n(t) + (1 - b_n) L_n(0)) \sqrt{z(t)} dW(t), \\ dz(t) &= \theta (1 - z(t)) dt + \eta_n \sqrt{z(t)} dZ(t), \end{aligned} \quad (16.115)$$

with correlation $\langle dW(t), dZ(t) \rangle = 0$. These SDEs are understood to be in the T_{n+1} -forward measure. The set of parameters (σ_n, b_n, η_n) defines the volatility smile in options on the rate $L_n(T_n)$ and is calibrated to market as described in Section 16.1.4. The variance of the Libor rate in the model (16.115) can easily be calculated:

Proposition 16.8.8. *Recall the definition (8.11) of $\Psi_{\bar{z}}(v, u; t)$. Then*

$$\text{Var}^{n+1}(L_n(T_n)) = \frac{L_n(0)^2}{b_n^2} \left[\Psi_{\bar{z}} \left((\sigma_n b_n)^2, 0; T_n \right) - 1 \right].$$

Proof. Direct calculations. \square

We again comment that the expression for the variance $\text{Var}^{n+1}(L_n(T_n))$ involves $L_n(0)$, which here makes (16.113) a quadratic equation in $L_n(0)$. In addition, it should be noted that the implied values of parameters (σ_n, b_n, η_n) also, in principle, depend on the parameter $L_n(0)$ through the calibration process. However, the loss of accuracy is negligible if (σ_n, b_n, η_n) are calibrated with the “previous” value of the forward rate $L_n(0)$, i.e. the value before the update of the convexity adjustment.

16.8.6 Forward Rate Correlations

Once forward rate variances have been computed, the computation of (16.113) can be completed provided that we can establish the correlations $\text{Corr}^B(L_j(T_j), L_m(T_m))$. There are many ways this can be done, but some type of model assumption will generally be required. For instance, if we have a calibrated LM model (16.59) lying around, we may calculate correlations in the Libor market model from the formula (14.35). Specifically, assuming $T_j \leq T_m$, we have for the model (16.59),

$$\text{Corr}^B(L_j(T_j), L_m(T_m)) = \frac{\int_0^{T_j} \lambda_j(s)^\top \lambda_m(s) ds}{\left(\int_0^{T_j} \|\lambda_j(s)\|^2 ds \right)^{1/2} \left(\int_0^{T_m} \|\lambda_m(s)\|^2 ds \right)^{1/2}}. \quad (16.116)$$

Extraction of the model parameters $\{\lambda_j(\cdot)\}$ from the market is described in Chapter 14. Since correlations do not change often, this calibration can be performed “off-line”, i.e. on an infrequent basis with the results reused as needed.

The approach above assumes that a full LM model has been implemented and calibrated. This may not always be practical, so let us consider a simplified method that retains the general spirit of a full LM model. First, we assume that the dynamics of Libor rates originate from a time-stationary Gaussian process of the mean-reverting type,

$$dL_i(t) = O(dt) + \sigma_0 e^{-\kappa(T_i - t)} dW_i(t), \quad i = j, m,$$

where $W_j(t)$ and $W_m(t)$ are scalar Q^B -Brownian motions with correlation

$$\langle dW_j(t), dW_m(t) \rangle = q(T_j - t, T_m - t) dt,$$

for some function $q : \mathbb{R}^2 \rightarrow [-1, 1]$. Representative examples of the correlation function q are listed in Section 14.3.2. Ignoring drift terms and assuming $T_j \leq T_m$, we have

$$\begin{aligned} \text{Cov}^B(L_j(T_j), L_m(T_m)) &= \sigma_0^2 \int_0^{T_j} e^{-\kappa(T_j - t)} e^{-\kappa(T_m - t)} q(T_j - t, T_m - t) dt, \\ \text{Var}^B(L_i(T_i)) &= \frac{\sigma_0^2}{2\kappa} (1 - e^{-2\kappa T_i}), \quad i = j, m. \end{aligned}$$

Hence

$$\text{Corr}^B(L_j(T_j), L_m(T_m)) = 2\kappa \frac{\int_0^{T_j} e^{-\kappa(T_j - t)} e^{-\kappa(T_m - t)} q(T_j - t, T_m - t) dt}{\sqrt{(1 - e^{-2\kappa T_j})(1 - e^{-2\kappa T_m})}}. \quad (16.117)$$

A special case arises when $q(T_j - t, T_m - t) = \rho_{j,m}$ and does not depend on t , in which case (16.117) reduces to

$$\begin{aligned} \text{Corr}^B(L_j(T_j), L_m(T_m)) \\ = \rho_{j,m} \frac{(e^{2\kappa T_j} - 1) e^{-\kappa(T_j + T_m)}}{\sqrt{(1 - e^{-2\kappa T_j})(1 - e^{-2\kappa T_m})}} = \rho_{j,m} \sqrt{\frac{e^{2\kappa T_j} - 1}{e^{2\kappa T_m} - 1}}. \quad (16.118) \end{aligned}$$

Formulas (16.117) and (16.118) do not require a full calibration of an LM model, only the estimation of basic forward rate correlations and a single mean reversion. We note that the role of the latter quantity is to govern the amount of de-correlation caused by the fact that L_j and L_m fix at different time. The mean reversion can potentially be estimated from market data (see Section 13.1.8), or could be a direct trader input.

16.8.7 The Formula

For convenience, let us now pull all previous results together into a single, easily referenced formula. First, let us summarize the notations. By $\{L_n(0)\}_{n=1}^{N-1}$ we denote the (unknown) sequence of forward rates for the tenor structure $\{T_n\}_{n=0}^N$, and by $\{F_n(0)\}_{n=1}^{N-1}$ we denote the (known) sequence of futures rates. For each n , $n = 1, \dots, N - 1$, let the triple (σ_n, b_n, η_n) be the set of parameters of the model (16.115) implied from market prices of options on the rate $L_n(T_n)$ of different strikes.

Theorem 16.8.9. For each n , $n = 0, \dots, N - 1$, the forward rate $L_n(0)$ is obtained from the futures rates $\{F_i(0)\}_{i=0}^n$ and forward rates for previous periods $\{L_i(0)\}_{i=0}^{n-1}$ by solving the following equation,

$$\begin{aligned} L_n(0) &= V(0) \left(1 + \frac{1}{2} \sum_{j,m=0}^n D_{j,m} \frac{L_j(0)L_m(0)}{b_j b_m} \right. \\ &\quad \times \left. \left(\Psi_{\bar{z}} \left((\sigma_j b_j)^2, 0; T_j \right) - 1 \right)^{1/2} \left(\Psi_{\bar{z}} \left((\sigma_m b_m)^2, 0; T_m \right) - 1 \right)^{1/2} c_{j,m} \right), \end{aligned} \quad (16.119)$$

with

$$V(0) = \left[\prod_{i=0}^n \frac{1 + \tau_i L_i(0)}{1 + \tau_i F_i(0)} \right] F_n(0),$$

$$\begin{aligned} D_{j,m} &= \left(-\frac{\tau_j}{1 + \tau_j F_j(0)} + \frac{1_{\{j=n\}}}{F_n(0)} \right) \left(-\frac{\tau_m}{1 + \tau_m F_m(0)} + \frac{1_{\{m=n\}}}{F_n(0)} \right) \\ &\quad + 1_{\{j=m\}} \left(\frac{\tau_j^2}{(1 + \tau_j F_j(0))^2} - \frac{1_{\{j=n\}}}{F_n(0)^2} \right), \end{aligned}$$

and

$$c_{j,m} = \text{Corr}^B(L_j(T_j), L_m(T_m))$$

as given by, for example, (16.116), (16.117) or (16.118). The function $\Psi_{\bar{z}}(v, u; t)$ is defined by (8.11) and is available in closed form per Proposition 8.3.8.

16.9 Convexity and Moment Explosions

When dealing with convexity products — Libor-in-arrears, CMS, ED futures, and so forth — we find (equations (16.36), (16.54) and (16.110)) that their values depend on the variance of some underlying rate, i.e. a *second moment* of the appropriate terminal distribution. Some care must be taken in the model setup to ensure that this second moment is well-behaved. For instance, if we have elected to work in a stochastic volatility setup, Proposition 8.3.10 and the discussion around it tell us that the second moment of the underlying in a stochastic volatility model may fail to exist, even for fairly reasonable values of model parameters. Should that occur, the theoretical convexity value will become infinite.

Intuitively, convexity value depends on prices of options at a continuum of strikes, from 0 to $+\infty$. In the market however, only prices of options over a finite range of strikes are observed, and infinite prices arise solely from the

model-based extrapolation of the volatility smile beyond the observable range. We can control the smile extrapolation by altering the model specification in the manner discussed after Proposition 8.3.10. Alternatively, at least for pricing vanilla derivatives, we can control smile extrapolation explicitly, e.g. by restricting the domain of integration in the replication method, as we have already discussed in Sections 16.4 and 16.6.1. In particular, when evaluating the variance of some rate $S(t)$ we would replace

$$\text{Var}(S(T)) = 2 \int_{-\infty}^{S(0)} \mathbb{E}(K - S(T))^+ dK + 2 \int_{S(0)}^{\infty} \mathbb{E}(S(T) - K)^+ dK$$

with

$$\text{Var}(S(T)) \approx 2 \int_{K_{\min}}^{S(0)} \mathbb{E}(K - S(T))^+ dK + 2 \int_{S(0)}^{K_{\max}} \mathbb{E}(S(T) - K)^+ dK$$

for some $-\infty < K_{\min} \leq S(0) \leq K_{\max} < \infty$, as justified by the fact that only options with market-observable strikes can be used in hedging. The same idea can be applied to all convexity products evaluated by the replication method. The extra parameters, K_{\min} and K_{\max} , could even be used to calibrate CMS (and other) convexity values to market, if these market values are available.

An alternative to outright cropping of the integration domain would be to institute an ad-hoc modification of the model-implied density of $S(T)$ for small/large values of the swap rate. This can be accomplished in many different ways, e.g. by artificially flattening out the implied volatility smile for large strikes¹⁷. The particulars of this scheme are left to the reader to explore, but let us note that any scheme to flatten out the implied volatility smile should, of course, be smooth as a function of strikes, to avoid generation of negative densities.

If we wish to control smile wing behavior in such a way that we are always certain that a valid density arises, we can also contemplate ad-hoc measures to modify model densities directly to prevent moment explosion (or to otherwise regularize the model). Assume that we have implemented a model where the density for $S(T)$ in its annuity measure \mathbb{Q}^A is $\psi(s)$. A call payout (as needed for a payer swaption) is therefore valued as

$$\mathbb{E}^A((S(T) - K)^+) = \int_K^{\infty} (s - K) \psi(s) ds.$$

Let us introduce some user-specified strikes K_{\min} and K_{\max} , and rewrite the expectation as (assuming $K_{\min} < K_{\max}$)

¹⁷We note in passing that this technique can also be used to control errors in volatility expansion formulas, such as the one used in SABR, which may yield negative densities in the wings unless some kind of regularization is performed.

$$\begin{aligned}
& \mathbb{E}^A \left((S(T) - K)^+ \right) \\
&= \int_{K_{\min}}^{K_{\max}} (s - K)^+ \psi(s) ds \\
&\quad + \int_{-\infty}^{K_{\min}} (s - K)^+ \psi(s) ds + \int_{K_{\max}}^{\infty} (s - K)^+ \psi(s) ds \\
&= \int_{K_{\min}}^{K_{\max}} (s - K)^+ \psi(s) ds \\
&\quad + Q^A(S(T) \leq K_{\min}) \int_{-\infty}^{K_{\min}} (s - K)^+ \psi(s | S(T) \leq K_{\min}) ds \\
&\quad + Q^A(S(T) > K_{\max}) \int_{K_{\max}}^{\infty} (s - K)^+ \psi(s | S(T) > K_{\max}) ds.
\end{aligned}$$

Here we have introduced conditional densities

$$\begin{aligned}
\psi(s | S(T) \leq K_{\min}) ds &= Q^A(S(T) \in [s, s + ds] | S(T) \leq K_{\min}), \\
\psi(s | S(T) > K_{\max}) ds &= Q^A(S(T) \in [s, s + ds] | S(T) > K_{\max}).
\end{aligned}$$

Suppose now that we wish to replace the density ψ with another density in the tails, i.e. for values of $S(T)$ outside of the range $[K_{\min}, K_{\max}]$. To do this, let us write

$$\begin{aligned}
& \mathbb{E}^A \left((S(T) - K)^+ \right) \\
&= \int_{K_{\min}}^{K_{\max}} (s - K)^+ \psi(s) ds \\
&\quad + Q^A(S(T) \leq K_{\min}) \int_{-\infty}^{K_{\min}} (s - K)^+ \psi_{\min}(s | S(T) \leq K_{\min}) ds \\
&\quad + Q^A(S(T) > K_{\max}) \int_{K_{\max}}^{\infty} (s - K)^+ \psi_{\max}(s | S(T) > K_{\max}) ds,
\end{aligned}$$

where we have introduced two new conditional densities ψ_{\min} and ψ_{\max} . Consider now some $K \in [K_{\min}, K_{\max}]$ and let us require that $\mathbb{E}^A((S(T) - K)^+)$ is unchanged by the introduction of ψ_{\min} and ψ_{\max} ; i.e., we insist that the smile in the strike region $[K_{\min}, K_{\max}]$ is preserved after manipulation of the tail densities. Additionally, we of course should demand that $\mathbb{E}^A(S(T)) = S(0)$. The first of these restrictions requires that

$$\begin{aligned}
Q^A(S(T) > K_{\max}) \int_{K_{\max}}^{\infty} (s - K) \psi_{\max}(s | S(T) > K_{\max}) ds \\
= \int_{K_{\max}}^{\infty} (s - K) \psi(s) ds, \quad K \in [K_{\min}, K_{\max}].
\end{aligned}$$

As

$$\int_{K_{\max}}^{\infty} \psi_{\max}(s | S(T) > K_{\max}) ds = 1,$$

it follows that

$$\begin{aligned} Q^A(S(T) > K_{\max}) & \left(\int_{K_{\max}}^{\infty} s \psi_{\max}(s | S(T) > K_{\max}) ds - K \right) \\ &= \int_{K_{\max}}^{\infty} s \psi(s) ds - K Q^A(S(T) > K_{\max}) \end{aligned}$$

or

$$\int_{K_{\max}}^{\infty} s \psi_{\max}(s | S(T) > K_{\max}) ds = \frac{\int_{K_{\max}}^{\infty} s \psi(s) ds}{Q^A(S(T) > K_{\max})}. \quad (16.120)$$

Insisting also that $E^A(S(T))$ is unchanged by the introduction of ψ_{\min} and ψ_{\max} then leads to

$$\int_{-\infty}^{K_{\min}} s \psi_{\min}(s | S(T) \leq K_{\min}) ds = \frac{\int_{-\infty}^{K_{\min}} s \psi(s) ds}{Q^A(S(T) \leq K_{\min})}. \quad (16.121)$$

The right-hand sides of (16.120)–(16.121) can be computed from the given model, yielding two consistency requirements any density modification in the tails must satisfy. If we, say, postulate that the conditional densities ψ_{\min} and ψ_{\max} originate from Black models with unknown constant volatilities σ_{\min} and σ_{\max} , respectively, the consistency requirements will allow us to back out σ_{\min} and σ_{\max} . Indeed, notice that in such a setup

$$\begin{aligned} \psi_{\min}(s | S(T) \leq K_{\min}) \\ = \frac{1}{S_0 \Phi(d(K_{\min})) \sigma_{\min} \sqrt{2\pi T}} \exp\left(-\frac{1}{2} d(s)^2\right), \quad s \leq K_{\min}, \end{aligned}$$

where

$$d(x) = \frac{\ln(x/S(0)) + \frac{1}{2}\sigma_{\min}^2 T}{\sigma_{\min} \sqrt{T}},$$

which can be used to show that

$$\int_{-\infty}^{K_{\min}} s \psi_{\min}(s | S(T) \leq K_{\min}) ds = S(0) \frac{\Phi(d(K_{\min}) - \sigma_{\min} \sqrt{T})}{\Phi(d(K_{\min}))}.$$

Similarly,

$$\int_{K_{\max}}^{\infty} s \psi_{\max}(s | S(T) > K_{\max}) ds = S(0) \frac{\Phi(-d(K_{\max}) + \sigma_{\max} \sqrt{T})}{\Phi(-d(K_{\max}))}.$$

which allows us to uncover σ_{\min} and σ_{\max} from (16.120)–(16.121). We point out that the scheme above can easily be modified to handle more complicated density tails, e.g. the ones from CEV or displaced diffusion models.

Multi-Rate Vanilla Derivatives

After our analysis of single-rate vanilla derivatives in the previous chapter, we now proceed to consider European-type payoffs that are linked to more than one swap or Libor rate. The most important member of this class is the CMS spread, but securities such as floating range accruals and floating range accruals on a CMS spread are also popular in the market.

Valuation of these *multi-rate vanilla derivatives* in a dynamic term structure model presents no conceptual difficulties. However, given the (relatively) high traded volume in derivatives of this type, application of a full term structure model may not be accurate or fast enough (or both) in practice. Hence, in this chapter we look for extensions of the vanilla models of Chapter 16 that allow us to price derivatives linked to multiple rates in an efficient manner. As in Chapter 16, convexity adjustments are an inherent part of valuation, but require only straightforward extensions of the methods developed for the single-rate case. New challenges do arise in the multi-rate setting, however, as we now face the need to specify and control the dependence structure between the rates involved. We spend most of the chapter discussing this issue in detail.

17.1 Introduction to Multi-Rate Vanilla Derivatives

We define a multi-rate vanilla derivative to be a derivative security with a European-type payoff linked to more than one market rate. Given a payment date T_p , a collection of swap or Libor rates $S_1(\cdot), \dots, S_d(\cdot)$, a collection of fixing dates t_1, \dots, t_d , and a d -argument function $f(s_1, \dots, s_d)$, a multi-rate derivative is defined by its payoff

$$V(T_p) = f(S_1(t_1), \dots, S_d(t_d)) \text{ paid at time } T_p. \quad (17.1)$$

Of course, we require $t_i \leq T_p$ for all $i = 1, \dots, d$. Our main focus in this chapter is on the case when all t_i 's are the same: $t_i = T$ for some $T \leq T_p$.

and all $i = 1, \dots, d$. That said, we do develop methods to deal with fixing dates that are not all exactly equal — a case important for floating range accruals, for example — but we still insist here on the fixing dates being “not too far” from each other. If fixing dates are wide apart, valuation is typically best handled in a dynamic term structure model (e.g., the LM model).

As was the case for their single-rate counterparts, multi-rate vanilla derivative payoffs (17.1) are rarely traded themselves; rather, they constitute building blocks for traded securities that are (additive) collections of multi-rate cash flows. Examples include the multi-rate exotic swaps (Section 5.13.3), e.g. swaps paying CMS spread or digital CMS spreads, and various flavors of curve caps. Also popular are the range accruals, see Section 5.13.4; the following variations are (among many others) covered by the techniques in this chapter:

1. Accrual is based on a single market rate, but payment is linked to a different market rate (floating range accruals).
2. The payment rate is either fixed or a function of a collection of market rates, and the observation rate is a difference of two market rates (CMS spread range accruals).
3. Dual range accruals.
4. Curve cap range accruals.

As mentioned above, we focus our attention on developing “vanilla” models for multi-rate derivatives, i.e. models that steer clear of defining dynamics for the whole yield curve and instead merely aim to specify, in the most direct way possible, the distribution of the collection of swap rates $S_1(t_1), \dots, S_d(t_d)$. The efficacy of this approach stems from the fact that the value of payoffs in the class (17.1) have limited, if any, dependence on the actual continuous-time dynamics of the yield curve, since only the joint distribution of terminal rate observations will enter valuation formulas. With this in mind, we may distill the essence of various vanilla-type methods to the following steps.

1. One-dimensional (terminal) distributions of all relevant market rates are extracted from market prices of swaptions.
2. All one-dimensional distributions are brought under the same equivalent martingale measure.
3. A dependence structure is imposed on the vector of market rates, while ensuring that market-implied marginal one-dimensional distributions are preserved.
4. Parameters used in the specification of the dependence structure are estimated historically or, if sufficient market information is available, implied from the prices of certain instruments.
5. A suitable numerical method is applied to integrate the payoff against a specified multi-dimensional distribution.

We start our discussion of multi-rate vanilla derivatives valuation with the first two items of this program.

17.2 Marginal Distributions and Reference Measure

The value of a cash flow with the payoff (17.1) at time $t = 0$ is given by

$$V(0) = \mathbb{E}(\beta(T_p)^{-1} f(S_1(t_1), \dots, S_d(t_d))), \quad (17.2)$$

where \mathbb{E} is the expectation operator for the risk-neutral measure Q . As discussed in Chapter 16, the distribution of each swap rate $S_i(t_i)$ under its annuity measure Q^{A_i} (a measure for which the annuity $A_i(t)$ linked to the rate $S_i(t)$ is a numeraire) can be deduced from market prices of swaptions across strikes. We emphasize that the measures Q^{A_i} are different for different swap rates, and there will generally exist no measure under which all S_i 's are martingales. In principle, one can choose any annuity measure and proceed to derive distributions of all rates under that measure — technical tools for this can be developed relatively easily by extending the results of Section 16.6.9. This approach suffers from a certain arbitrariness, and it is typically both more natural and more convenient to work with the T_p -forward measure. Moreover, we already know how to translate a swap rate distribution under the annuity measure into its distribution under the T_p -forward measure, see Section 16.6.9.

Changing to T_p -forward measure in (17.2), we obtain

$$V(0) = P(0, T_p) \mathbb{E}^{T_p}(f(S_1(t_1), \dots, S_d(t_d))). \quad (17.3)$$

As we recall from Section 16.6.9, the distribution of $S_i(t_i)$ under Q^{T_p} is linked to the (market-implied) distribution of $S_i(t_i)$ under Q^{A_i} via the density relationship

$$Q^{T_p}(S_i(t_i) \in ds) = \frac{A_i(0)}{P(0, T_p)} \alpha_i(s) Q^{A_i}(S_i(t_i) \in ds), \quad (17.4)$$

where the annuity mapping function $\alpha_i(s)$ is defined as

$$\alpha_i(s) = \mathbb{E}^{A_i}\left(\left.\frac{P(t_i, T_p)}{A_i(t_i)}\right| S_i(t_i) = s\right), \quad i = 1, \dots, d. \quad (17.5)$$

In deriving the annuity mapping functions we can follow Chapter 16 and impose a functional relationship between $P(t_i, T_p)/A_i(t_i)$ and $S_i(t_i)$. As we may do so independently for each $i = 1, \dots, d$, the application of the relevant method(s) from Chapter 16 is straightforward.

The formulas (17.4)–(17.5) are exact as written, but we would virtually always apply an approximation for $\alpha_i(s)$ in (17.5). As far as such approximations are concerned, we should note that the techniques discussed in Chapter

16 will be associated with a certain degree of inconsistency in the multi-rate context. For instance, in the common case where $t_i = T$ for some $T \leq T_p$ and all $i = 1, \dots, d$, we see that $P(T, T_p)/A_i(T)$ for a given i is a function of all swap rates $S_1(T), \dots, S_d(T)$. Therefore, the calculation of $\alpha_i(s)$ should in principle involve the dependence structure of all swap rates in the payoff. Although we shall introduce such a dependence structure between all swap rates (see later in Section 17.3) to calculate the value of the derivative cash flow, we typically do *not* use this dependence structure when estimating $\alpha_i(s)$ by (17.5). Instead, we would normally content ourselves with the simpler methods of Chapter 16, such as the linear TSR model. While acknowledging the inconsistency we, among others, realize considerable practical advantages of separating the specification of measure changes via (17.4)–(17.5) from that of the dependence structure; consequently, we adopt this approach throughout this chapter.

With marginal distributions of each $S_i(t_i)$ in (17.3) specified by (17.4), the value $V(0)$ will be strongly sensitive to the dependence structure imposed by the model on the random variables $S_1(t_1), \dots, S_d(t_d)$. Specifying and controlling such dependence is at the heart of the problem of valuing multi-rate vanilla cash flows, with the so-called *copula method* being a popular choice for the job. We introduce this method next.

17.3 Dependence Structure via Copulas

17.3.1 Introduction to Gaussian Copula Method

The copula approach, popularized in financial applications by its widespread usage in credit derivatives modeling, is a method of constructing a joint distribution of random variables consistently with pre-specified one-dimensional marginal distributions. While a thorough treatment of the subject is well beyond the scope of this book — the reader is referred to Nelsen [2006] for that — we proceed to present salient points in the next few sections.

To warm up, let us start with the so-called *Gaussian copula method*, the most common, and easily understood, type of the copula methods in use. We have already seen a particular application of this method in Section 16.7. Let us assume that one-dimensional cumulative distribution functions (CDFs) $\Psi_1(\cdot), \dots, \Psi_d(\cdot)$ have been given, and we are tasked with constructing a multi-dimensional random vector (X_1, \dots, X_d) with a measure of control over the dependence of the random variables X_i , but constrained so that each variable X_i has CDF $\Psi_i(\cdot)$, $i = 1, \dots, d$. The Gaussian copula method accomplishes this by specifying

$$X_i = \Psi_i^{-1}(\Phi(Z_i)), \quad i = 1, \dots, d, \quad (17.6)$$

where $\Phi(\cdot)$ is the standard one-dimensional Gaussian CDF and (Z_1, \dots, Z_d) is a multi-dimensional, normalized¹ Gaussian vector with the correlation matrix R . Clearly the CDF of each X_i thus defined is $\Psi_i(\cdot)$ (see Section 3.1.1.1), and the correlation matrix R provides a way to control the co-dependence structure in the vector.

Let us denote the joint CDF of (X_1, \dots, X_d) as constructed above by $\Psi_{\text{gauss}}(x_1, \dots, x_d)$, and the joint d -dimensional (Gaussian) CDF of (Z_1, \dots, Z_d) by $\Phi_d(z_1, \dots, z_d; R)$. Then it follows from (17.6) that

$$\Psi_{\text{gauss}}(x_1, \dots, x_d) = \Phi_d(\Phi^{-1}(\Psi_1(x_1)), \dots, \Phi^{-1}(\Psi_d(x_d)); R). \quad (17.7)$$

From the joint CDF of (X_1, \dots, X_d) we easily obtain the joint probability density function (PDF) as

$$\psi_{\text{gauss}}(x_1, \dots, x_d) \quad (17.8)$$

$$= \frac{\partial^d}{\partial x_1 \dots \partial x_d} \Psi_{\text{gauss}}(x_1, \dots, x_d) \quad (17.9)$$

$$\begin{aligned} &= \frac{\partial^d}{\partial z_1 \dots \partial z_d} \Phi_d(z_1, \dots, z_d; R) \Big|_{z_i = \Phi^{-1}(\Psi_i(x_i)) \forall i} \times \prod_{i=1}^d \frac{\Psi'_i(x_i)}{\Phi'(\Phi^{-1}(\Psi_i(x_i)))} \\ &= \phi_d(\Phi^{-1}(\Psi_1(x_1)), \dots, \Phi^{-1}(\Psi_d(x_d)); R) \times \prod_{i=1}^d \frac{\psi_i(x_i)}{\phi(\Phi^{-1}(\Psi_i(x_i)))}, \end{aligned}$$

where $\phi(z)$ and $\phi_d(z_1, \dots, z_d; R)$ are the one- and d -dimensional Gaussian PDFs, respectively, and $\psi_i(x_i)$ is the one-dimensional PDF of X_i , $i = 1, \dots, d$.

With the joint PDF $\psi_{\text{gauss}}(x_1, \dots, x_d)$ of (X_1, \dots, X_d) available, the undiscounted² value of a derivative with a payoff $f(X_1, \dots, X_d)$ is given by the multi-dimensional integral

$$V = \int \dots \int f(x_1, \dots, x_d) \psi_{\text{gauss}}(x_1, \dots, x_d) dx_1 \dots dx_d. \quad (17.10)$$

Another, sometimes more useful, expression may be obtained from (17.6), as the derivative security with the payoff $f(X_1, \dots, X_d)$ can also be considered to have the payoff $f(\Psi_1^{-1}(\Phi(Z_1)), \dots, \Psi_d^{-1}(\Phi(Z_d)))$. Therefore its value is given by

$$V = \int \dots \int f(\Psi_1^{-1}(\Phi(z_1)), \dots, \Psi_d^{-1}(\Phi(z_d))) \phi_d(z_1, \dots, z_d; R) dz_1 \dots dz_d. \quad (17.11)$$

¹I.e. $E(Z_i) = 0$, $\text{Var}(Z_i) = 1$, $i = 1, \dots, d$.

²In this chapter we will often consider undiscounted values of derivatives, as we mostly work with cash flows that pay at a single payment time T_p . As it should always be clear whether discounting is applied or not, we do not introduce new notation.

The Gaussian copula method can be implemented easily, is well-understood, and widely used. Still, it suffers from its share of problems, chief among them being its limited control over the shape of the joint distribution, a point that we shall address in more detail later in this chapter.

17.3.2 General Copulas

To develop co-dependence structures more general than those implied by the Gaussian copula method, consider first rewriting (17.7) as

$$\Psi_{\text{gauss}}(x_1, \dots, x_d) = C_{\text{gauss}}(\Psi_1(x_1), \dots, \Psi_d(x_d); R), \quad (17.12)$$

where

$$C_{\text{gauss}}(u_1, \dots, u_d; R) \triangleq \Phi_d(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_d); R).$$

The function $C_{\text{gauss}}(u_1, \dots, u_d; R)$ is easily seen to define a specific multi-dimensional distribution function (as defined in, for example, Nelsen [2006]) for a vector of d random variables with marginal distributions that are all *uniform* on $[0, 1]$. It turns out that this concept can be generalized nicely.

Definition 17.3.1. Consider a function $C : [0, 1]^d \rightarrow [0, 1]$. $C(u_1, \dots, u_d)$ is said to be a d -dimensional copula function if it defines a valid joint distribution function for a d -dimensional vector of random variables, with each variable being uniformly distributed on $[0, 1]$.

The requirement that a copula defines a joint distribution function for a vector of uniformly distributed random variables puts a number of strong constraints on the form of the function C . For instance, it is clear that C must be increasing in its arguments. Also, if the i -th argument of C is 0, C itself must return 0, irrespective of the values of the remaining $d - 1$ arguments. Further, if all but the i -th argument of C are set to 1, C must return the i -th argument itself. The last two simple relations (which we invite the reader to verify) can be summarized as follows:

$$C(u_1, u_2, \dots, u_{i-1}, 0, u_{i+1}, \dots, u_d) = 0, \quad (17.13)$$

$$C(1, 1, \dots, 1, u_i, 1, 1, \dots, 1) = u_i. \quad (17.14)$$

To give a few introductory examples, notice that a particularly simple copula arises if all d uniform variables underlying the copula are independent. The resulting *independence copula* C_{ID} is

$$C_{\text{ID}} = \prod_{i=1}^d u_i. \quad (17.15)$$

To state the copula C_D that defines *perfect dependence*, introduce d uniform random variables U_1, U_2, \dots, U_d and set $U_1 = U_2 = \dots = U_d$. In this case we get (with P being a probability measure)

$$\begin{aligned}
C_D(u_1, u_2, \dots, u_d) &= P(U_1 \leq u_1, U_2 \leq u_2, \dots, U_d \leq u_d) \\
&= P(U_1 \leq u_1, U_1 \leq u_2, \dots, U_1 \leq u_d) \\
&= P\left(U_1 \leq \min_{i=1,\dots,d} u_i\right) \\
&= \min_{i=1,\dots,d} u_i.
\end{aligned} \tag{17.16}$$

For obvious reasons, the *perfect anti-dependence* copula C_{AD} can only be stated for the case $d = 2$. In this case, we write $U_2 = 1 - U_1$ for two uniform random variables U_1, U_2 , and get

$$\begin{aligned}
C_{AD}(u_1, u_2) &= P(U_1 \leq u_1, U_2 \leq u_2) \\
&= P(U_1 \leq u_1, 1 - U_1 \leq u_2) \\
&= P(1 - u_2 \leq U_1 \leq u_1) \\
&= (u_1 + u_2 - 1)^+.
\end{aligned} \tag{17.17}$$

The anti-dependence copula in (17.17) cannot be extended to $d > 2$, but we can still define a function

$$G_{AD}(u_1, u_2, \dots, u_d) = \left(\sum_{i=1}^d u_i + 1 - d \right)^+, \tag{17.18}$$

such that $G_{AD} = C_{AD}$ for the case $d = 2$. Although G_{AD} defined this way is not itself a copula for $d \geq 3$, it turns out that this function can be used to bound any valid copula function. Specifically, one can prove the following result.

Theorem 17.3.2. *Any valid d -dimensional copula function C must satisfy the Frechet bounds*

$$G_{AD}(u_1, u_2, \dots, u_d) \leq C(u_1, u_2, \dots, u_d) \leq C_D(u_1, u_2, \dots, u_d).$$

The real strength of the copula function concept originates with the fact that it allows us to separate co-dependence information from marginal distributions. Specifically, given a copula function $C(u_1, \dots, u_d)$ and a collection of marginal CDFs $\Psi_1(x_1), \dots, \Psi_d(x_d)$, we can construct a d -dimensional joint distribution function $\Psi_C(x_1, \dots, x_d)$ with marginals $\Psi_1(x_1), \dots, \Psi_d(x_d)$ by a formula similar to (17.12),

$$\Psi_C(x_1, \dots, x_d) = C(\Psi_1(x_1), \dots, \Psi_d(x_d)). \tag{17.19}$$

The simple proof of the fact that $\Psi_C(x_1, \dots, x_d)$ defined by (17.19) is a true d -dimensional distribution function is left to the reader. By the so-called *Sklar's theorem*, the opposite is also true: for any d -dimensional distribution

function there exists a copula function such that the joint distribution function can be represented in the form (17.19).³

The joint PDF ψ_C associated with the CDF Ψ_C in (17.19) is given by (compare to (17.9))

$$\begin{aligned}\psi_C(x_1, \dots, x_d) &= \frac{\partial^d}{\partial x_1 \dots \partial x_d} \Psi_C(x_1, \dots, x_d) \\ &= c(\Psi_1(x_1), \dots, \Psi_d(x_d)) \prod_{i=1}^d \psi_i(x_i),\end{aligned}\quad (17.20)$$

where

$$c(u_1, \dots, u_d) = \frac{\partial^d}{\partial u_1 \dots \partial u_d} C(u_1, \dots, u_d), \quad (17.21)$$

is the *copula density* and $\psi_i(\cdot)$'s are the marginal PDFs. For the Gaussian copula the copula density is given by

$$c_{\text{gauss}}(u_1, \dots, u_d; R) = \frac{\phi(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_d); R)}{\prod_{i=1}^d \phi(\Phi^{-1}(u_i))}. \quad (17.22)$$

17.3.3 Archimedean Copulas

With technical background material now out of the way, let us proceed to examine some concrete examples of copula functions, beyond the Gaussian class. One choice that is quite popular in the copula literature is the so-called *Archimedean* class of copulas. This class requires specification of a *generator function* $\omega: [0, 1] \rightarrow \mathbb{R}$ satisfying

$$\lim_{x \rightarrow 0} \omega(x) = +\infty, \quad \omega(1) = 0, \quad \omega'(x) < 0, \quad \omega''(x) > 0.$$

From a generator function, a corresponding Archimedean copula $C_{\text{arch}}(u_1, \dots, u_d; \omega)$ can be defined by the relation

$$C_{\text{arch}}(u_1, \dots, u_d; \omega) = \omega^{-1} \left(\sum_{i=1}^d \omega(u_i) \right).$$

It is a trivial exercise to show that $C_{\text{arch}}(u_1, \dots, u_d; \omega)$ is indeed a copula function.

The generator function is often indexed with a parameter, specifying a parametric family of Archimedean copulas. Of particular note are the following two families:

³If the marginal distribution functions $\Psi_1, \Psi_2, \dots, \Psi_d$ are all continuous, then the copula function is *unique*.

- *Clayton copula:*

$$\omega_{\text{clayton}}(u; \theta) = u^{-\theta} - 1, \quad \theta > 0.$$

- *Gumbel copula:*

$$\omega_{\text{gumbel}}(u; \theta) = (-\ln u)^{\theta}, \quad \theta > 0.$$

The corresponding copulas are given by

$$C_{\text{clayton}}(u_1, \dots, u_d; \theta) = \left(\sum_{i=1}^d u_i^{-\theta} - d + 1 \right)^{-1/\theta},$$

$$C_{\text{gumbel}}(u_1, \dots, u_d; \theta) = \exp \left(- \left(\sum_{i=1}^d (-\ln u_i)^{\theta} \right)^{1/\theta} \right).$$

In the special case of $\theta = 1$, the Gumbel copula becomes

$$C_{\text{gumbel}}(u_1, \dots, u_d; 1) = \prod_{i=1}^d u_i,$$

which is the independence copula C_{ID} introduced earlier.

A quick graphing exercise shows that as θ is raised, both the Clayton and Gumbel copulas assign increasing probability mass around the point $(0, \dots, 0)$; in terms of the joint distribution of the market rates, this corresponds to an increase in the probability of a joint down-move of the rates.

17.3.4 Making Copulas from Other Copulas

With the Archimedean copulas introduced in the previous section, we can use the parameter θ of both the Clayton or Gumbel copula to control the probability of a joint down-move of interest rates, but we have no direct control over other moves of interest such as a joint *up*-move of the rates. This is easily fixed, however, by an application of the following general result.

Lemma 17.3.3. *If $C(u_1, \dots, u_d)$ is a copula, then the function obtained by reflecting C in dimension i ,*

$$\bar{C}(u_1, \dots, u_i, \dots, u_d; \{i\}) \triangleq C(u_1, \dots, 1, \dots, u_d) - C(u_1, \dots, 1 - u_i, \dots, u_d), \quad (17.23)$$

is also a copula for any $i = 1, \dots, d$. The density of the reflected copula is given by

$$\bar{c}(u_1, \dots, u_i, \dots, u_d; \{i\}) = c(u_1, \dots, 1 - u_i, \dots, u_d).$$

Proof. Trivial consequence of the fact that if U follows a uniform $[0, 1]$ distribution, then so does $1 - U$. \square

By repeated application of the lemma, it is easy to see that it generalizes to multiple indices. Specifically, if we denote by $\bar{C}(\dots; \{i_1, \dots, i_M\})$ a function obtained by repeating the mapping (17.23) for all i_m , $m = 1, \dots, M$, this is still a copula. Focusing on the two-dimensional case $d = 2$ and choosing the Clayton copula for concreteness, we have

$$\begin{aligned}\bar{C}_{\text{clayton}}(u_1, u_2; \theta; \{1, 2\}) &= C_{\text{clayton}}(1, 1; \theta) - C_{\text{clayton}}(1, 1 - u_2; \theta) \\ &\quad - C_{\text{clayton}}(1 - u_1, 1; \theta) + C_{\text{clayton}}(1 - u_1, 1 - u_2; \theta),\end{aligned}$$

and now the parameter θ controls the probability of a joint *up*-move in the two market rates. In the copula $\bar{C}_{\text{clayton}}(u_1, u_2; \theta; \{1\})$ the parameter θ controls the joint probability of an up-move of the first rate and a down-move in the second rate, and in the copula $\bar{C}_{\text{clayton}}(u_1, u_2; \theta; \{2\})$, the parameter θ controls the joint probability of a down-move of the first rate and an up-move in the second rate.

Another way of creating copulas uses the observation that a convex combination of copulas is also a copula.

Lemma 17.3.4. *Let us denote $u = (u_1, \dots, u_d)^\top$, and let there be given M different d -dimensional copulas $C_1(u), \dots, C_M(u)$, as well as a collection of non-negative weights w_1, \dots, w_M such that $\sum_{m=1}^M w_m = 1$. The linear combination, or mixture,*

$$C_{\text{mix}}(u) = \sum_{m=1}^M w_m C_m(u),$$

is also a copula.

Proof. Trivial. \square

One can interpret mixture copulas as representations of the idea that different dependence structures of the random variables are realized in different states of the world, with these states having probabilities w_m , $m = 1, \dots, M$. To give a simple example, consider a Gaussian copula setting where there are two states of the correlation matrix: R_{hi} (“excited state”) and R_{normal} (“normal state”). Assuming that the probabilities of these states are w_{normal} and $w_{\text{hi}} = 1 - w_{\text{normal}}$, respectively, we would define

$$\begin{aligned}C_{\text{multigauss}}(u_1, \dots, u_d; w_{\text{normal}}, w_{\text{hi}}, R_{\text{normal}}, R_{\text{hi}}) \\ = w_{\text{normal}} C_{\text{gauss}}(u_1, \dots, u_d; R_{\text{normal}}) + w_{\text{hi}} C_{\text{gauss}}(u_1, \dots, u_d; R_{\text{hi}}).\end{aligned}\tag{17.24}$$

If, say, R_{hi} is supposed to reflect an unlikely crash state, we could set $R_{\text{hi}} \gg R_{\text{normal}}$ (in, e.g., an element-wise sense) and $w_{\text{hi}} \ll w_{\text{normal}}$.

To motivate the next method for constructing copulas, we note that for any function $q(u)$ we have a trivial identity

$$u = q(u) \times \frac{u}{q(u)}.$$

Hence, if we have two (2-dimensional) copulas $C_1(u, v)$ and $C_2(u, v)$ and functions $q(u), p(v)$ such that $q(1) = p(1) = 1$, then the function

$$C(u, v) = C_1(q(u), p(v)) C_2\left(\frac{u}{q(u)}, \frac{v}{p(v)}\right)$$

would satisfy requirement (17.14), as

$$\begin{aligned} C(u, 1) &= C_1(q(u), 1) C_2\left(\frac{u}{q(u)}, 1\right) = q(u) \frac{u}{q(u)} = u, \\ C(1, v) &= C_1(1, p(v)) C_2\left(1, \frac{v}{p(v)}\right) = p(v) \frac{v}{p(v)} = v. \end{aligned}$$

While we have not demonstrated that C is a copula, it turns out that this is the case if certain conditions are imposed on q (and p). The relevant result (suitably extended to dimension d) for these so-called *product copulas* can be found in Liebscher [2008]:

Theorem 17.3.5. *Let C_1, \dots, C_M be d -dimensional copulas, and let $q_{m,i} : [0, 1] \rightarrow [0, 1]$, $m = 1, \dots, M$, $i = 1, \dots, d$, be functions that are either strictly increasing or identically equal to 1. Suppose that $\prod_{m=1}^M q_{m,i}(u) = u$ for $u \in [0, 1]$, $i = 1, \dots, d$, and $\lim_{u \rightarrow 0+} q_{m,i}(u) = q_{m,i}(0)$. Then*

$$C(u_1, \dots, u_d) = \prod_{m=1}^M C_m(q_{m,1}(u_1), \dots, q_{m,d}(u_d))$$

is a copula.

In this theorem, consider taking $M = 2$, $C_1 = \prod_{i=1}^d u_i$ (the independence copula), $C_2 = C_{\text{gauss}}(u_1, \dots, u_d; R)$ (Gaussian copula with correlation matrix R), $q_{1,i}(u) = u^{1-\theta_i}$, and $q_{2,i}(u) = u^{\theta_i}$ for $\theta_i \in [0, 1]$, $i = 1, \dots, d$. We then obtain a copula that we find particularly useful for multi-rate derivatives:

Corollary 17.3.6. *Let R be a $d \times d$ correlation matrix and $\theta = (\theta_1, \dots, \theta_d)^T \in [0, 1]^d$ a d -dimensional vector of parameters. Then the power Gaussian function*

$$C_{\text{PG}}(u_1, \dots, u_d; R, \theta) \triangleq \left(\prod_{i=1}^d u^{1-\theta_i} \right) C_{\text{gauss}}(u_1^{\theta_1}, \dots, u_d^{\theta_d}; R) \quad (17.25)$$

is a copula.

Section 17.4.3 demonstrates that the power Gaussian copula provides a parsimonious, yet flexible way of specifying dependence structure for CMS spread options.

17.4 Copula Methods for CMS Spread Options

We recall from Section 5.13.3 that the payoff of a spread option is given by

$$(S_1(T) - S_2(T) - K)^+$$

at time $T_p \geq T$, where $S_1(T)$, $S_2(T)$ are two swap rates of different tenors fixing at time T . The (undiscounted) value of a spread option is therefore given by

$$V(0; T, K) = E^{T_p} \left((S_1(T) - S_2(T) - K)^+ \right). \quad (17.26)$$

CMS spread options (Section 5.13.3) are, arguably, the most liquid of the multi-rate vanilla derivatives. Euro- and USD-denominated spread options are traded among brokers, assuring reasonable visibility into market prices.

17.4.1 Normal Model for the Spread

Market quotes of CMS spread options often come in the form of implied Normal (also known as Bachelier, Gaussian or, simply, basis-point) volatilities of the spread. As we have already seen in Section 14.4.3, the implied Normal spread volatility $\sigma_N(T, K)$ (for a given pair of swap rates) is defined by equating the undiscounted market price of a spread option to its Normal model price at the given volatility, i.e.

$$V_{\text{mkt}}(0; T, K) = c_N(0, E^{T_p}(S_1(T) - S_2(T)); T, K; \sigma_N(T, K)),$$

where $c_N(t, S; T, K; \lambda)$ is the Normal pricing formula with volatility λ defined by (7.16).

A few features of this formula are worth pointing out. First, the Normal model, rather than the Black model, is used as a convention for quoting spread options because the spread process $S_1(T) - S_2(T)$ can become negative, something which is disallowed in a log-normal setting. Indeed, spread options that correspond to zero strike, $K = 0$, are among the most liquid — and for the Black model, zero-strike implied volatility is simply not defined. Second, we notice that in order to back out the implied volatility, we must use the *convexity-adjusted* forward of the spread $E^{T_p}(S_1(T) - S_2(T))$. When trading with each other, dealers therefore need to agree on the convexity adjustments of the underlying swap rates before they can agree on the implied volatility. Finally, we note that in practice the implied volatility of a particular CMS spread will always depend on both the strike and the expiry, reflecting the fact that market distributions of spreads are not perfectly Gaussian. Much of the work in modeling spread options is focused on properly capturing the market-implied distribution of the spread.

The Normal model provides a convenient common language for quoting spread options, but it is not well-suited for risk management. Apart from the standard risk management issues that arise when using models with

strike-indexed parameters (see the related discussion in Section 16.1.1), we highlight the fact that strike-indexed spread volatilities generate no explicit link to the marginal distribution parameters — including volatilities — of the underlying swap rates. As a consequence, spread options will, unreasonably, show no vega (volatility sensitivity, see Section 8.9) to swaption parameters, making them difficult, or even impossible, to hedge. A related issue is the absence of direct information useful for the pricing of *other* payoffs that depend on the same two swap rates. For example, it is unclear how to translate a strike-dependent spread volatility into a parameter that could be used to value a spread option with non-standard gearing, i.e. a payoff

$$(S_1(T) - gS_2(T) - K)^+, \quad (17.27)$$

where $g \neq 1$.

The two issues above can be addressed with a certain degree of success by copula methods. We shall describe this approach momentarily, but let us note that if our ultimate goal is to build dynamic term structure models that are consistent with the prices of vanilla spread options (and European swaptions, of course), we will ultimately need to abandon the copula approach — see Section 17.8 and beyond.

17.4.2 Gaussian Copula for Spread Options

Arguably, the simplest copula-based spread option valuation method is obtained by specifying a two-dimensional Gaussian copula, a parameterization that depends on a single correlation parameter ρ . A spread option is then valued using the following procedure.

1. For each of the swap rates $S_i(T)$, $i = 1, 2$, a market-implied density under its own annuity measure, $\psi_i^{A_i}(x)$, is derived from swaption prices. We may perform this derivation non-parametrically (see Section 16.6.9) or rely on a vanilla model calibrated to swaption prices.
2. Using Proposition 16.6.4, each $\psi_i^{A_i}(x)$, $i = 1, 2$, is converted into the PDF $\psi_i^{T_p}(x)$ of the i -th swap rate under the T_p -forward measure .
3. For a given correlation ρ , the joint probability density function $\psi^{T_p}(x_1, x_2; \rho)$ of $(S_1(T), S_2(T))$ is defined by (17.20), where the copula density $c_{\text{gauss}}(\cdot)$ for the Gaussian copula is given by (17.22). (In this case, notice that R is a 2×2 matrix with 1 on the diagonal and ρ off-diagonal.)
4. The payoff $(x_1 - x_2 - K)^+$ is integrated against the density $\psi^{T_p}(x_1, x_2; \rho)$ to obtain the undiscounted model price of the spread option:

$$V_{\text{mdl}}(0; T, K, \rho) = \int \int (x_1 - x_2 - K)^+ \psi^{T_p}(x_1, x_2; \rho) dx_1 dx_2.$$

See Sections 17.6.1, 17.6.2 for details on numerical implementation.

To calibrate the model to market, we would also perform the following step:

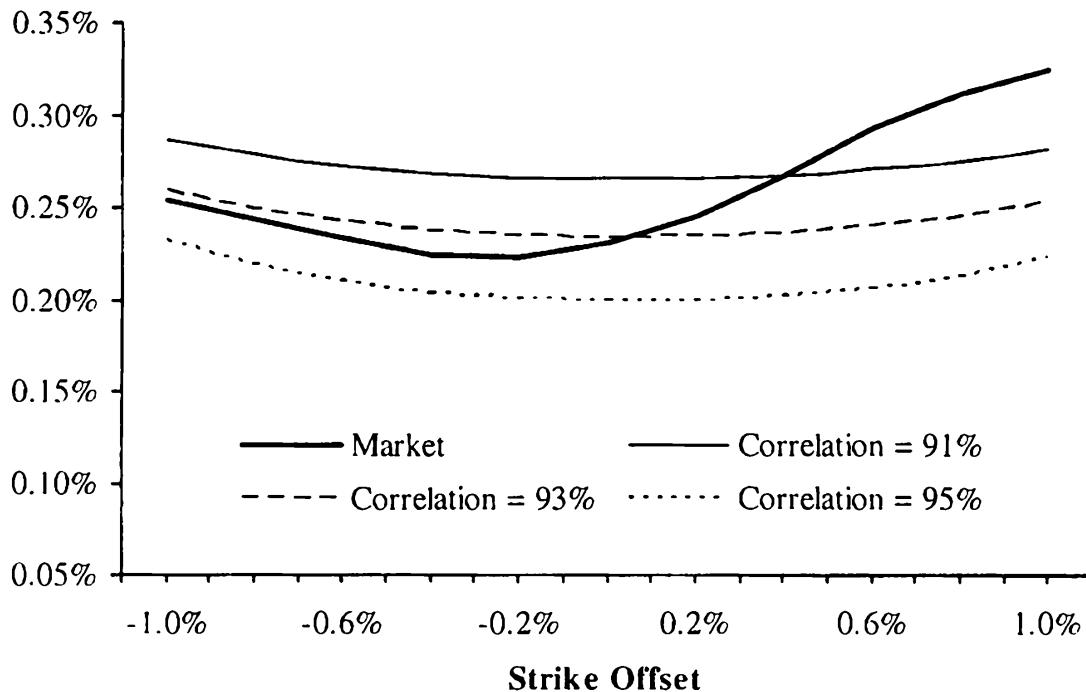
- 5 Find $\rho = \rho(T, K)$, the *implied spread correlation*, such that the model price of the spread option matches its market price,

$$V_{\text{mdl}}(0; T, K, \rho(T, K)) = V_{\text{mkt}}(0; T, K).$$

Clearly, with ρ fixed, changes to marginal distributions of the two swap rates would affect their joint distribution via (17.20), ultimately impacting the spread option value. Also, for a given ρ we can integrate any payoff, including (17.27), against the joint density, allowing us to value all payoffs linked to the two swap rates consistently. Hence, as advertised the copula method can overcome several of the issues we identified earlier for the simple Normal spread model. On the other hand, were we to calibrate (as in Step 5) the Gaussian copula model to market values of spread options with a fixed maturity date but different strikes, we would typically obtain a *different* value of the implied spread correlation for each strike, an effect sometimes known as the *correlation smile*. In other words, a simple Gaussian copula has insufficient flexibility to capture market-observed distributions of CMS spreads. To elaborate a bit further on this, in Figure 17.1 we show typical shapes of Normal volatilities of spreads across strikes as implied by different values of the Gaussian copula correlation ρ . We see that ρ basically can only shift the implied spread option volatility smile in parallel, making it impossible to match the market-implied volatility smile listed in the figure. To properly match the market, we apparently need to go beyond Gaussian copulas and look for more flexible alternatives.

Before starting our search for an alternative copula, let us briefly clarify one point about Gaussian copulas. While the implied Gaussian copula correlation is often used to characterize the dependence of the swap rates, it is important to realize that this parameter is not well-defined unless marginal distributions of the swap rates (under some common measure!) are clearly specified. Above we defined the implied correlation by i) implying the marginal distributions under the annuity measures from swaptions; and ii) transplanting these distribution into the T_p -forward measure using a given annuity-mapping function (e.g. one defined by a the linear TSR model). Were we to change any of these “ingredients”, we would need to use a different correlation in the copula to obtain the same spread option value. For example, the following marginals will also lead to reasonable definitions of the implied spread correlation (all under the T^p -forward measure):

1. The marginal distribution for each swap rate $S_i(T)$, $i = 1, 2$, is log-normal with mean $E^{T_p}(S_i(T))$, and volatility set equal to the implied Black volatility of the at-the-money (strike = $S_i(0)$) swaption.
2. The marginal distribution for each swap rate $S_i(T)$, $i = 1, 2$, is log-normal with mean $E^{T_p}(S_i(T))$, and volatility set equal to the implied Black volatility of the at-the-money (strike = $E^{T_p}(S_i(T))$) CMS cap.

Fig. 17.1. Implied Normal Spread Volatility

Notes: Implied Normal spread volatility for a 5 year spread option on the difference between 10 year and 2 year swap rates. The x -axis shows the strike as an offset to the CMS-adjusted forward value of the spread. The “Market” data was observed in November 2007. The spread volatility curves implied by a Gaussian copula are shown for three levels of copula correlation (91%, 93%, and 95%), assuming market-implied marginal distributions per Section 17.4.2.

3. The marginal distribution for each swap rate $S_i(T)$, $i = 1, 2$, is Gaussian with mean $E^{T_p}(S_i(T))$, and volatility equal to the implied Normal volatility of the at-the-money (strike = $S_i(0)$) swaption.
4. The marginal distribution for each swap rate $S_i(T)$, $i = 1, 2$, is Gaussian with mean $E^{T_p}(S_i(T))$, and volatility equal to the implied Normal volatility of the at-the-money (strike = $E^{T_p}(S_i(T))$) CMS cap.

This, far from exhaustive, list of alternatives is intended to give the reader a flavor of issues that can arise when two parties are communicating price information in terms of implied spread correlations — clearly they would have to be very precise about all relevant details to avoid misunderstandings⁴. Communication issues aside, the choice of conventions for specifying the meaning of implied Gaussian copula correlations largely comes down to personal preferences and consistency with the rest of one’s modeling setup. We leave it to the reader to ponder pros and cons of various approaches; not

⁴To demonstrate the dependence of implied correlation on the marginal distribution of swap rates, Appendix 17.A uses a setting with displaced diffusion processes to quantify the effect on ATM spread option prices due to changes in the swap rate volatility skew.

surprisingly, the way we defined implied spread correlations to begin with is our personal choice. It is worth pointing out, however, that in the last case listed above (and *only* in the last case), the spread correlations can be extracted from the implied spread volatilities directly by simple algebraic manipulations. If we denote the Gaussian copula correlation in case 4 by $\rho_N(T, K)$, the relevant at-the-money CMS cap volatilities by $\sigma_{N,i}$, $i = 1, 2$, and the normal spread volatility by $\sigma_N(T, K)$, then it is easily seen that

$$\sigma_N(T, K)^2 = \sigma_{N,1}^2 + \sigma_{N,2}^2 - 2\sigma_{N,1}\sigma_{N,2}\rho_N(T, K).$$

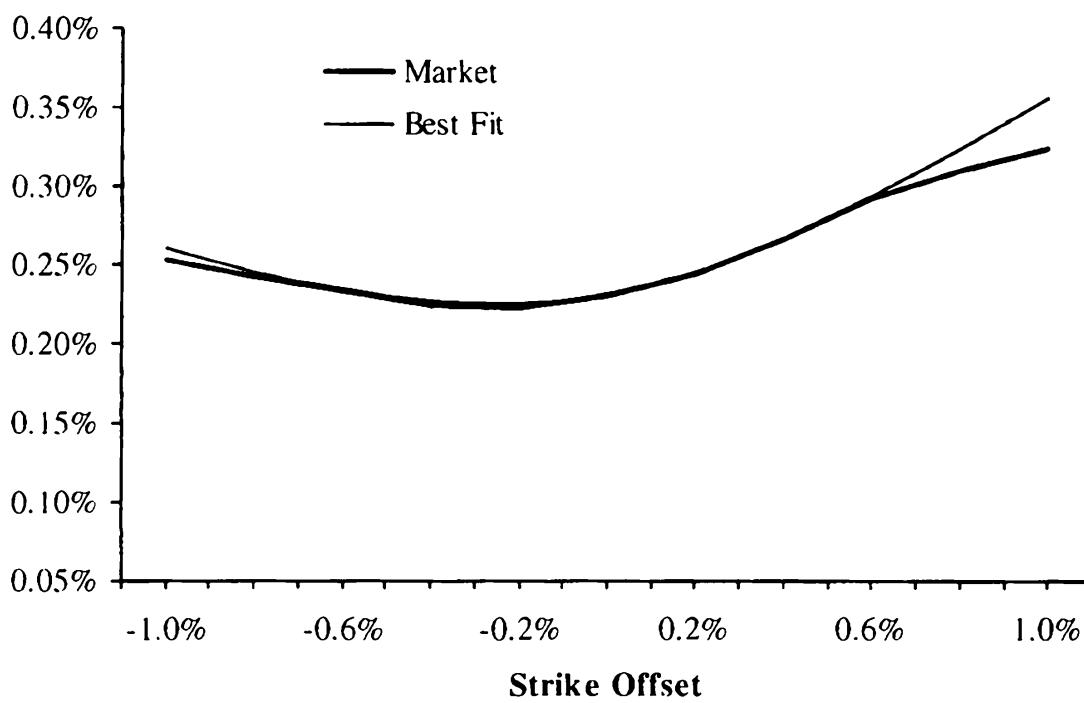
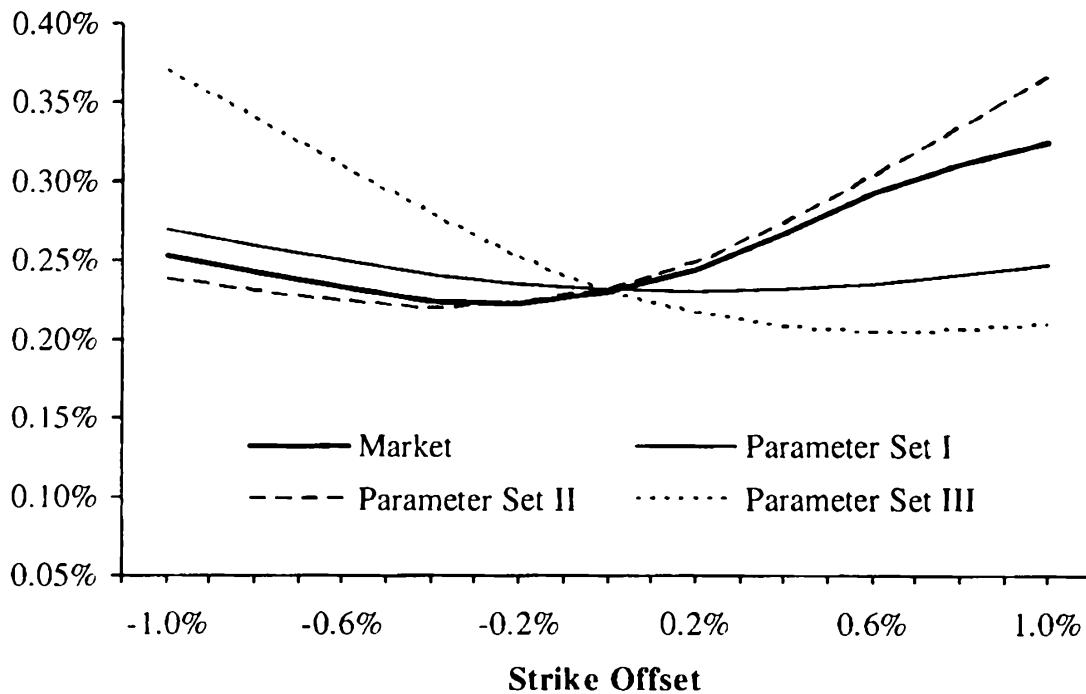
17.4.3 Spread Volatility Smile Modeling with the Power Gaussian Copula

After our small detour into various definitions of implied spread correlations, we now proceed with the task of identifying a copula that would allow us to match the market-implied spread volatility smile as closely as possible. As we saw in Figure 17.1, the Gaussian copula provides us with the ability to change the overall level of the implied spread volatility smile, but lacks controls over its slope or curvature. Mixtures of Gaussian copulas as in (17.24) allow for more flexibility, but ultimately only provide the mechanism to control the curvature of the implied spread volatility smile and not its slope, a fact that is easy to verify experimentally. While we can consider adding standard Archimedean copulas from Section 17.3.3, either stand-alone, “reflected” as in Lemma 17.3.3, or as parts of mixtures, we generally find that this does not provide sufficient flexibility either. On the other hand, the power Gaussian copula (17.25), despite its parsimonious specification, turns out to provide direct control over the relevant features of the spread volatility smile.

In the two-dimensional case, the Gaussian power copula is given by

$$C_{PG}(u, v; \rho, \theta_1, \theta_2) = u^{1-\theta_1} v^{1-\theta_2} C_{gauss}(u^{\theta_1}, v^{\theta_2}; \rho), \quad (17.28)$$

where ρ is a correlation coefficient and $\theta_1, \theta_2 \in [0, 1]$. We would expect the correlation ρ to move the implied spread volatility smile up and down, just as in the (pure) Gaussian case of Figure 17.1. It turns out that the remaining two parameters θ_1, θ_2 provide good control over the slope and curvature of the smile. Starting from $\theta_1 = \theta_2 = 1$, the base Gaussian case, we find that as θ_1 decreases from 1 towards 0, the spread smile rotates counter-clock-wise, and as the parameter θ_2 decreases, the smile rotates in the opposite direction. The effects are clearly visible in the first graph in Figure 17.2. By decreasing both parameters at the same time curvature is added to the smile, and a good fit to market volatilities can be achieved, as the second graph in Figure 17.2 demonstrates.

Fig. 17.2. Implied Normal Spread Volatility

Notes: Implied Normal spread volatility for the spread option in Figure 17.1 in the power Gaussian copula (17.28). Three parameter scenarios as well as best fit parameters from Table 17.1 are shown.

17.4.4 Copula Implied From Spread Options

As we have shown, the power Gaussian copula is capable of reproducing a wide range of market-observed shapes of spread volatility smiles. Yet, clearly, it cannot reproduce at least some of the spread volatility smiles exactly (as,

	Set I	Set II	Set III	Best Fit
ρ	92.6%	95.6%	95.8%	95.5%
θ_1	100.0%	90.0%	100.0%	91.0%
θ_2	100.0%	100.0%	90.0%	99.0%

Table 17.1. Power Gaussian copula (17.28) parameter sets for Figure 17.2.

for example, seen in the wings of the graph in Figure 17.2). Nor, in theory, could the same feat be accomplished by any other copula function with a finite number of parameters. It is, then, natural to wonder whether it is possible to come up with a copula that would reproduce market volatilities (or, equivalently, values) of spread options *exactly* for all values of spread strikes. Of course, spread options are never traded in the whole continuum of strikes, yet the question remains valid in the idealized case of knowing the full distribution of the spread. We call such a copula a *spread-implied copula*.

As we write this book there are no definitive results on spread-implied copulas, so our treatment will necessarily be brief. From general dimensionality analysis it is clear that, in general, there should be a continuum of copulas that would match a set of spread option values of all strikes (and, of course, given marginal distributions of individual swap rates). It is also clear that some collections of spread option values are fundamentally incompatible with given marginal distributions. As a somewhat trivial example, consider marginal distributions that correspond to non-random swap rates (i.e. marginal PDFs are delta functions). Then, clearly, most exogenously specified values of spread options will be incompatible with such marginals irrespective of what dependence structure we would specify.

While it is relatively easy to come up with some examples where spread option values are inconsistent with marginal distributions, the precise results on the existence of spread-implied copulas are unknown to us. Pragmatically, however, there are a number of constructive algorithms for creating “candidate” spread-implied copulas; upon construction these functions can be verified to be true copulas (or not, as the case might be).

One such algorithm, recently proposed by Austing [2010] in the context of a related problem in FX cross-smile modeling, uses the fact that the joint distribution of the swap rates $S_1(T)$, $S_2(T)$ can be obtained from the prices of the so-called *best-of-calls options*, i.e. options with the payoff

$$\max \left((S_1(T) - K_1)^+, (S_2(T) - K_2)^+ \right)$$

for all strikes K_1 and K_2 . Importantly, prices of such best-off call options can be parameterized consistently with given values of options on each individual swap rate (for all strikes) and all spread options. We refer the reader to the original paper for details, and instead take a somewhat different tack.

Of central importance to our discussion is the result that links a copula function to values of spread options that we develop later in the book, see Corollary 17.6.2. For convenience we pre-announce it here; the formula (17.38) states that the spread option values $V(0; T, K)$ are given by

$$V(0; T, K) = \int_{-\infty}^{\infty} (1_{\{x>0\}} - C(\Psi_1(x), \Psi_2(x - K))) dx - E^{T_r}(S_2(T)) - K \quad (17.29)$$

for all values of strike K . Here Ψ_1 and Ψ_2 are the marginal CDFs of the swap rates $S_1(T)$ and $S_2(T)$ and C is the copula function. If spread option values $V(0; T, K)$ are given for all K , then we can treat (17.29) as an (integral) equation for the unknown copula function C .

One can imagine a number of approaches for attacking this equation. We can, for example, discretize $V(0; T, \cdot)$ and $C(\cdot, \cdot)$ on a grid to obtain a linear system for the grid values of the copula function. The constraints on C to be a copula are then given by certain linear inequalities, and the resulting problem can be solved by linear algebra methods. One should be mindful that (17.29) does not introduce enough constraints to give a unique solution, so they should be supplemented with other, exogenous, conditions. We leave it to the reader to explore these ideas.

As (17.29) is underspecified, another line of attack would involve parameterizing C by a one-dimensional family and then solving for the parameter function. For example, one can take an Archimedean copula from Section 17.3.3 with an unspecified parameter function $\omega(\cdot)$ that one then would solve for (numerically) from (17.29). In the same vain, one can take a product copula from Theorem 17.3.5 with one of the $q_{m,i}(\cdot)$ unspecified and solve for it from (17.29).

A rather simple alternative to these, admittedly quite computationally expensive, methods is based on the idea of mixing two copulas together with the weight that is a function of (essentially) spread option strike. The resulting function is not guaranteed to be a copula (something that should be checked post-construction) but this deficiency is somewhat compensated by a rather simple numerical algorithm, an algorithm that we now proceed to outline.

Let us choose two copulas, $C_{lo}(u, v)$ and $C_{hi}(u, v)$, and let us define the “copula” C by

$$\begin{aligned} C(u, v; \alpha(\cdot)) &= \alpha (\Psi_1^{-1}(u) - \Psi_2^{-1}(v)) C_{lo}(u, v) \\ &\quad + (1 - \alpha (\Psi_1^{-1}(u) - \Psi_2^{-1}(v))) C_{hi}(u, v), \end{aligned} \quad (17.30)$$

where $\alpha(\cdot)$ is an unknown mixing weight. Substituting this expression into (17.29) we obtain a simple equation on the weight function $\alpha(K)$ that we can solve to yield

$$\alpha(K) = \frac{V(0; T, K) - V_{hi}(0; T, K)}{V_{lo}(0; K, K) - V_{hi}(0; T, K)},$$

where V_{lo} and V_{hi} are the spread option values that correspond to the copulas C_{lo} and C_{hi} . The copulas C_{lo} and C_{hi} should be chosen so that the spread option values V are spanned by V_{lo} and V_{hi} ; for example we can take the anti-dependence copula C_{AD} (see (17.17)) as C_{lo} and the perfect dependence copula C_D (see (17.16)) as C_{hi} , or we can take Gaussian copulas with sufficiently low/high correlations. The right choice of C_{lo} and C_{hi} will guarantee that $C(u, v) \in [0, 1]$; also by construction the marginals conditions (17.13)–(17.14) are always satisfied. Where C may fail to be a real copula is in being a true two-dimensional distribution function; whether this is satisfied or not will have to be checked *post-factum*.

17.5 Rates Observed at Different Times

In Section 17.4 we assumed that all swap rates fix at the same time T . This assumption is, however, not required for the copula method to work, although certain complications do arise for securities with multiple fixing dates. To illustrate, consider a floating range accrual which, according to the definition in Section 5.13.4, pays the amount

$$S_1(T) \times \frac{\#\{t \in [T, T + \tau] : S_2(t) \in [l, u]\}}{\#\{t \in [T, T + \tau]\}}$$

at time $T_p = T + \tau$, where $\#\{\cdot\}$ is the number of business days that satisfy the specified trigger condition. Clearly, a range accrual can be decomposed into a series of floating digital options, i.e. contracts with the payoff

$$S_1(T) \times 1_{\{S_2(t) \in [l, u]\}},$$

for $t \in [T, T + \tau]$, paid at T_p . The (undiscounted) value of such digital is

$$V_{\text{digi}}(0) = E^{T_p} (S_1(T) \times 1_{\{S_2(t) \in [l, u]\}}).$$

The distributions of $S_1(T)$, $S_2(t)$ in their annuity measures can be mapped into the distributions under the T_p -forward measure using standard techniques, at which point one could, in principle, apply the copula method. The main complication here is not so much the mechanics of the copula method, but rather the meaning of the parameters for the dependence structure. If we take the Gaussian copula as an example, the correlation parameter of the copula would specify the dependence between $S_1(T)$ and $S_2(t)$ — two swap rates observed at *different* times. Clearly we cannot use the same correlation parameter as we would use to characterize the dependence between $S_1(T)$ and $S_2(T)$ (i.e. when the rates are observed on the same date)⁵. Instead,

⁵This is most easily seen by assuming that the two rates are in fact the same. In this case the correlation between $S_1(T)$ and $S_2(T)$ is obviously 1, whereas the correlation between $S_1(T)$ and $S_2(t)$ would be less than 1 as the increment $S_2(t) - S_2(T)$ would typically be only weakly dependent on $S_1(T)$.

the correct correlation parameter should originate from the *terminal* co-dependence between the two rates ($S_1(T), S_2(t)$) and should reflect both the correlation of rates observed at the same time *and* their inter-temporal de-correlation (which is typically quite significant).

One can attempt to deal with the issue above by specifying copula correlations (or, in general, copula parameters) that are functions of time t (in addition to them being functions of time T). Clearly, this is not particularly satisfactory as a large number of parameters would need to be kept and updated. Independent marking of such parameters would impose a heavy operational burden and, importantly, could introduce hard-to-trace arbitrage possibilities into the model. A better alternative, in our view, is to maintain only one copula correlation that corresponds to the dependence of the two rates ($S_1(T), S_2(T)$) observed on the same date, and then devise rules that would link other correlations to this “anchor” correlation. As we have done in the past, we may look at term structure models (that, by definition, are self-consistent in this regard) for inspiration. By resorting to approximations we limit the applicability of the method to only relatively small mismatches between the observation dates (about a year or so, probably) — for anything longer we would strongly recommend a direct application of a suitable multi-factor term structure model.

One reasonable approximation can be derived by assuming that $S_2(t) - S_2(T)$ is independent from $S_1(T)$ — a very respectable approximation — and that we can approximate the dynamics of $S_2(t)$ by a one-factor Gaussian mean-reverting process with constant volatility and mean reversion parameter κ_{S_2} , the same approximations that we employed in Section 16.8.6. Proceeding as in that section, we obtain (compare to (16.118))

$$\text{Corr}(S_1(T), S_2(t)) = \text{Corr}(S_1(T), S_2(T)) \times \sqrt{\frac{1 - e^{-2\kappa_{S_2}T}}{1 - e^{-2\kappa_{S_2}t}}}.$$

With this parameterization, correlation of $(S_1(T), S_2(T))$ defines the overall level of correlations for rates observed on different dates, while the mean reversion κ_{S_2} defines speed of further de-correlation arising from fixing date mismatches, providing a parsimonious yet flexible description of the whole universe of various correlations.

17.6 Numerical Methods for Copulas

Let us now turn our attention to issues of numerical implementation of valuation methods for copulas. The (undiscounted) value of a derivative with the payoff $f(S_1(t_1), \dots, S_d(t_d))$ paid at time T_p , where $S_i(t_i)$ is a swap rate observed at time t_i , $i = 1, \dots, d$, is equal to

$$V(0) = E^{T_p}(f(S_1(t_1), \dots, S_d(t_d))). \quad (17.31)$$

Letting $\psi(x_1, \dots, x_d)$ be the joint probability density of the swap rates $S_1(t_1), \dots, S_d(t_d)$ under the T_p -forward measure⁶, then the value can be represented as an integral

$$V(0) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, \dots, x_d) \psi(x_1, \dots, x_d) dx_1 \dots dx_d.$$

If the dependence structure between the swap rates is defined by a copula $C(u_1, \dots, u_d)$ then, according to (17.20),

$$\begin{aligned} V(0) &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, \dots, x_d) \\ &\quad \times c(\Psi_1(x_1), \dots, \Psi_d(x_d)) \left(\prod_{i=1}^d \psi_i(x_i) \right) dx_1 \dots dx_d, \end{aligned} \quad (17.32)$$

where the copula density $c(u_1, \dots, u_d)$ is defined by (17.21) and where we have denoted marginal PDFs and CDFs of the swap rates under the T_p -forward measure by $\psi_i(x)$ and $\Psi_i(x)$, $i = 1, \dots, d$, respectively. The formula (17.32) is the basic valuation formula for the copula method; changing to variables $u_i = \Psi_i(x_i)$ we can rewrite it as another useful formula,

$$V(0) = \int_0^1 \cdots \int_0^1 f(\Psi_1^{-1}(u_1), \dots, \Psi_d^{-1}(u_d)) c(u_1, \dots, u_d) du_1 \dots du_d. \quad (17.33)$$

17.6.1 Numerical Integration Methods

A quadrature rule approximates an integral with a finite sum of the values of the integrand over a suitably chosen grid covering the domain of integration. While not strictly required, the integration grid is often chosen to be a direct product of one-dimensional integration grids, so that we have something like

$$\begin{aligned} &\int \cdots \int g(y_1, \dots, y_d) dy_1 \dots dy_d \\ &\approx \sum_{m_1=1}^{M_1} \cdots \sum_{m_d=1}^{M_d} \mu_{m_1, \dots, m_d} g(y_{1,m_1}, \dots, y_{d,m_d}) \end{aligned} \quad (17.34)$$

for the grid $\{y_{1,m_1}\} \times \dots \times \{y_{d,m_d}\}$ and weights $\{\mu_{m_1, \dots, m_d}\}$. The presence of nested sums makes (17.34) impractical in high dimensions, as the number of points grows exponentially in the number of dimensions (“curse of dimensionality”). However, in the practically important case where d is

⁶Note we drop the superscript T_p for notational convenience for the duration of this section.

small (say, $d = 2$ or 3), the integrals in (17.32) or (17.33) can be computed quite efficiently with various schemes of the type (17.34). A good selection of methods is reviewed in Press et al. [1992], and many schemes are implemented in numerical software packages. With such pre-canned routines readily available, numerical integration for the copula method may therefore seem as straightforward as calling a suitable black-box procedure for the integral we are interested in; however, a robust, efficient implementation requires a bit more thought.

The first decision we need to make is which of the two integration formulas (17.32) and (17.33) to use. In (17.33) the limits of integration are finite, which simplifies discretization, and the marginal PDFs and CDFs $\psi_i(x)$, $\Psi_i(x)$ are not required when evaluating the integrand. On the other hand, (17.33) requires an efficient algorithm for calculating the inverses of marginal CDFs $\Psi_i^{-1}(u)$. Ultimately, whether (17.32) or (17.33) is most convenient will often depend on the specifics of the model at hand. For concreteness, we here choose (17.33) for our discussion.

Inverse CDFs that appear in (17.33) are rarely available in closed form and typically must be calculated numerically. For the sake of efficiency, these inverse CDFs should always be pre-computed before the main integration starts. This could be (for a given dimension i) as simple as caching $v_j = \Psi_i(\xi_j)$ for ξ_j on a given grid $\xi_1 < \dots < \xi_J$ that spans the domain of $\Psi_i(\cdot)$. In the main integration computation one might then approximate

$$\Psi_i^{-1}(u) \approx \frac{v_{n+1} - u}{v_{n+1} - v_n} \xi_n + \frac{u - v_n}{v_{n+1} - v_n} \xi_{n+1},$$

with n such that $v_n \leq u < v_{n+1}$. A more refined approach would augment this by inserting extra points in the intervals $[\xi_n, \xi_{n+1}]$ for which gaps $v_{n+1} - v_n$ are larger than a pre-specified tolerance $\epsilon > 0$.

With the inverse CDFs pre-computed and stored in caches, each evaluation of the integrand in (17.33) would consist of d lookups in inverse CDF caches, together with a (usually quite straightforward) evaluation of the payoff $f(\cdot)$ and a computation of the copula density function $c(\cdot)$. As the copula density function in most cases is known analytically, cache lookups will often dominate the evaluation time, suggesting that the integration grid be organized in such a way that cache lookups are efficient.

As we sometimes prefer not to use adaptive integration schemes (see Section 23.2.1 for explanation), we need to pay attention to the smoothness properties of the integrand, in particular the payoff f (as the copula density function c is usually smooth enough). Clearly, for most payoffs f would be either discontinuous (digital options on CMS spread, say) or, if continuous, then not differentiable (put and call options on CMS spread, say). Many relevant strategies for dealing with non-smooth payoffs are discussed in Chapter 23, and we urge the reader to get acquainted with the material in that chapter before proceeding with an actual implementation. For the purposes of our discussion here, we just observe that in the integral (17.33)

or, rather, the generic scheme (17.34), we often must treat the innermost integration (summation) differently from the outer integrals. To explain, let us set $d = 2$ and consider the nested integral

$$\int_0^1 \left(\int_0^1 g(u_1, u_2) du_1 \right) du_2,$$

where we use the short-hand notation

$$g(u_1, u_2) = f(\Psi_1^{-1}(u_1), \Psi_2^{-1}(u_2))c(u_1, u_2).$$

If g is non-smooth, numerical computation of the inner integral

$$\int_0^1 g(u_1, u_2) du_1$$

will generally require us to identify the points where $g(\cdot, u_2)$ is singular and either include them in the integration grid or increase the density of grid points around these points (see Chapter 23). It is therefore critical that we use an integration scheme that allows us complete freedom in locating the integration nodes⁷. On the other hand, when we calculate

$$\int_0^1 G(u_2) du_2,$$

where $G(u_2)$ is (the numerically calculated value of) the integral $\int_0^1 g(u_1, u_2) du_1$, the integrand G is often smooth, as the singularities of g would have been integrated out in the inner integral. Hence, for the outer integration, we can use fast schemes suitable for smooth integrands. For integrating over a finite interval (such as $[0, 1]$), the *Gauss-Lobatto* quadrature⁸ (see Kythe and Schäferkotter [2004]) is a good choice.

While integrating in one of the dimensions may cure singularities, this cannot be the case if the non-smooth features of the payoff are defined solely by the “outer” variable; in the example above that would be the case if we had $g(u_1, u_2) = 1_{\{u_2 > K\}}$ or something similar. This situation can be handled by switching the order of integration and integrating in variable u_2 first. For most cases, in fact, integrating in one particular dimension gives a smoother function than in any other, so an advanced integration routine would have logic to determine the dimension on which to perform the innermost quadrature. Of course there are situations when the payoff is non-smooth in both directions; careful handling of singularities in both integrals is then required.

⁷A trapezoidal scheme is a possible candidate. Most Gaussian quadrature rules, however, are not ideal, since their integration grids are not directly user-specifiable, but emerge as roots of a particular function.

⁸Sometimes also called *Gauss-Legendre* quadrature after the name of the family of polynomials whose roots define the integration grid.

Before concluding our remarks on integration methods for copulas, we briefly consider the specialization of the method for Gaussian copulas (see Section 17.3.1). For $C = C_{\text{gauss}}$, we recall the alternative valuation formula (17.11),

$$\begin{aligned} V(0) &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(\Psi_1^{-1}(\Phi(z_1)), \dots, \Psi_d^{-1}(\Phi(z_d))) \\ &\quad \times \phi(z_1, \dots, z_d; R) dz_1 \dots dz_d. \end{aligned}$$

Using this representation can lead to numerical improvements, since each nested integral is now represented as an integral over a Gaussian density. For such integrals, the Gauss-Hermite quadrature method (already mentioned in Section 12.3.4.3; see also Press et al. [1992]) is often very efficient. Of course, we would recommend this quadrature rule only if the integrand is sufficiently smooth (e.g., for outer integrals). Also, our recommendation to cache inverse CDFs still apply.

17.6.2 Dimensionality Reduction for CMS Spread Options

The discussion in the previous section is generic and applies to arbitrary multi-rate payoffs. In the important special case of CMS spread options, we can achieve further efficiencies of numerical implementation by reformulating the problem as one involving only *one-dimensional* integrals. The following proposition (following Dhaene and Goovaerts [1996] and Berrahoui [2005]) contains the relevant result.

Proposition 17.6.1. *Let us denote*

$$\gamma(x, K) = \frac{d}{dx} \Psi(x, x - K),$$

where $\Psi(x_1, x_2)$ is the joint CDF of $(S_1(T), S_2(T))$ under Q^{T_r} . Then, the undiscounted value of a spread option as defined by (17.26) is given by

$$V(0; T, K) = \int_{-\infty}^{+\infty} x \gamma(x, K) dx - \int_{-\infty}^{+\infty} x \psi_2(x) dx - K, \quad (17.35)$$

where $\psi_2(x)$ is the one-dimensional marginal PDF of $S_2(T)$ under Q^{T_r} .

Proof. We have

$$\begin{aligned} V(0; T, K) &= \int_{-\infty}^{+\infty} \left(\int_{-\infty}^{x_1 - K} (x_1 - x_2 - K) \psi(x_1, x_2) dx_2 \right) dx_1 \\ &= \int_{-\infty}^{+\infty} x_1 \left(\int_{-\infty}^{x_1 - K} \psi(x_1, x_2) dx_2 \right) dx_1 \\ &\quad - \int_{-\infty}^{+\infty} (x_2 + K) \left(\int_{x_2 + K}^{+\infty} \psi(x_1, x_2) dx_1 \right) dx_2. \end{aligned}$$

Recall that $\psi(x_1, x_2) = \frac{\partial^2}{\partial x_1 \partial x_2} \Psi(x_1, x_2)$. Therefore

$$\begin{aligned} V(0; T, K) &= \int_{-\infty}^{+\infty} x_1 \left(\frac{\partial}{\partial x_1} \Psi(x_1, x_1 - K) - \frac{\partial}{\partial x_1} \Psi(x_1, -\infty) \right) dx_1 \\ &\quad - \int_{-\infty}^{+\infty} (x_2 + K) \left(\frac{\partial}{\partial x_2} \Psi(+\infty, x_2) - \frac{\partial}{\partial x_2} \Psi(x_2 + K, x_2) \right) dx_2. \end{aligned}$$

Note that

$$\frac{\partial}{\partial x_1} \Psi(x_1, -\infty) = 0, \quad \frac{\partial}{\partial x_2} \Psi(+\infty, x_2) = \psi_2(x_2).$$

Then

$$\begin{aligned} V(0; T, K) &= \int_{-\infty}^{+\infty} x_1 \frac{\partial}{\partial x_1} \Psi(x_1, x_1 - K) dx_1 \\ &\quad + \int_{-\infty}^{+\infty} (x_2 + K) \frac{\partial}{\partial x_2} \Psi(x_2 + K, x_2) dx_2 - \int_{-\infty}^{+\infty} x \psi_2(x) dx - K. \quad (17.36) \end{aligned}$$

We also have

$$\gamma(x, K) = \frac{\partial}{\partial x_1} \Psi(x, x - K) + \frac{\partial}{\partial x_2} \Psi(x, x - K)$$

and

$$\begin{aligned} \int_{-\infty}^{+\infty} x \gamma(x, K) dx &= \int_{-\infty}^{+\infty} x \frac{\partial}{\partial x_1} \Psi(x, x - K) dx \\ &\quad + \int_{-\infty}^{+\infty} x \frac{\partial}{\partial x_2} \Psi(x, x - K) dx. \end{aligned}$$

By substituting $x = x_1$ in the first integral and $x = x_2 + K$ in the second, the proposition now follows from (17.36). \square

In the case when the joint CDF is generated by a copula $C(u_1, u_2)$, i.e. when $\Psi(x_1, x_2) = C(\Psi_1(x_1), \Psi_2(x_2))$, we have

$$\begin{aligned} \gamma(x, K) &= \frac{\partial}{\partial u_1} C(\Psi_1(x), \Psi_2(x - K)) \psi_1(x) \\ &\quad + \frac{\partial}{\partial u_2} C(\Psi_1(x), \Psi_2(x - K)) \psi_2(x - K). \end{aligned}$$

Finally, let us present the result of Proposition 17.6.1 in a somewhat different form, already used in the discussion in Section 17.4.4.

Corollary 17.6.2. *The undiscounted value $V(0; T, K)$ at time $t = 0$ of a spread option with strike K at expiry T as defined by (17.26) is given by*

$$V(0; T, K) = \int_{-\infty}^{\infty} (1_{\{x>0\}} - \Psi(x, x - K)) dx - E^{T_p}(S_2(T)) - K. \quad (17.37)$$

Alternatively it can be expressed in terms of the copula function,

$$V(0; T, K) = \int_{-\infty}^{\infty} (1_{\{x>0\}} - C(\Psi_1(x), \Psi_2(x - K))) dx - E^{T_p}(S_2(T)) - K. \quad (17.38)$$

Proof. Integrating by parts we obtain

$$\begin{aligned} \int_{-\infty}^{\infty} x \gamma(x, K) dx &= \int_{-\infty}^{\infty} x \left(\frac{d}{dx} \Psi(x, x - K) \right) dx \\ &= \int_{-\infty}^0 x d\Psi(x, x - K) - \int_0^{\infty} x d(1 - \Psi(x, x - K)) \\ &= - \int_{-\infty}^0 \Psi(x, x - K) dx + \int_0^{\infty} (1 - \Psi(x, x - K)) dx, \end{aligned}$$

and (17.37) follows. The result (17.38) follows from the definition of the copula function. \square

17.6.3 Dimensionality Reduction for Other Multi-Rate Derivatives

The formula (17.35) serves as a base for deriving similar one-dimensional representations for other derivatives. For example, differentiating the definition

$$V_{\text{spread}}(0; T, K) = E^{T_p} \left((S_1(T) - S_2(T) - K)^+ \right)$$

with respect to the strike and exchanging the order of the expected value operator and differentiation, we obtain

$$E^{T_p} (1_{\{S_1(T) - S_2(T) > K\}}) = - \frac{\partial}{\partial K} V_{\text{spread}}(0; T, K).$$

On the left-hand side we have the value of a *digital* spread option, and the expression on the right-hand side may be computed as the one-dimensional integral obtained by differentiating (17.35) with respect to K ,

$$\begin{aligned} E^{T_p} (1_{\{S_1(T) - S_2(T) > K\}}) &= - \frac{\partial}{\partial K} V_{\text{spread}}(0; T, K) \\ &= - \int_{-\infty}^{+\infty} x \frac{\partial \gamma(x, K)}{\partial K} dx + 1. \end{aligned}$$

To obtain further results, we first generalize Proposition 17.6.1 to spread options with arbitrary gearings.

Corollary 17.6.3. *Let $a_1, a_2 > 0$. Then*

$$\begin{aligned} \mathbb{E}^{T_r} \left((a_1 S_1(T) - a_2 S_2(T) - K)^+ \right) \\ = \int_{-\infty}^{+\infty} x \left(\frac{d}{dx} \Psi(x/a_1, (x-K)/a_2) \right) dx \\ - \int_{-\infty}^{+\infty} x \left(\frac{d}{dx} \Psi(+\infty, x/a_2) \right) dx - K, \quad (17.39) \end{aligned}$$

where $\Psi(x_1, x_2)$ is, as before, the joint CDF of $(S_1(T), S_2(T))$ under \mathbb{Q}^{T_r} .

We can now differentiate the formula (17.39) with respect to a_1 and a_2 , yielding one-dimensional integral representations for the values of derivatives with the payoffs

$$S_1(T) \mathbf{1}_{\{a_1 S_1(T) - a_2 S_2(T) > K\}}, \quad S_2(T) \mathbf{1}_{\{a_1 S_1(T) - a_2 S_2(T) > K\}}. \quad (17.40)$$

We leave the detailed derivation of these integrals as an exercise to the reader, and just note that the payoffs in (17.40) are those of floating digital spread options. These options are not only important by themselves but are also components of floating spread range accruals, as explained in Section 17.5. As long as the payment rate of the range accrual is either S_1 or S_2 (i.e. equal to one of the rates that define the spread), the value of this security can be obtained by one-dimensional integration. We notice that the specification (17.40) also includes (non-spread) floating digitals, as we can produce the payoff

$$S_1(T) \mathbf{1}_{\{S_2(T) > K\}}$$

by setting $a_1 = 0, a_2 = -1$.

The valuation expression for a floating digital is, in fact, easy — and instructive — to derive directly. The key here is the following important lemma.

Lemma 17.6.4. *If the distribution of $S_1(T), S_2(T)$ under \mathbb{Q}^{T_r} is given by the joint CDF $\Psi(x_1, x_2)$ with the copula function $C(u_1, u_2)$, then the conditional CDF of $S_2(T)$ given $S_1(T)$ is directly computable from the copula,*

$$\mathbb{Q}^{T_r}(S_2(T) < x_2 | S_1(T) = x_1) = \frac{\partial C}{\partial u_1}(\Psi_1(x_1), \Psi_2(x_2)),$$

where Ψ_1, Ψ_2 are the marginal CDFs of S_1 and S_2 .

Proof. We have

$$\begin{aligned} \mathbb{Q}^{T_r}(S_2(T) < x_2 | S_1(T) = x_1) \\ = \frac{\mathbb{E}^{T_r}(\delta(S_1(T) - x_1) \mathbf{1}_{\{S_2(T) < x_2\}})}{\mathbb{E}^{T_r}(\delta(S_1(T) - x_1))} = \frac{\partial \Psi(x_1, x_2) / \partial x_1}{\psi_1(x_1)}, \end{aligned}$$

where $\delta(\cdot)$ is the Dirac delta function and ψ_1 is the PDF of S_1 . On the other hand,

$$\frac{\partial \Psi(x_1, x_2)}{\partial x_1} = \frac{\partial C(\Psi_1(x_1), \Psi_2(x_2))}{\partial x_1} = \frac{\partial C(\Psi_1(x_1), \Psi_2(x_2))}{\partial u_1} \psi_1(x_1),$$

and the result follows. \square

Conditioning the payoff of the floating digital on $S_1(T)$ and applying the lemma, we obtain

$$\begin{aligned} \mathbb{E}^{T_p}(S_1(T)1_{\{S_2(T)>K\}}) &= \int_{-\infty}^{\infty} x Q^{T_p}(S_2(T) > K | S_1(T) = x) \psi_1(x) dx \\ &= \int_{-\infty}^{\infty} x \left(1 - \frac{\partial C}{\partial u_1}(\Psi_1(x), \Psi_2(K)) \right) \psi_1(x) dx. \end{aligned} \quad (17.41)$$

Remark 17.6.5. The result of Lemma 17.6.4 provides an alternative, perhaps more elegant, route to the proof of Proposition 17.6.1 (or, more generally, Corollary 17.6.3). To see this, we write for the spread option value $V(0; T, K)$,

$$\begin{aligned} V(0; T, K) &= \mathbb{E}^{T_p} (S_1(T)1_{\{S_1(T)-S_2(T)-K \geq 0\}}) \\ &\quad - \mathbb{E}^{T_p} (S_2(T)1_{\{S_1(T)-S_2(T)-K \geq 0\}}) \\ &\quad - KE^{T_p} (1_{\{S_1(T)-S_2(T)-K \geq 0\}}). \end{aligned}$$

Conditioning the first term on $S_1(T)$ and the second on $S_2(T)$ (and the third one on either $S_1(T)$ or $S_2(T)$) and using the result of Lemma 17.6.4, a one-dimensional integral representation for the spread option is obtained.

17.6.4 Dimensionality Reduction by Conditioning

Reducing the dimensionality of integrals is often an effective technique for improving computational performance and/or for extending the domain of applicability of direct integration methods. For some payoffs this can be achieved by application of the method in Sections 17.6.2 and 17.6.3; others require different approaches. In Lemma 17.6.4 we demonstrated a particular application of the principle of *conditioning*, an idea that has quite general applicability. The gist of the method is simple: if we can calculate (or approximate) in closed form the expectation of the payoff conditioned on some b -dimensional subset of the d variables in the payoff definition, then we can reduce the dimension of the valuation integral from d to b .

Gaussian random variables are particularly amendable to conditioning methods, courtesy of Lemma 14.6.5 that we have already applied in the context of Brownian bridge calculations, and for calculating swaption values in a two-dimensional Gaussian model (see Section 12.1.6.1). The lemma helps us to calculate the distribution of one Gaussian random variable conditioned

on another; importantly, the conditional distribution remains Gaussian. Let us demonstrate how the lemma could be applied to the problem of dimensionality reduction, by considering a simple floating digital payoff⁹

$$S_1(T) \mathbf{1}_{\{S_2(T) > K\}}, \quad (17.42)$$

where S_1 and S_2 are two swap rates, for notational simplicity assumed to pay at T . The log-normal model would specify

$$S_i(T) = \tilde{S}_i e^{\sigma_i Z_i - \sigma_i^2/2}, \quad i = 1, 2, \quad (17.43)$$

where $\tilde{S}_i = E^T(S_i(T))$ are CMS-adjusted forward swap rates, σ_i are unscaled¹⁰ log-normal volatilities of the swap rates, and the vector $(Z_1, Z_2)^\top$ is Gaussian with zero mean, unit variance and correlation ρ (all under Q^T). The (undiscounted) value of the derivative with payoff (17.42) is given by

$$V(0) = \tilde{S}_1 E^T \left(e^{\sigma_1 Z_1 - \sigma_1^2/2} \mathbf{1}_{\{Z_2 > \ln(K/\tilde{S}_2)/\sigma_2 + \sigma_2/2\}} \right). \quad (17.44)$$

A direct evaluation of (17.42) would require two-dimensional integration, but it is easy to see that by applying the result of Lemma 14.6.5 we can reduce the dimension of the integral to one. In particular, if we condition on Z_2 , we have,

$$\begin{aligned} V &= \tilde{S}_1 E^T \left(E^T \left(e^{\sigma_1 Z_1 - \sigma_1^2/2} \mathbf{1}_{\{Z_2 > \ln(K/\tilde{S}_2)/\sigma_2 + \sigma_2/2\}} \middle| Z_2 \right) \right) \\ &= \tilde{S}_1 E^T \left(\mathbf{1}_{\{Z_2 > \ln(K/\tilde{S}_2)/\sigma_2 + \sigma_2/2\}} E^T \left(e^{\sigma_1 Z_1 - \sigma_1^2/2} \middle| Z_2 \right) \right). \end{aligned} \quad (17.45)$$

Since

$$Z_1 | Z_2 \sim \mathcal{N}(Z_2 \rho, 1 - \rho^2), \quad (17.46)$$

we can evaluate the conditional expected value analytically,

$$E^T \left(e^{\sigma_1 Z_1 - \sigma_1^2/2} \middle| Z_2 \right) = e^{\sigma_1 \rho Z_2 - \sigma_1^2 \rho^2/2}, \quad (17.47)$$

which gives us

$$V(0) = \tilde{S}_1 e^{-\sigma_1^2 \rho^2/2} E^T \left(e^{\rho \sigma_1 Z_2} \mathbf{1}_{\{Z_2 > \ln(K/\tilde{S}_2)/\sigma_2 + \sigma_2/2\}} \right).$$

Now we need to integrate the modified payoff against the distribution of Z_2 only, which is a one-dimensional integration. Of course, in the log-normal model, this will give the same result as the formula (17.41).

⁹Note that in Section 17.6.3 we already derived a one-dimensional integral representation for this payoff, but we use it anyway to demonstrate the main idea of the method we develop here.

¹⁰Not annualized, i.e. not divided by \sqrt{T} .

As should be clear from the above, successful applications of conditioning methods rely to a large degree on the availability of the appropriate analytical tools such as Lemma 14.6.5. For general copulas we can calculate conditional CDFs quite easily, as Lemma 17.6.4 demonstrated. This can take us a long way, but to gain analytical tractability, we may ultimately need to consider a narrower range of copulas. In particular, if one is content to use Gaussian copulas, or combinations thereof, then the techniques discussed above can be applied virtually unchanged. For example, conditional CDFs in Gaussian copulas are available in closed form by direct application of Lemma 14.6.5. We leave the (simple) derivation of this generic result to the reader, while focusing here on the applications of the method to some specific payoffs.

As we recall, the Gaussian copula method essentially replaces (17.43) with a more generic expression

$$S_i(T) = \Lambda_i(Z_i), \quad i = 1, 2, \quad (17.48)$$

where $\Lambda_i(x) \triangleq \Psi_i^{-1}(\Phi(x))$ are “mapping functions” from Gaussian variates to market rates. We have already shown (Lemma 17.6.4) how to calculate the value of the payoff (17.42) by conditioning on S_1 . Conditioning on S_2 , as done for the Gaussian case in (17.45), is also straightforward. Specifically, we obtain

$$V(0) = E^T (1_{\{\Lambda_2(Z_2) > K\}} E^T (\Lambda_1(Z_1) | Z_2)),$$

and we observe that it is necessary to calculate, or approximate, $E^T(\Lambda_1(Z_1) | Z_2)$. A simple approximation may be obtained by replacing $\Lambda_1(\cdot)$ with a quadratic function:

$$\Lambda_1(x) \approx \Lambda_1(0) + \Lambda'_1(0)x + \frac{1}{2}\Lambda''_1(0)x^2,$$

so that we would have

$$E^T (\Lambda_1(Z_1) | Z_2) \approx \Lambda_1(0) + \Lambda'_1(0)E^T (Z_1 | Z_2) + \frac{1}{2}\Lambda''_1(0)E^T (Z_1^2 | Z_2),$$

where, from (17.46),

$$E^T (Z_1 | Z_2) = \rho Z_2, \quad E^T (Z_1^2 | Z_2) = 1 + \rho^2(Z_2^2 - 1).$$

This could be refined slightly by expanding not around $x = 0$ but around $x^* = \Lambda_1^{-1}(\tilde{S}_1)$, where \tilde{S}_1 is the forward CMS-adjusted swap rate.

A more general approach of expanding $\Lambda_1(x)$ into a Taylor series of arbitrary order is also possible and not much more work, as all required terms of the type $E^T(Z_1^r | Z_2)$, $r \geq 1$, are easily available from (17.46). Finally, a related, and very accurate, method is based on approximating $\Lambda_1(x)$ with a truncated cosine series. If we choose a range $[-z_{\max}, z_{\max}]$ such that $Q^T(Z_1 \in [-z_{\max}, z_{\max}]) = 1 - \epsilon$ for small $\epsilon > 0$ (for example we can take z_{\max} equal to 3 or 4), then, since $\Lambda_1(x)$ is typically continuous, we can approximate it for $x \in [-z_{\max}, z_{\max}]$ arbitrarily closely by a sum of the type

$$\Lambda_1(x) \approx \text{Re} \left(\sum_{m=0}^M w_m e^{\lambda_m x - \lambda_m^2 / 2} \right), \quad (17.49)$$

where

$$\lambda_m = \frac{\pi i m}{2z_{\max}}, \quad m = 0, \dots, M,$$

and $i = \sqrt{-1}$. The weights w_m could be quickly computed by an inverse discrete cosine transform, see Press et al. [1992]. The required conditional expected value is then given by (see (17.47))

$$\mathbb{E}^T (\Lambda_1(Z_1) | Z_2) \approx \text{Re} \left(\sum_{m=0}^M w_m e^{\lambda_m \rho Z_2 - \lambda_m^2 \rho^2 / 2} \right).$$

The expansion methods do not always work; for example, a payment amount of a floating digital could be a function of the rate S_1 rather than equal to the rate itself. To consider a typical example, let us analyze the expectation

$$V(0) = \mathbb{E} ((S_1 - u)^+ 1_{\{S_2 > K\}}),$$

essentially a call option on S_1 conditioned on S_2 being above some level. Conditioning on Z_2 we see that we need to calculate the expected value

$$\mathbb{E} ((\Lambda_1(Z_1) - u)^+ | Z_2). \quad (17.50)$$

In the case where S_1 is log-normal, this conditional expected value would present no difficulties and would be given by the Black formula (with some parameters dependent on Z_2). Matters are more complicated in the general case, however. To proceed, we recall that $Z_1 = \rho Z_2 + \bar{\rho} \xi$, $\bar{\rho} = (1 - \rho^2)^{1/2}$, where ξ is a standard Gaussian random variable independent of Z_2 . Hence, for a fixed value of Z_2 , the task of evaluating (17.50) reduces to that of calculating

$$\mathbb{E} ((\Lambda_1(a + b\xi) - u)^+) \quad (17.51)$$

for some a, b .

To proceed further, it is important to recall that the values of options of the type $e(v) = \mathbb{E}((\Lambda_1(\xi) - v)^+)$ are easily available to us, as these are just the option values in the marginal model we use for the rate S_1 . With that in mind, we can calculate the option value in (17.51) by replicating the payoff with $e(v)$, $v \in \mathbb{R}$, following the ideas of Section 16.6.1. However, it is not clear if we can achieve substantial computational savings — the original aim of the conditioning method — along this route as the replication method requires calculation of an integral. To save computational effort, suppose we limit ourselves to a single strike in the replication. A natural choice of that strike is such v that the two payoffs

$$(\Lambda_1(a + bx) - u)^+, \quad (\Lambda_1(x) - v)^+$$

have a discontinuity at exactly the same point x . Such a point is given by

$$v^* = \Lambda_1(x^*), \quad x^* = (\Lambda_1^{-1}(u) - a) / b.$$

Then, by matching the slope of both payoffs at the critical point x^* , we obtain a single-strike “replication”,

$$(\Lambda_1(a + bx) - u)^+ \approx \frac{b\Lambda'_1(a + bx^*)}{\Lambda'_1(x^*)} (\Lambda_1(x) - \Lambda_1(x^*))^+$$

which, when applied to the conditional expected value in (17.50), gives us the following approximation

$$\mathbb{E}((\Lambda_1(Z_1) - u)^+ | Z_2) \approx \bar{\rho} \frac{\Lambda'_1(\rho Z_2 + \bar{\rho}x^*(Z_2))}{\Lambda'_1(x^*(Z_2))} e(\Lambda_1(x^*(Z_2))).$$

where the critical point $x^*(Z_2)$ depends on Z_2 and is given by

$$x^*(Z_2) = (\Lambda_1^{-1}(u) - \rho Z_2) / \bar{\rho}.$$

The accuracy of the method obviously depends on how far the mapping function $\Lambda_1(\cdot)$ is from identity.

17.6.5 Dimensionality Reduction by Measure Change

Methods based on conditioning are not the only choice for dimensionality reduction. Another useful approach is based in the idea of performing calculations under a different measure to simplify the payoff. The main technical tool here is the following specialization of the Girsanov theorem (see Theorem 1.5.1).

Lemma 17.6.6. *If X is a d -dimensional Gaussian vector with mean μ and covariance matrix Σ , v is a d -dimensional vector and $f(\cdot)$ is a function $\mathbb{R}^d \rightarrow \mathbb{R}$, then*

$$\mathbb{E}\left(e^{v^\top(X-\mu)-v^\top\Sigma v/2} f(X)\right) = \widehat{\mathbb{E}}(f(X)),$$

where in measure $\widehat{\mathbb{Q}}$, vector X has mean $\mu + \Sigma v$ and covariance matrix Σ .

Proof. See Section 3.5 of Karatzas and Shreve [1997]. \square

To see how this method works, let us continue with the payoff (17.42) and return to the log-normal model (17.43) for the moment. Applying Lemma 17.6.6 to the problem (17.44), we immediately obtain the following one-dimensional representation:

$$V(0) = \tilde{S}_1 \widehat{\mathbb{E}}\left(1_{\{Z_2 > \ln(K/\tilde{S}_2)/\sigma_2 + \sigma_2/2\}}\right), \quad (17.52)$$

where Z_2 under \hat{Q} is Gaussian but has a different drift,

$$Z_2 \sim \mathcal{N}(\sigma_1 \rho, 1).$$

As desired, the problem has been reduced to that of a one-dimensional integration (in fact, in this case, the expectation is available in closed form). While this example is rather simple, many other payoffs are amendable to the same type of treatment.

As was the case for the conditioning method, the measure change method extends to the general Gaussian copula setup. In the model (17.48), with the approximation (17.49), we have for the value of the payoff (17.42),

$$\begin{aligned} V(0) &= \mathbb{E}^T (\Lambda_1(Z_1) \mathbf{1}_{\{\Lambda_2(Z_2) > K\}}) \\ &\approx \operatorname{Re} \left(\mathbb{E}^T \left(\left(\sum_{m=0}^M w_m e^{\lambda_m Z_1 - \lambda_m^2 / 2} \right) \mathbf{1}_{\{\Lambda_2(Z_2) > K\}} \right) \right) \\ &= \operatorname{Re} \left(\sum_{m=0}^M w_m \mathbb{E}^T \left(e^{\lambda_m Z_1 - \lambda_m^2 / 2} \mathbf{1}_{\{\Lambda_2(Z_2) > K\}} \right) \right). \end{aligned}$$

Applying Lemma 17.6.6 to each term in turn, we obtain

$$V(0) \approx \operatorname{Re} \left(\sum_{m=0}^M w_m \hat{\mathbb{E}}^m \left(\mathbf{1}_{\{\Lambda_2(Z_2) > K\}} \right) \right), \quad (17.53)$$

where under measure $\hat{\mathbb{Q}}^m$, Z_2 is Gaussian with variance 1 and the mean $\lambda_m \sigma_1 \rho$ for each $m = 1, \dots, M$. While it may seem strange to have Gaussian random variables with a (purely) imaginary mean, each term $\hat{\mathbb{P}}^m(\Lambda_2(Z_2) > K)$ in (17.53) is actually well-defined as an analytic continuation of the Gaussian CDF into a complex domain. Let us demonstrate on a simple example. Let Z be a standard Gaussian random variable. Then

$$\begin{aligned} \mathbb{E}^T \left(e^{iZ+1/2} \mathbf{1}_{\{Z>K\}} \right) &= \frac{1}{\sqrt{2\pi}} \int_K^\infty e^{ix+1/2} e^{-x^2/2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_K^\infty e^{-(x-i)^2/2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{K-i}^{\infty-i} e^{-z^2/2} dz, \end{aligned}$$

where the last integral is understood to be over the contour $\{x - i : x \in [K, \infty)\}$ in the complex plane and is well-defined, see Chapter 7 of Abramowitz and Stegun [1965].

The method of (17.53) requires m one-dimensional integrations, which is an improvement — as long as m is not too large — over one *two*-dimensional integration that would be required for a standard valuation of (17.42).

17.6.6 Monte Carlo Methods

Let us return to the general problem of calculating the value of a multi-rate payoff in (17.31). We observe again that if the dimensionality d of the payoff is higher than 3 (after potential dimensionality reductions by conditioning and measure change methods), nested numerical integration of the type (17.34) becomes unattractive compared to direct Monte Carlo simulation of the d -dimensional joint distribution. A particularly simple Monte Carlo scheme is obtained if the copula used is Gaussian (with correlation matrix R). Recalling (17.6), it should be clear that we can calculate the (undiscounted) value of the derivative as

$$V(0) \approx \frac{1}{N} \sum_{n=1}^N f(\Psi_1^{-1}(\Phi(Z_{n,1})), \dots, \Psi_d^{-1}(\Phi(Z_{n,d}))), \quad (17.54)$$

where $\mathbf{Z}_1, \dots, \mathbf{Z}_N$, with $\mathbf{Z}_n = (Z_{n,1}, \dots, Z_{n,d})$, are N independent samples from a d -dimensional Gaussian distribution.

The case of a non-Gaussian copula is conceptually similar. Given a copula function $C(u_1, \dots, u_d)$, assume momentarily that we know how to generate a random sample for the vector $\mathbf{U} = (U_1, \dots, U_d)$, where each U_i has a uniform distribution on $[0, 1]$, and the dependence structure is given by the copula C ,

$$Q^{T_n}(U_1 < u_1, \dots, U_d < u_d) = C(u_1, \dots, u_d). \quad (17.55)$$

Then, the value of the derivative in (17.33) is given by

$$V(0) \approx \frac{1}{N} \sum_{n=1}^N f(\Psi_1^{-1}(U_{n,1}), \dots, \Psi_d^{-1}(U_{n,d})), \quad (17.56)$$

where $\mathbf{U}_n = (U_{n,1}, \dots, U_{n,d})$ has the same distribution as \mathbf{U} for each $n = 1, \dots, N$, and all \mathbf{U}_n are independent. Calculations with formulas (17.54) or (17.56) are straightforward; we only remind the reader that inverse CDFs $\Psi_i^{-1}(u)$ should be pre-computed and cached before the main simulation.

The success of the numerical implementation of the scheme (17.56) hinges on our ability to simulate a random sample from a given copula as in (17.55). We have demonstrated how to do this for the Gaussian copula, so let us consider the Archimedean copulas that were also introduced in Section 17.3. The simulation algorithm for a bivariate Archimedean copula with the generator function $\omega(\cdot)$ can be based on the following result from Nelsen [2006].

Lemma 17.6.7. *Let (U_1, U_2) be a random vector with uniform marginals and joint distribution function $C_{\text{arch}}(u_1, u_2; \omega(\cdot))$. Define two new random variables*

$$R = \omega(U_1) / (\omega(U_1) + \omega(U_2)), \quad F = C_{\text{arch}}(U_1, U_2; \omega(\cdot)).$$

Then, the joint distribution function of (R, F) is given by

$$P(R < r, F < f) = r \times A_\omega(f), \quad A_\omega(f) \triangleq 1 - \omega(f)/\omega'(f).$$

Here, R and F are independent and R is uniformly distributed on $[0, 1]$.

With the help of this lemma, a sample (U_1, U_2) from an Archimedean copula can be generated with the following algorithm:

1. Simulate two independent random variables R and W , uniformly distributed on $[0, 1]$.
2. Set $F = A_\omega^{-1}(W)$, where $A_\omega(f) = 1 - \omega(f)/\omega'(f)$.
3. Set $U_1 = \omega^{-1}(R\omega(F))$ and $U_2 = \omega^{-1}((1 - R)\omega(F))$.

A multi-dimensional extension of this algorithm exists; see Wu et al. [2006] for details.

From our basic ability to simulate Gaussian and Archimedean copulas, we can devise simulation schemes for the “aggregate” copulas outlined in Section 17.3.4. First, let us consider the reflection method of Lemma 17.3.3. If $(U_1, \dots, U_i, \dots, U_d)$ is a sample from some copula $C(u_1, \dots, u_d)$, then, clearly, $(U_1, \dots, 1 - U_i, \dots, U_d)$ is a sample from the reflected copula $\bar{C}(u_1, \dots, u_d; \{i\})$; it follows that simulating reflected copulas is straightforward.

Simulation of a convex linear combination of copulas as given by Lemma 17.3.4 is also easy. To state the algorithm, let $\mathbf{U}^m = (U_1^m, \dots, U_d^m)$, $m = 1, \dots, M$, be a collection of independent samples from the copulas C_m ; and let W be a discrete random variable with the distribution $P(W = m) = w_m$, $m = 1, \dots, M$. If W is independent of all \mathbf{U}^m , then the sample

$$(U_1, \dots, U_d) = \mathbf{U} \triangleq \mathbf{U}^W = \sum_{m=1}^M \mathbf{U}^m 1_{\{W=m\}}$$

is a sample from the mixture copula

$$\sum_{m=1}^M w_m C_m(u_1, u_2, \dots, u_d).$$

Finally, let us turn to the product copulas as defined by Theorem 17.3.5. While they may appear more complicated than mixture copulas, product copulas are, in fact, quite straightforward to simulate. To present the basic idea, we first observe that if random variables X_m , $m = 1, \dots, M$, are independent, then

$$P\left(\max_{m=1, \dots, M} X_m < x\right) = P(X_1 < x, \dots, X_M < x) = \prod_{m=1}^M P(X_m < x). \tag{17.57}$$

As above, let $\mathbf{U}^m = (U_1^m, \dots, U_d^m)$, $m = 1, \dots, M$, be a collection of independent samples from the copulas C_m . As pointed out by Liebscher [2008], if we define

$$\mathbf{U} = (U_1, \dots, U_d), \quad U_i \triangleq \max_{m=1, \dots, M} q_{m,i}^{-1}(U_i^m), \quad i = 1, \dots, d,$$

it follows from (17.57) that

$$\begin{aligned} & P(U_1 < u_1, \dots, U_d < u_d) \\ &= P\left(\max_{m=1, \dots, M} q_{m,1}^{-1}(U_1^m) < u_1, \dots, \max_{m=1, \dots, M} q_{m,d}^{-1}(U_d^m) < u_d\right) \\ &= \prod_{m=1}^M P(U_1^m < q_{m,1}(u_1), \dots, U_d^m < q_{m,d}(u_d)), \end{aligned}$$

which is the product copula we wished to produce. We leave it to the reader to write down a step-by-step simulation scheme for the product copula, using the result above.

17.7 Limitations of the Copula Method

The copula method has gained widespread acceptance for multi-rate derivatives, in large part due to the ease with which a multivariate distribution consistent with market-observed marginal distributions can be constructed and parameterized. Before the reader starts assuming that copulas are a panacea, we should warn that copula applications have their share of limitations. First and foremost, we emphasize that copulas generally do not result in a dynamic model for the yield curve, nor are they consistent with the most popular classes of such models. As we have seen, the copula method allows us to easily ascribe a joint terminal distribution to a collection of CMS rates, for the purpose of pricing European multi-rate options. In the general case it is difficult, if not impossible, to find a dynamic term structure model that would be consistent with the joint terminal distributions produced by the copula method. To see why this might be problematic, consider a path-dependent exotic security for which a multi-rate vanilla is part of the payoff specification; in such a case, application of a copula for the “embedded” vanilla option would inherently be inconsistent with the dynamic model itself. Such an inconsistency may, among other ills, cause internal arbitrages in the model and nonsensical hedge ratios.

Also, before one gets carried away by the simplicity and convenience of copula parameterizations, it should be remembered that copulas in practical use are not usually chosen for their links to observed financial relationships, but instead for their technical properties and ease of implementation. While we (as should be clear to the reader by now) always welcome a dose of

pragmatism in financial modeling, choosing tools simply because they are easy to use, rather than because they make sense, is obviously a problematic idea. A related issue is the fact that many copulas have parameters which are either devoid of meaning, or are consistently misinterpreted by users. For instance, the parameters in a Gaussian copula, the entries in the correlation matrix R , are obviously *not* the actual correlations of the interest rates being modeled¹¹ — instead, the (linear) correlations between rates in the copula will depend on the marginal distribution of the rates. In effect, the meaning of the copula parameters (the matrix R) will change when the marginal distributions of the rates change. See Appendix 17.A for a concrete demonstration of this effect.

As a last warning about copulas, let us note that a typical copula uses just a few parameters to capture the often complex dependence structure of various rates over periods of time. In practice, the parameterization of any model almost invariably defines the rules for parameter interpolation: if a is a particular model parameter, then it is natural for the users to use constant or linear interpolation in a , either explicitly in a risk management system or implicitly in their heads, to fill in values between observations. In a copula setting, this can become a problem as naive interpolation of parameters across option expiries may lead to a distorted picture of how the dynamic dependence of two rates evolves¹². Potential problems of this kind are touched upon in Section 17.5.

Although care must be taken to avoid some of the pitfalls above, copulas clearly have a place in modeling of vanilla derivatives and, in any case, have managed to become a de-facto standard for some important products, such as European CMS spread options. Our message in this section is simply that the limitations of the method should be clearly understood, to ensure that it is applied effectively and appropriately. On this note, let us stop our discussion of copulas and look at alternative ways of introducing dependence between market rates.

17.8 Stochastic Volatility Modeling for Multi-Rate Options

As described in Chapter 16, a stochastic volatility model is often our preferred choice for pricing of single-rate derivatives such as swaptions and CMS-linked products. In the context of multi-rate derivatives, having the distribution of

¹¹ Unless the marginal distributions of the rates are Gaussian, of course.

¹² At the time of writing, the Gaussian copula correlations implied from at-the-money spread options were largely independent of the expiry, an observation that is inconsistent with predictions of almost all multi-factor term structure models. Some observers attribute this market feature to the proliferation of copulas with time-independent parameters.

each rate described by a stochastic volatility model gives us an opportunity to define co-dependence between these rates by techniques other than the copula method. Broadly, if each swap rate involved in the payoff of a given multi-rate derivative has its own asset process and its own stochastic variance process (such as in (16.75)–(16.76)), then the co-dependence structure between rates can be controlled by correlating the Brownian motions that drive the asset and stochastic variance processes.

A stochastic volatility model for a given swap rate is often formulated in the annuity measure specific to that rate, in which case a translation into a common measure will be required in the multi-rate setup. Leaning on the general discussion in Section 17.2, we choose as common measure the forward measure associated with the payment date T_p . Conveniently, the problem of translation of dynamics from annuity to forward measures has already been considered in Chapter 16 where two different approaches were suggested. We consider both in turn.

17.8.1 Measure Change by Drift Adjustment

In Section 16.6.11 we derived a change of drifts associated with a shift of measure for SDEs driving a stochastic volatility model. Applying Proposition 16.6.8 to each swap rate, we obtain the following dynamics for a collection of d swap rates $(S_1(\cdot), \dots, S_d(\cdot))$ in the T_p -forward measure,

$$dS_i(t) = \lambda_i \varphi_i(S_i(t)) \sqrt{z_i(t)} \left(dW_i^{T_p}(t) + v^{S_i}(t) dt \right), \quad (17.58)$$

$$dz_i(t) = \theta(1 - z_i(t)) dt + \eta_i \psi_i(z_i(t)) \left(dZ_i^{T_p}(t) + v^{z_i}(t) dt \right), \quad z_i(0) = 1, \quad (17.59)$$

where $i = 1, \dots, d$ and the different parameters are explained in Section 16.6.11. Individual swap rate parameters $(\lambda_i, \varphi_i(\cdot), \eta_i)$ are obtained by the standard European swaption calibration for each swap rate (Section 16.1.4), and the drifts $v^{z_i}(t)$, $v^{S_i}(t)$ follow by the measure-change arguments of Proposition 16.6.8 (see also Corollary 16.6.9 for an important special case). The dependence structure between the swap rates may then be defined by correlations

$$\langle dW_i^{T_p}(t), dW_j^{T_p}(t) \rangle = \rho_{i,j} dt, \quad \langle dZ_i^{T_p}(t), dZ_j^{T_p}(t) \rangle = R_{i,j} dt, \quad i, j = 1, \dots, d,$$

where we keep asset/volatility correlations for each swap rate at zero, i.e.

$$\langle dW_i^{T_p}(t), dZ_i^{T_p}(t) \rangle = 0, \quad i = 1, \dots, d,$$

to replicate the setup of Proposition 16.6.8 exactly. Alternatively, we can calibrate these correlations together with other marginal parameters (which

would require a small extension of Proposition 16.6.8). Valuation of multi-rate derivatives in the model (17.58) requires Monte Carlo simulation, as the dimensionality of the model is high and the complexity of drifts does not lend itself easily to closed-form approximations.

17.8.2 Measure Change by CMS Caplet Calibration

The need to perform Monte Carlo simulation in the model (17.58) does not automatically render the approach unsuited for practical purposes — the Monte Carlo simulation of a d -asset stochastic volatility model is fairly quick, especially for important special cases of $d = 2, 3$. Nevertheless, it is a drawback. One way to develop more efficient schemes relies on the approach in Section 16.6.10, where we suggested translating stochastic volatility parameters from the annuity measure to the forward measure by first pricing CMS caplets in the model defined in the annuity measure, and then calibrating a new stochastic volatility model in the T_p -measure to these prices. See in particular (16.77) and the surrounding discussion. We remind the reader that this approach to a measure change is largely ad-hoc.

Let us assume that parameters have been suitably adjusted to incorporate the measure translation for each swap rate $S_i(t)$, $i = 1, \dots, d$. We therefore have available a CMS-adjusted forward rate $\tilde{S}_i = E^{T_p}(S_i(T))$, and a triple of T_p -measure parameters (λ_i, b_i, η_i) for each rate, where we dropped tildes from the notation of (16.77) for improved readability. A d -rate model can then be formulated by correlating all the driving Brownian motions,

$$\begin{aligned} dS_i(t) &= \lambda_i \left(b_i S_i(t) + (1 - b_i) \tilde{S}_i \right) \sqrt{z_i(t)} dW_i^{T_p}(t), \quad S_i(0) = \tilde{S}_i, \quad (17.60) \\ dz_i(t) &= \theta (1 - z_i(t)) dt + \eta_i \sqrt{z_i(t)} dZ_i^{T_p}(t), \quad z_i(0) = 1, \end{aligned}$$

$i = 1, \dots, d$, with the correlations defined by a $2d \times 2d$ correlation matrix R in the block form

$$\text{Corr} \left(\begin{pmatrix} dW^{T_p}(t) \\ dZ^{T_p}(t) \end{pmatrix}, \begin{pmatrix} dW^{T_p}(t) \\ dZ^{T_p}(t) \end{pmatrix} \right) = R, \quad R \triangleq \begin{pmatrix} R^{WW} & R^{WZ} \\ (R^{WZ})^\top & R^{ZZ} \end{pmatrix}, \quad (17.61)$$

where $W^{T_p}(t) = (W_1^{T_p}(t), \dots, W_d^{T_p}(t))^\top$, $Z^{T_p}(t) = (Z_1^{T_p}(t), \dots, Z_d^{T_p}(t))^\top$, and the matrices R^{WW} , R^{WZ} , R^{ZZ} are $d \times d$. We emphasize that the parameters (λ_i, b_i, η_i) , $i = 1, \dots, d$, are not obtained by a standard European swaption calibration for each swap rate, but by the more complicated two-step calibration described in Section 16.6.10.

With the joint dynamics of all swap rates under the same measure, the model (17.60) presents a straightforward extension of a one-factor displaced Heston model to d dimensions. For standard payoffs such as spread options or, more generally, options on the weighted average of d rates, (17.60) is simple enough for us to derive efficient closed-form approximations by

Markovian projection methods, an exercise that we postpone to Appendix A. For more complicated derivatives we can instead resort to Monte Carlo simulation. Each individual swap rate can be efficiently simulated using the methods from Section 9.5, for instance the Quadratic-Exponential (QE) discretization scheme for stochastic variance of Section 9.5.3.3, and the simplified Broadie-Kaya algorithm of Section 9.5.5.2 for the swap rate. To correlate different swap rates, as well as swap rate variances (and swap rates to variances of other swap rates) we just draw correlated Gaussian random variables to drive the discretization schemes¹³. For non-linear schemes such as QE, we note that this approach is not exact as the correlation between increments of swap rates and variances will not be exactly equal to the correlations of driving Gaussian variables. Nevertheless, numerical tests on the QE scheme show that this seemingly naive approximation is often of very good quality even for relatively coarse time discretizations. It is worth pointing out, however, that construction of accurate discretization schemes for multi-dimensional Heston-style stochastic volatility SDEs is an area of ongoing research.

17.8.3 Impact of Correlations on the Spread Smile

Apart from $\langle dW_i^{T''}(t), dZ_i^{T''}(t) \rangle$, $i = 1, \dots, d$, correlation parameters (17.61) in the model (17.60) do not affect the marginal distributions of each rate. They are, however, expected to affect the joint distribution of the rates; so, in a way, they define a “copula” function for the swap rates. To gain some intuition for the impact of various correlation parameters on the joint distribution, let us for illustrative purposes consider a model suitable for a CMS spread option with payoff

$$(S_1(T) - S_2(T) - K)^+, \text{ paid at } T_p, \quad (17.62)$$

i.e., a version of (17.60) with $d = 2$. The correlation matrix then has the form,

$$R = \begin{pmatrix} 1 & R_{12}^{WW} & R_{11}^{WZ} & R_{12}^{WZ} \\ R_{12}^{WW} & 1 & R_{21}^{WZ} & R_{22}^{WZ} \\ R_{11}^{WZ} & R_{21}^{WZ} & 1 & R_{12}^{ZZ} \\ R_{12}^{WZ} & R_{22}^{WZ} & R_{12}^{ZZ} & 1 \end{pmatrix},$$

where only the parameters R_{12}^{WW} , R_{12}^{WZ} , R_{21}^{WZ} , R_{12}^{ZZ} are “free” in the sense of not affecting the marginal distributions of the two swap rates.

Of the various entries in R , the “spot-spot” correlation R_{12}^{WW} is the most obvious in its effect on option value. Specifically, as the spread option value depends strongly on the variance of the swap rate difference, increasing

¹³If the QE scheme is used, whenever a uniform random variable U is required, we would write $U = \Phi(Z)$ with Z being a standard Gaussian random variable. This way, the QE scheme can be driven solely by Gaussian random variables.

(decreasing) the correlation between the Brownian motions driving the rates will decrease (increase) the option value. In the model (17.60), R_{12}^{WW} is typically the primary determinant of the spread option value or, equivalently, the overall level of the spread volatility smile (see Section 17.4.1). The effect of other correlation entries in R is more subtle and hard to grasp without resorting to numerical experiments. To conduct such an experiment, let us look at a CMS spread option with expiry $T=5$ years, with the model parameters in Table 17.2.

	Rate 1	Rate 2
CMS-adjusted swap rate \tilde{S}_i	4.97%	4.60%
Volatility λ_i	11.8%	13.2%
Skew b_i	100%	70%
Mean reversion of variance θ	10%	10%
Volatility of variance η_i	120%	120%

Table 17.2. Model Parameters for Heston model for CMS Spread

The base case for the correlation matrix is given in (17.63).

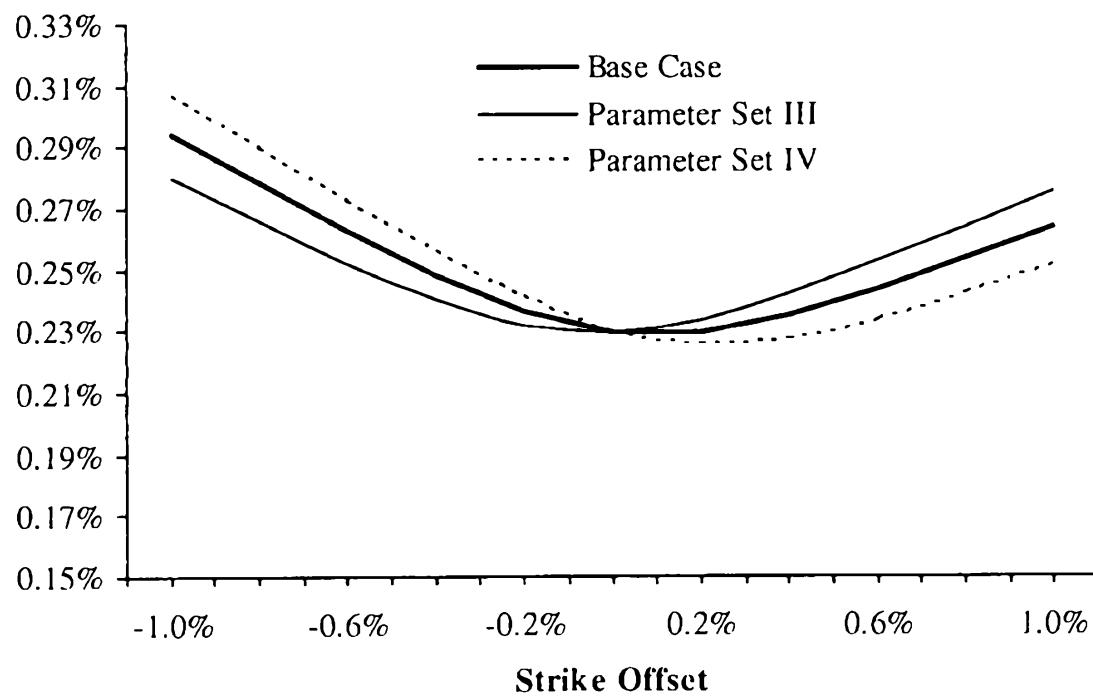
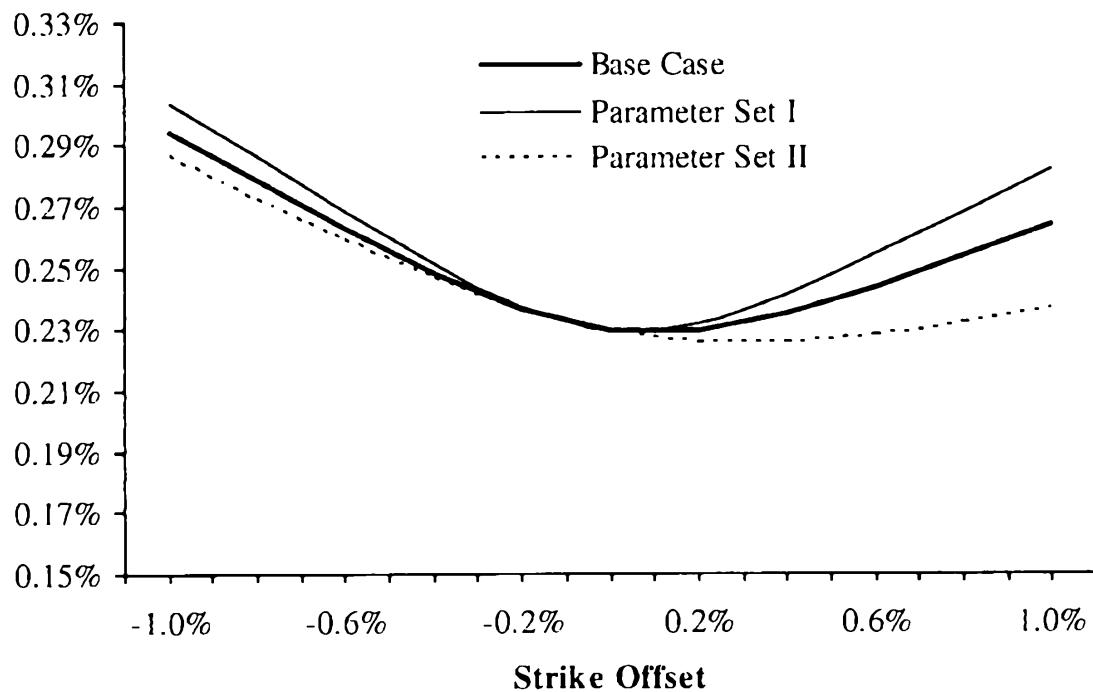
$$R = \begin{pmatrix} 100\% & 95\% & -25\% & -25\% \\ 95\% & 100\% & -20\% & -25\% \\ -25\% & -20\% & 100\% & 95\% \\ -25\% & -25\% & 95\% & 100\% \end{pmatrix}. \quad (17.63)$$

As is evident from Figure 17.3, the “vol-vol” correlation R_{12}^{ZZ} will move the overall level of the volatility smile up and down; after correcting for this level effect with the “spot-spot” correlation R_{12}^{WW} , increasing “vol-vol” correlation will allow one to add curvature to the spread volatility smile. Also, in Figure 17.3 we see that “spot1-vol2” correlation R_{12}^{WZ} affects the slope of the spread volatility smile. Interestingly, for this particular set of parameters, the “spot2-vol1” correlation R_{21}^{WZ} has very little impact on the spread volatility smile (as a consequence, the effect is not shown).

	Base Case	Set I	Set II	Set III	Set IV
R_{12}^{WW}	95%	92%	97%	95%	95%
R_{12}^{WZ}	-25%	-25%	-25%	-30%	-20%
R_{21}^{WZ}	-25%	-25%	-25%	-25%	-25%
R_{12}^{ZZ}	95%	97%	93%	95%	95%

Table 17.3. Correlation Parameter Sets for Figure 17.3

Fig. 17.3. Implied Normal Spread Volatility



Notes: Implied Normal spread volatility for the spread option (17.62) in the model described in Section 17.8.3, in particular in Table 17.2 and equation (17.63). The correlation scenarios shown in the graphs are listed in Table 17.3. “Strike offset” is the difference between the strike and the expected value of the CMS spread.

17.8.4 Connection to Term Structure Models

As we discussed previously in Section 17.7, making the copula method of Section 17.3 consistent with a typical dynamic term structure model is a

difficult objective that may not be possible to achieve effectively. In this respect, multi-rate vanilla models of the type (17.60) have a clear edge, as stochastic volatility is our preferred method of adding volatility smile capabilities to full term structure models. In particular, the multi-stochastic volatility LM model of Section 15.7, when applied to spread options, takes the form (17.60) (see Proposition 15.7.1), and hence the methods of Section 17.8 could be applied.

As multi-stochastic volatility models have yet to enter the mainstream, a more pressing task would be the construction of efficient methods for spread option pricing in simpler term structure models, e.g. ordinary (mono) stochastic volatility LM models. A crude approach for this was introduced in Section 14.4.3.2; the next section discusses several refinements.

17.9 CMS Spread Options in Term Structure Models

There are several reasons why we would want to efficiently calculate values of multi-rate derivatives — and in particular CMS spread options — in a term structure model. For example, for pricing an exotic derivative on underlyings that involve multi-rate payoffs we may wish to check how closely the term structure model values the underlying compared to the market; any observed differences can be used to correct the price of the exotic produced by the term structure model (see Chapter 21 for more on this topic). Another need arises in calibration, as we sometimes include CMS spread options in the calibration set as a source of market-implied correlations between swap rates (see Section 14.5.9). The efficiency requirements imposed by both of these applications typically rule out Monte Carlo methods, so here we seek to develop closed-form approximations for a few specific term structure models.

17.9.1 Libor Market Model

We first consider Libor market (LM) models, as these are often used as workhorses for pricing exotics on multi-rate underlyings. For simplicity, let us focus on a CMS spread option with the payoff (17.62), although similar techniques can be applied to other payoffs.

In Section 14.4.3.2 we presented a simple method for pricing CMS spread options in the LM model, based on a Gaussian approximation to the spread process (and little concern for the effects of measure changes). Here we develop a more sophisticated, and more accurate, approach that utilizes copula techniques. In fact, there is not much left to do at this point, as we have already developed almost all the “ingredients” we need for the method. In Section 16.6.6 we derived the annuity mapping function consistent with the LM model; using the machinery of Section 16.6.9, the mapping function could be turned into a CDF of each swap rate in the T_p -forward measure. Thus, according to Section 17.4, all that remains to do

is to find a copula that is consistent with the dependence structure of the swap rates in the LM model. Considering, for concreteness, the stochastic volatility LM model specification (14.15)–(14.16) (also see (16.59)), we recall the results of Section 14.4.3.1 and, in particular, (14.34) that tell us that the instantaneous correlation between two swap rates in the stochastic volatility LM model is, in fact, deterministic¹⁴ and is the same as in a displaced log-normal LM model. In a displaced log-normal LM model, swap rates are approximately functions of Gaussian variables (see Appendix 17.A), so we can approximate the dependence structure of two swap rates in the SV-LM model by a Gaussian copula. The correlation to be used in the copula may then be approximated as the term correlation of the swap rates $\rho_{\text{term}}(0, T)$, as given by (14.35). We formalize this discussion as a proposition.

Proposition 17.9.1. *The undiscounted value $V(0)$ of a CMS spread option with the payoff (17.62) in the stochastic volatility LM model (14.15)–(14.16) is approximately given by the two-dimensional Gaussian copula integral*

$$V(0) \approx \int \int \left(\left(\Psi_1^{T_p} \right)^{-1}(\Phi(z_1)) - \left(\Psi_2^{T_p} \right)^{-1}(\Phi(z_2)) - K \right)^+ \times \phi(z_1, z_2; \rho_{\text{term}}(0, T)) dz_1 dz_2,$$

where $\phi(z_1, z_2; R)$ is a two-dimensional Gaussian density with correlation R , $\rho_{\text{term}}(0, T)$ is given by (14.35), and $\Phi(z)$ is the standard Gaussian CDF.

Remark 17.9.2. In Proposition 17.9.1, $\Psi_i^{T_p}(s)$ are the T_p -forward measure CDFs of the swap rates $S_i(T)$, $i = 1, 2$. From Proposition 16.6.4 these CDFs can be obtained from the annuity-measure swap rate CDFs $\Psi_i^{A_i}(s)$ and the annuity mapping functions $\alpha_i(s)$, $i = 1, 2$. The annuity mapping function $\alpha_i(s)$ may be computed as in Proposition 16.6.3 and the CDFs $\Psi_i^{A_i}(s)$, $i = 1, 2$, can be obtained by differentiating (in strike) European swaption values in the LM model, per formulas in Proposition 14.4.3 or Section 15.2.

The result in Proposition 17.9.1 trivially generalizes to general multi-rate cash flows. Of course, for spread options specifically, we can make the algorithm more efficient by utilizing one-dimensional integration formulas from Section 17.6.2. If the speed of this approach is too slow, we can trade accuracy for speed by utilizing some of the ideas of Section 17.8.2. Recall that in the SV-LM model, the distribution of each $S_i(T)$ in its corresponding annuity measure is given by the SV model. If we can approximate the marginal distributions of $S_1(T)$, $S_2(T)$ in Q^{T_p} by the SV model as well, then the distribution of $S_1(T) - S_2(T)$ could be quite effectively approximated by the same distribution, see Appendix A, leading to a closed-form approximation to the spread option value. Section 16.6.10 outlines the numerical approach

¹⁴This is a consequence of using just a single stochastic variance scaling applied to all Libor rates.

to the measure change that preserves the SV distribution class; in a pinch, we can just reuse the SV parameters from the annuity measure distributions while adjusting the forward swap rates with the CMS convexity adjustments,

$$S_i(0) \rightarrow E^{T_p}(S_i(T)), \quad i = 1, 2. \quad (17.64)$$

While Section 16.6.10 warns against performing the measure shift in the SV model by the sole change of the forward (17.64), the impact of such laissez-faire approach on the quality of CMS spread option valuation in LM models is likely to be muted, given all the other approximations we are making on the way.

Finally, let us note that Antonov and Arneguy [2009] present an alternative approximation idea that is based on working in a measure in which the spread $S_1(t) - S_2(t)$ is a martingale. While such a measure cannot be easily characterized by a numeraire, it can still be defined by the drift change in the Brownian motions driving the model's SDEs. We refer the interested reader to the source paper, as the approach is too involved to be described here in a few lines.

17.9.2 Quadratic Gaussian Model

Having dealt with the LM models, we now turn our attention to the multi-factor quadratic Gaussian (QG) models of Section 12.3. Interestingly, the most productive approach here is quite different from that used for the LM models — reflecting the different approaches for European swaption pricing in the two models, see Section 12.3.4. For the purpose of elaborating on this observation, we continue examining the CMS spread option (17.62).

As a start, we recall that in a QG model, a swap rate $S(T)$ is a deterministic function of the state vector $z(T)$, $S(T) = S(T, z(T))$. In one of the approximations to the swap rate we developed in Section 12.3.4, we replaced the function with a quadratic form of the state vector, see (12.92). Let us denote the quadratic approximations to the two rates involved in the payoff (17.62) by $S_{1,q}(T, z)$ and $S_{2,q}(T, z)$,

$$S_{i,q}(T, z) = z^\top \gamma_{S_i} z + h_{S_i}^\top z - E^A (z(T)^\top \gamma_{S_i} z(T) + h_{S_i}^\top z(T)) + S_i(0), \quad i = 1, 2.$$

Then, the undiscounted value of the spread option is given approximately by

$$V(0) \approx E^{T_p} \left((S_{1,q}(T, z(T)) - S_{2,q}(T, z(T)) - K)^+ \right).$$

Two points should now be obvious. One is that the difference of two quadratic forms $S_{1,q}(T, z) - S_{2,q}(T, z)$ is itself a quadratic form in z ,

$$\begin{aligned} S_{1,q}(T, z) - S_{2,q}(T, z) &= z^\top (\gamma_{S_1} - \gamma_{S_2}) z + (h_{S_1} - h_{S_2})^\top z \\ &- E^A \left(z(T)^\top (\gamma_{S_1} - \gamma_{S_2}) z(T) + (h_{S_1} - h_{S_2})^\top z(T) \right) + S_1(0) - S_2(0). \end{aligned}$$

Another is that the distribution of $z(T)$ in the T_p -forward measure is known — it is Gaussian with a known mean $m^{T_p}(0, T, 0)$ and covariance matrix $\nu^{T_p}(0, T, 0)$, see Proposition 12.3.4. Hence, the problem of pricing a spread option in the QG model is almost identical to the problem of pricing a European swaption, and any of the methods of Section 12.3.4.3 could be applied (with the most efficient method probably being the two-dimensional integration method in Theorem 12.3.7). Further details are available in Piterbarg [2009b].

17.A Appendix: Implied Correlation in Displaced Log-Normal Models

17.A.1 Preliminaries

The purpose of this appendix is to briefly examine how marginal distributions affect spread option prices in a Gaussian copula. We shall consider only marginal distributions originating from displaced log-normal dynamics, so first let us recall that a (one-dimensional) process of the form

$$dX(t) = \lambda((1-b)X_0 + bX(t)) dW(t), \quad X(0) = X_0,$$

with $b > 0$ has the solution

$$X(T) = \frac{X_0}{b} \left(\exp(-b^2 \lambda^2 T/2 + b\lambda W(T)) - 1 + b \right).$$

Given this result, let us consider two swap rates $S_1(t)$ and $S_2(t)$ in the T -forward measure, with

$$\mathbb{E}^T(S_1(T)) = \mathbb{E}^T(S_2(T)) = S_0 > 0,$$

and set, for $b > 0$,

$$S_i(T) = \frac{S_0}{b} \left(\exp(-b^2 \lambda_i^2 T/2 + b\lambda_i Z_i \sqrt{T}) - 1 + b \right), \quad b > 0, \quad i = 1, 2, \tag{17.65}$$

where Z_1 and Z_2 are two standard Gaussian random variables with constant correlation ρ . As $S_1(T)$ and $S_2(T)$ are monotonic functions of correlated Gaussian variables, it is clear¹⁵ that their dependency is generated by a two-dimensional Gaussian copula with correlation parameter ρ .

From the definition (17.65) it follows that

$$\text{Var}(S_i(T)) = \frac{S_0^2}{b^2} \left(e^{b^2 \lambda_i^2 T} - 1 \right), \quad i = 1, 2,$$

¹⁵A copula is easily shown to be invariant with respect to monotonic transformations of the underlying variables.

and

$$\text{Cov}(S_1(T), S_2(T)) = \frac{S_0^2}{b^2} \left(e^{\rho b^2 \lambda_1 \lambda_2 T} - 1 \right). \quad (17.66)$$

Therefore

$$\text{Corr}(S_1(T), S_2(T)) = \frac{e^{\rho b^2 \lambda_1 \lambda_2 T} - 1}{\sqrt{e^{b^2 \lambda_1^2 T} - 1} \sqrt{e^{b^2 \lambda_2^2 T} - 1}}. \quad (17.67)$$

We notice that $\text{Corr}(S_1(T), S_2(T))$ is, of course, not equal to ρ , but instead is given by a more complicated expression that in most practically relevant cases will decrease in b , for fixed λ_1 and λ_2 . Based solely on this observation, we might therefore expect that spread option prices would increase in the skew parameter b .

17.A.2 Implied Log-Normal Correlation

The correlation in (17.67) is not a particularly market-oriented way of characterizing spread option value. As we have seen earlier in this chapter, a better measure may be to use implied spread volatility. To offer a slightly different perspective, in this appendix we instead work with an implied log-normal correlation ρ_{LN} , defined as the value of the copula correlation ρ that will match the ATM spread option¹⁶ in the true model, after b has been set to 1 (but with ATM option prices kept constant). Examining how ρ_{LN} depends on the skew b will give us a convenient scalar measure of how skew affects co-dependence in a spread option setting.

To compute ρ_{LN} , first consider the zero-strike payout

$$V(0) = \mathbb{E}^T((S_1(T) - S_2(T))^+).$$

Writing $V(0) = \mathbb{E}^T(S_1(T)(1 - S_2(T)/S_1(T))^+)$ shows (after a measure change, see Section 17.6.5) that

$$V(0) = \frac{S_0}{b} (2\Phi(d) - 1), \quad d = \frac{\sigma b}{2} \sqrt{T}, \quad (17.68)$$

where $\sigma = \sqrt{\lambda_1^2 + \lambda_2^2 - 2\rho\lambda_1\lambda_2}$. We note in passing that when $b = 1$ this formula is known as the *Margrabe formula*, see Margrabe [1978]. We also note that, for $i = 1, 2$,

$$\mathbb{E}^T((S_i(T) - S_0)^+) = \frac{S_0}{b} (2\Phi(y_i) - 1), \quad y_i = \frac{\lambda_i b}{2} \sqrt{T}. \quad (17.69)$$

Suppose now that we observe implied Black ATM volatility for $S_1(T)$ and $S_2(T)$ to be σ_1 and σ_2 , respectively. Using the Black option pricing formula, for $i = 1, 2$ we therefore must have

¹⁶We can extend the definition to handle non-ATM strikes, but this shall not be needed for the purposes of this appendix.

$$\begin{aligned} \mathbb{E}^T \left((S_i(T) - S_0)^+ \right) &= S_0 \Phi \left(\frac{1}{2} \sigma_i \sqrt{T} \right) - S_0 \Phi \left(-\frac{1}{2} \sigma_i \sqrt{T} \right) \\ &= S_0 \left(2 \Phi \left(\frac{\sigma_i}{2} \sqrt{T} \right) - 1 \right). \end{aligned} \quad (17.70)$$

Suppose also that we have best-fit b to some value different from 1. To preserve ATM option prices, we equate (17.69) and (17.70),

$$\frac{1}{b} \left(2 \Phi \left(\frac{\lambda_i b}{2} \sqrt{T} \right) - 1 \right) = 2 \Phi \left(\frac{\sigma_i}{2} \sqrt{T} \right) - 1, \quad i = 1, 2.$$

We can solve these equations in closed form for λ_i to yield

$$\lambda_i = \frac{2}{b \sqrt{T}} \Phi^{-1} \left(\left(\Phi \left(\frac{\sigma_i}{2} \sqrt{T} \right) - \frac{1}{2} \right) b + \frac{1}{2} \right), \quad i = 1, 2. \quad (17.71)$$

In most cases, this equation results in $\lambda_1 \approx \sigma_1$ and $\lambda_2 \approx \sigma_2$. Assuming that the copula correlation in our model has been set to ρ , by our definition of ρ_{LN} , we get from (17.68) that

$$\frac{1}{b} (2 \Phi(d) - 1) = 2 \Phi \left(\frac{1}{2} \sqrt{\sigma_1^2 + \sigma_2^2 - 2 \rho_{LN} \sigma_1 \sigma_2} \sqrt{T} \right) - 1,$$

or

$$\sqrt{\sigma_1^2 + \sigma_2^2 - 2 \rho_{LN} \sigma_1 \sigma_2} = \frac{2}{\sqrt{T}} \Phi^{-1} \left(\frac{1}{2b} (2 \Phi(d) - 1) + \frac{1}{2} \right), \quad (17.72)$$

from which we can extract an analytical expression for $\rho_{LN} = \rho_{LN}(b, \sigma_1, \sigma_2, \rho, T)$. We omit further details in the interest of brevity, but just notice the curious fact that $\rho_{LN}(b, \sigma_1, \sigma_1, 0.5, T) = 0.5$, independently of b, σ_1, T .

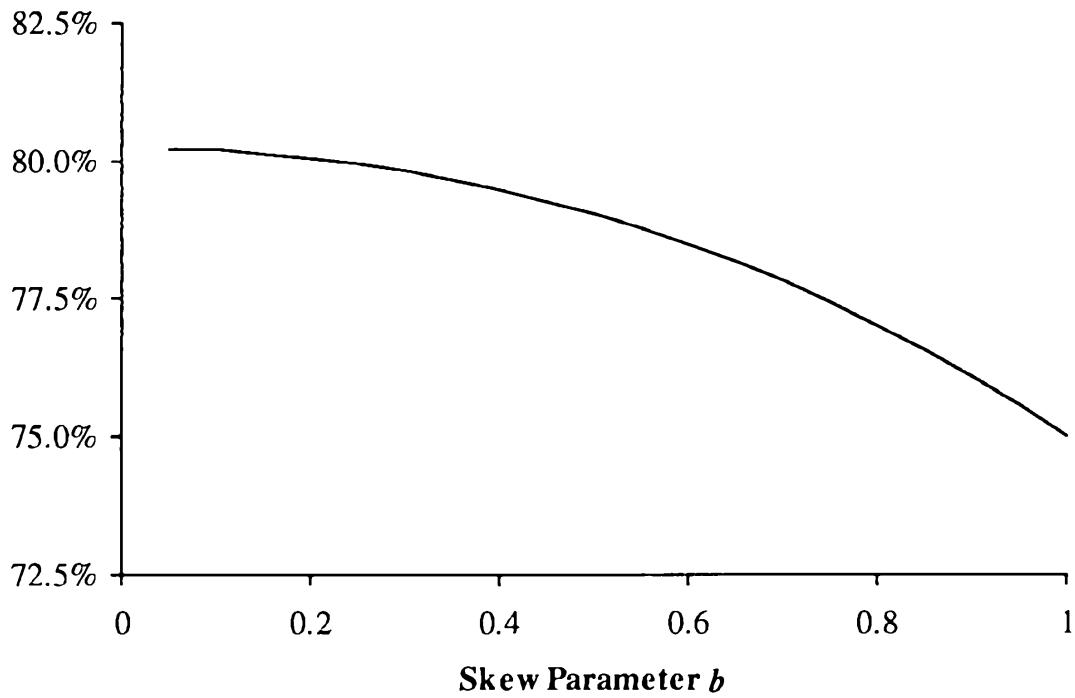
17.A.3 A Few Numerical Results

Going forward, we assume that log-normal volatilities are $\sigma_1 = \sigma_2 = 25\%$. First, we set $T = 10$ years and $\rho = 0.75$. Figure 17.4 shows ρ_{LN} as a function of b . Notice the implied correlation ρ_{LN} here decreases in b ; in other words, tilting the skew downwards (i.e. lowering b) has the effect of increasing the effective spread option correlation.

Next, we freeze $b = 0.25$ and $\rho = 0.75$ and let T vary. Figure 17.5 shows the resulting effects on ρ_{LN} . The effect of skew on implied correlation ρ_{LN} is clearly quite sensitive to maturity T , with the difference $|\rho - \rho_{LN}|$ increasing in T .

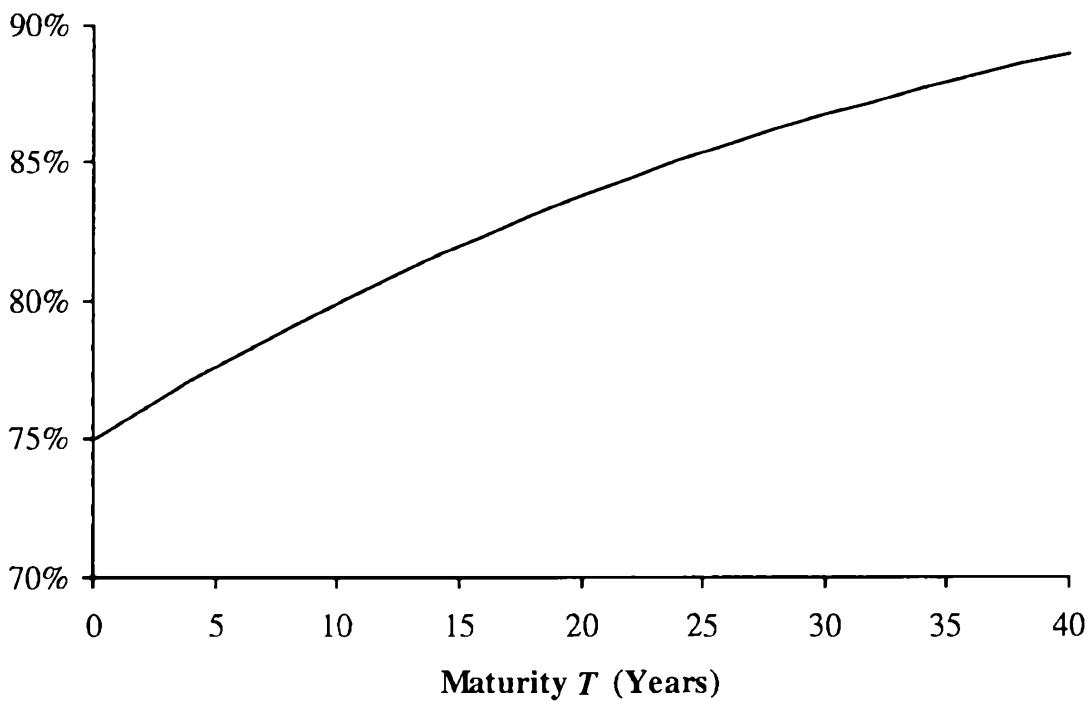
While we do not show the results, we note that lowering volatility (σ_1 and σ_2) will have qualitatively the same effect as lowering maturity. In low-volatility regimes (and for short-dated options), the practice of assuming that the market-implied ATM spread option correlation is independent of skew may therefore be defensible.

Fig. 17.4. Log-Normal Correlation ρ_{LN}



Notes: Model parameters are: $T = 10$, $\rho = 0.75$, and $\sigma_1 = \sigma_2 = 0.25$.

Fig. 17.5. Log-Normal Correlation ρ_{LN}



Notes: Model parameters are: $b = 0.25$, $\rho = 0.75$, and $\sigma_1 = \sigma_2 = 0.25$.

Callable Libor Exotics

Having discussed relatively simple vanilla securities in Chapters 16 and 17, we now move to the other extreme of the product spectrum and consider issues of valuation, calibration and risk management of *callable Libor exotics* (CLEs). CLEs were defined in Section 5.14 and constitute the most complicated class of interest rate derivatives traded in the market. The material in this chapter deals with CLEs in general; later chapters will take a more in-depth look at the idiosyncrasies of some of the most popular and/or challenging securities inside the CLE class.

18.1 Model Calibration for Callable Libor Exotics

Due to their inherent complexity, CLEs will have non-trivial dependencies on the dynamics of market rates and will require sophisticated term structure models for valuation and risk management. As discussed in earlier chapters, computation of non-vanilla prices in such models virtually always requires Monte Carlo simulation, which in turn introduces complications in determining the optimal exercise rules for CLEs. We shall spend a good portion of this chapter on a detailed discussion of algorithms for this particular problem. However, before a model can even be considered for CLE valuation, it needs to be calibrated, i.e. the *volatility structure* of the model needs to be parameterized to match available market information relevant for valuation.

At this point we should emphasize that CLEs and other exotic interest rate derivatives are different from many of vanilla securities considered previously in that their values are not observable in the market, a consequence of the fact that most exotic derivatives are sold to clients rather than traded between dealers. This fact leads to fundamental differences in the way models are used for valuation of vanilla vs. exotic derivatives. Whereas for many vanilla derivatives a model is primarily used to generate a hedging strategy and to perform inter- and extrapolation on market-observable prices, for

exotics the values are fundamentally derived from a model¹. Given the absence of market prices of exotics themselves, models for these securities have to be calibrated indirectly, to other market (and perhaps also non-market) information that is deemed relevant for the class of exotics under consideration. Speaking very loosely, the purpose of a model for CLEs can therefore be characterized as performing a sophisticated extrapolation of information from a series of “spanning” vanilla markets to compute a meaningful exotics price, something we already alluded to in Section 14.5.5.

18.1.1 Risk Factors for CLEs

In order to make sure that a CLE model captures as much relevant market information as possible, it is important to understand what we can actually rely on when calibrating a model. There are essentially three sources of potentially relevant information. The first, and arguably the most important, source is the market prices of liquid vanilla interest rate derivatives. The second source is historical information about market quantities such as volatilities and correlations of various market rates that may not be available for observation directly from quoted market prices. The third, and somewhat more amorphous, source is the modeler’s beliefs on what constitutes reasonable behavior of the model parameters. The last category, for example, includes views on how time-stationary or how smooth model parameters should be.

To perform model calibration, we typically choose particular targets from all three sources above, and ultimately set model parameters in a way to match those targets in our model. In order to be able to choose relevant targets, it is important to identify risk factors that affect the valuation of the particular CLE in question. This part of the analysis is necessarily product specific, so let us look at a particular example to demonstrate some salient points. We choose a reasonably representative (albeit not very complicated) example, namely a callable inverse floater, or CIF, with exercise dates

$$T_1 < T_2 < \dots < T_{N-1}$$

and structured coupon

$$C_n = \min(\max(6\% - L_n(T_n), 0\%), 4\%),$$

where $L_n(t)$ is a forward Libor rate fixing at T_n . In the language of Section 5.14, we see that the coupon strike is equal to 6%, the cap is at 4% and the floor is at 0%. It should be clear that the coupon can be decomposed into a portfolio of a long floorlet with strike 6% and a short floorlet with strike 2%. Hence, the callable inverse floater can be thought of as a Bermudan-style option on a combination of floors and a Libor leg.

¹A practice sometimes known as *mark-to-model*.

The underlying swap of this CIF consists of a collection of floorlets, which are vanilla derivatives with observable market prices. Market-implied volatilities² as always serve as a convenient representation of these market values. For classification purposes that will be clear in a moment, we call such volatilities *spot volatilities*, as we can observe them on the valuation date, or “on the spot”, from the market. Clearly, these volatilities affect the value of the exotic swap that underlies the CIF and hence have a direct impact on the value of the CIF. For that reason we should include them as targets for model calibration. Importantly, as caplet volatilities virtually always exhibit a volatility smile, volatilities for both 2% and 6% strikes are potentially relevant.

While the underlying exotic swap has no dependencies on spot volatilities of market rates other than those underlying the floorlets inherent in the CIF, the callability structure of the CIF nevertheless introduces additional dependence on other vanilla derivatives. To see this, suppose that we, say, ignored all exercise CIF dates but one, in which case the CIF would degenerate into a European-style option to enter the underlying swap. Even though the underlying swap is not vanilla (i.e. not a standard fixed-for-floating rate swap), it is clearly related to one and our single-call CIF is therefore related to the European swaption on that swap. More generally, the (multi-call) CIF will depend on the (spot) implied volatility of the swap rate that fixes on T_i , and runs for the period $[T_i, T_N]$, for each $i = 1, \dots, N - 1$. It is less clear what strike we should use to define this implied volatility; while all strikes are in fact relevant, as shall be clear later, sometimes one may choose to simplify and just use the ATM swaption volatility.

To summarize our discussion so far, we have argued that the value of the CIF depends on market-implied spot Libor rate volatilities for all expiries until T_{N-1} (of two different strikes), and spot swap rate volatilities for those swaptions for which expiry plus tenor is equal to T_N — the so called *core*, *diagonal*, or *coterminal* swaptions. When calibrating a model to price this CIF, we would seek (as a minimum) to calibrate it to these market volatilities.

While there are no more obvious spot volatility targets that should be included in CIF calibration, there are other volatility-related quantities that will affect the valuation of the CIF and other similar CLEs. Let us imagine that we have arrived at exercise date T_n , $n < N - 1$, and assume that the CIF has not been exercised in the past. At this point in time, the option holder needs to decide whether to exercise and receive the underlying value $U_n(T_n)$ (see (5.24)) or whether to continue to hold on to the CIF (with a view of perhaps exercising it later). The time T_n value of the remaining part of the underlying swap depends on caplet volatilities as observed at time T_n . As these volatilities will be known only at time T_n , we call them *forward volatilities* of Libor rates. Likewise, the option to hold a Bermudan-style

²As expressed, say, Black volatilities.

option on the underlying $U_{n+1}(T_n)$, will depend on forward volatilities of swap rates, specifically core swaption volatilities observed at time T_n . As the exercise decision at time T_n evidently will depend on the forward volatility structure at time T_n , the time 0 value of the CIF will depend on it as well.

At time 0, forward volatilities are generally unknown random variables, but any model will impose certain dynamics on the volatility structure of interest rates, which will have a direct impact on the model value of the CIF (and other CLEs). Ideally, we should make sure that the model projections for forward volatilities are in line with market-implied information. The fact that such information is typically either not available or difficult to extract adds significantly to the challenge of model calibration for CLEs. Frequently, it will be necessary to lean on historical data and to impose exogenous assumptions that the model builder might feel are financially reasonable.

In Section 13.1.8.1 we identified inter-temporal correlations — i.e. correlations of core swap rates observed on different fixing dates — as value drivers for Bermudan swaptions. Inter-temporal correlations are closely linked to the forward volatilities discussed above, and the two concepts can often be used interchangeably. To demonstrate, note that the (inter-temporal) correlation between two core swap rates, $S_{n,N-n}(T_n)$ and $S_{m,N-m}(T_m)$ ($n < m$) is given by

$$\begin{aligned} \text{Corr}(S_{n,N-n}(T_n), S_{m,N-m}(T_m)) &= \text{Cov}(S_{n,N-n}(T_n), S_{m,N-m}(T_m)) \\ &\quad \times \text{Var}(S_{n,N-n}(T_n))^{-1/2} \text{Var}(S_{m,N-m}(T_m))^{-1/2}. \end{aligned}$$

As $S_{m,N-m}(T_m) - S_{m,N-m}(T_n)$ is often only mildly dependent on $S_{m,N-m}(T_n)$ and $S_{n,N-n}(T_n)$, we can rewrite this as

$$\begin{aligned} \text{Corr}(S_{n,N-n}(T_n), S_{m,N-m}(T_m)) &\approx \text{Cov}(S_{n,N-n}(T_n), S_{m,N-m}(T_m)) \text{Var}(S_{n,N-n}(T_n))^{-1/2} \\ &\quad \times (\text{Var}(S_{m,N-m}(T_n)) + \text{Var}(S_{m,N-m}(T_m) - S_{m,N-m}(T_n)))^{-1/2} \\ &\approx \text{Corr}(S_{n,N-n}(T_n), S_{m,N-m}(T_n)) \\ &\quad \times \left(1 + \frac{\text{Var}(S_{m,N-m}(T_m) - S_{m,N-m}(T_n))}{\text{Var}(S_{m,N-m}(T_n))}\right)^{-1/2} \end{aligned}$$

and we see that, with (spot) correlation $\text{Corr}(S_{n,N-n}(T_n), S_{m,N-m}(T_n))$ and (spot) variance $\text{Var}(S_{m,N-m}(T_n))$ fixed, the inter-temporal correlation is directly linked to the forward variance (square of forward volatility) $\text{Var}(S_{m,N-m}(T_m) - S_{m,N-m}(T_n))$. This equivalence will play a role when we discuss the local projection method in Section 18.4.

Another aspect of the model behavior worth emphasizing here is the volatility smile dynamics imposed by the model. In Sections 8.8 and 16.1.1 we considered the hedging impact of joint moves in rates and volatility smiles. For CLEs, volatility smile dynamics affect not only the hedging strategy, but

also the valuation of CLEs, due to the forward volatility effect we discussed earlier. We notice that for many CLEs the effective coupon strikes at the exercise boundary often end up being deeply in- or out-of-the-money, since interest rate levels at which option exercise is optimal are usually significantly different from the levels of rates at inception of the trade.

18.1.2 Model Choice and Calibration

As argued above, the value of a typical CLE depends strongly on volatility smile dynamics and on market (spot) implied volatilities for a wide selection of vanilla options, often across a range of strikes. As such the appropriate choice of a model for CLE valuation should typically involve the following criteria.

- Ability to calibrate to a large collection of vanilla options across expiries, tenors, and strikes.
- Reasonable and controllable dynamics of the volatility structure.
- Multi-factor interest rate dynamics, especially for CLEs on multi-rate underlyings.

A combination of these requirements often rules out simpler, low-dimensional models, especially for more complicated CLEs. We typically recommend using either Libor market, multi-factor quasi-Gaussian, or multi-factor quadratic Gaussian models. While simpler low-dimensional models might succeed in fitting spot volatility information for selected Libor and swap rates, they often achieve this only by using values of model parameters that would imply unrealistic evolution of the volatility structure and, therefore, unreasonable pricing of CLEs (but see Section 18.4).

Another observation favors models that can calibrate to a large collection of European swaptions. A swap rate can be seen to be, approximately, a weighted average of Libor rates (see Section 14.4.2). Hence, an implied European swaption volatility contains some amount of information on market-consistent correlations between Libor (and other swap) rates. Extracting this information is complicated, but in a sense is what a Libor market model calibrated to, say, the whole swaption grid (including caps) is designed to do. In the same spirit, a model calibrated to all swaptions and caps can be thought of as giving us the best available *implied* forward volatility structure, i.e. market-consistent information of what the most likely behavior of the volatility structure through time would be.

Another dimension to consider here is the volatility smile dynamics. Obviously, for the model to generate unspanned stochasticity in volatility movements, stochastic volatility (SV) version of the model would be required. Given the discussion in Section 8.8 and the importance for CLE pricing of controlling the dynamics of the volatility smile, SV models typically have a clear edge over their local volatility counterparts in CLE pricing applications.

With the choice of the model settled, calibration typically proceeds by identifying proper targets and fitting them with model parameters. For Libor market models we, in fact, have already covered relevant issues quite extensively, and the reader is advised to revisit Section 14.5.5. For other types of multi-factor models, the issues are similar and can be resolved in the same spirit. The mechanics of calibration for relevant models are covered in their respective chapters. For example, smile calibration of LM models is discussed in Section 15.2.

18.2 Valuation Theory

18.2.1 Preliminaries

While various measures could be used for CLE valuation, for concreteness we choose to use the spot Libor measure Q^B with the discrete money market numeraire $B(t)$ defined on a tenor structure

$$0 = T_0 < T_1 < \dots < T_N, \quad \tau_n = T_{n+1} - T_n.$$

For notational simplicity, we let $E(\cdot)$ be the expected value operator in measure Q^B .

We recall from Section 5.14 that a callable Libor exotic is a Bermudan-style option on an exotic swap that specifies an exchange of structured coupons C_n for Libor rates L_n , fixing on T_n and paying on T_{n+1} , $n = 1, \dots, N-1$. We denote by X_n the net payment seen by the structured coupon receiver,

$$X_n = \tau_n \times (C_n - L_n(T_n)). \quad (18.1)$$

Also, we let $U_n(t)$ be the n -th exercise value, i.e. the value of all future payments if the callable Libor exotic is exercised at time T_n . Clearly

$$U_n(t) = B(t) \sum_{i=n}^{N-1} E_t (B(T_{i+1})^{-1} X_i). \quad (18.2)$$

For completeness we set

$$U_N(t) \equiv 0.$$

If a callable Libor exotic is exercised on T_n , the holder will receive $U_n(T_n)$, the present value of the remainder of the underlying exotic swap.

Finally, let $H_n(t)$ be the value at time t of a callable Libor exotic where exercise opportunities have been restricted to the dates $\{T_{n+1}, \dots, T_{N-1}\}$. In particular, $H_0(0)$ is then the time 0 value of the CLE. Each H_n is called a *hold value*, since $H_n(T_n)$ is the value of the choice of *not* exercising on date T_n , i.e. continuing to “hold” the derivative. We must necessarily have

$$H_0(t) \geq H_1(t) \geq \dots \geq H_{N-2}(t). \quad (18.3)$$

18.2.2 Recursion for Callable Libor Exotics

Let \mathcal{T}_n be a set of all stopping time indices that exceed n , $n \in \{0, \dots, N-1\}$, i.e. a set of random variables taking values in the set $\{n+1, \dots, N\}$ such that for any k and any $\xi \in \mathcal{T}_n$,

$$\{\xi = k\} \in \mathcal{F}_{T_k}.$$

The sequence of random variables $H_0(T_0), \dots, H_{N-1}(T_{N-1})$ (as in Section 18.2.1) defines the Snell envelope for the sequence of (discounted) exercise values,

$$H_n(T_n) = B(T_n) \sup_{\xi \in \mathcal{T}_n} \mathbb{E}_{T_n} (B(T_\xi)^{-1} U_\xi(T_\xi)), \quad n = 0, \dots, N-1. \quad (18.4)$$

By the general theory of optimal stopping (see Chapter 1.10), the random time index η_n that maximizes the right-hand side of (18.4) is given by

$$\eta_n(\omega) = \min \{k > n : U_k(T_k) \geq H_k(T_k)\} \wedge N, \quad (18.5)$$

and we set

$$\eta \triangleq \eta_0 = \min \{k \geq 1 : U_k(T_k) \geq H_k(T_k)\} \wedge N.$$

With this definition, the value of a callable contract can be re-written as

$$H_0(0) = \mathbb{E} (B(T_\eta)^{-1} U_\eta(T_\eta)) = \mathbb{E} \left(\sum_{n=\eta}^{N-1} B(T_{n+1})^{-1} X_n \right). \quad (18.6)$$

The Hamilton-Jacobi-Bellman equation that corresponds to the optimal stopping problem (18.4) can be solved by backward induction. In particular, we have, for $n = N-1, \dots, 1$,

$$H_{n-1}(T_{n-1}) = B(T_{n-1}) \mathbb{E}_{T_{n-1}} (B(T_n)^{-1} \max(H_n(T_n), U_n(T_n))), \quad (18.7)$$

subject to the terminal condition $H_{N-1} \equiv 0$. The recursion starts at the final time $n = N-1$ and progresses backward in time until we obtain the desired time 0 security value $H_0(0)$.

The financial meaning of the recursion above is straightforward. If a callable contract H_0 has not been exercised up to and including time T_n , then it is worth the hold value $H_n(T_n)$. If the callable contract is exercised at time T_n its value is equal to $U_n(T_n)$. Assuming optimal exercise, the value of the callable Libor exotic H_0 at time T_n is the maximum of the two,

$$\max(H_n(T_n), U_n(T_n)).$$

The value of this payoff at time T_{n-1} is then

$$B(T_{n-1}) \mathbb{E}_{T_{n-1}} (B(T_n)^{-1} \max(H_n(T_n), U_n(T_n))).$$

But clearly, this is the value of the CLE that can only be exercised at times T_n and beyond, i.e. of the CLE H_{n-1} , as specified in (18.7).

18.2.3 Marginal Exercise Value Decomposition

Before discussing techniques to numerically implement the valuation equations of the previous section, let us briefly review an important decomposition result for CLEs that follows from the recursion (18.7). After a slight rewrite, we obtain

$$\begin{aligned} H_{n-1}(T_{n-1}) - B(T_{n-1})\mathbb{E}_{T_{n-1}}(B(T_n)^{-1}H_n(T_n)) \\ = B(T_{n-1})\mathbb{E}_{T_{n-1}}\left(B(T_n)^{-1}(U_n(T_n) - H_n(T_n))^+\right). \end{aligned}$$

Note that

$$B(T_{n-1})\mathbb{E}_{T_{n-1}}(B(T_n)^{-1}H_n(T_n)) = H_n(T_{n-1}),$$

so that

$$H_{n-1}(T_{n-1}) - H_n(T_{n-1}) = B(T_{n-1})\mathbb{E}_{T_{n-1}}\left(B(T_n)^{-1}(U_n(T_n) - H_n(T_n))^+\right).$$

Taking discounted expectations to time 0 we obtain

$$H_{n-1}(0) - H_n(0) = \mathbb{E}\left(B(T_n)^{-1}(U_n(T_n) - H_n(T_n))^+\right)$$

and, summing up from $n = 1$ to $N - 1$,

$$H_0(0) - H_{N-1}(0) = \sum_{n=1}^{N-1} \mathbb{E}\left(B(T_n)^{-1}(U_n(T_n) - H_n(T_n))^+\right).$$

Since $H_{N-1}(0) = 0$ and $H_0(0)$ is the time 0 value of the CLE, we have established the following proposition, which is essentially inspired by an integral representation of an American option from Jamshidian [1992], as presented in Proposition 1.10.7.

Proposition 18.2.1. *The time 0 value $H_0(0)$ of a callable Libor exotic is equal to the sum of European options on the difference between the exercise and hold values at all exercise dates,*

$$H_0(0) = \sum_{n=1}^{N-1} \mathbb{E}\left(B(T_n)^{-1}(U_n(T_n) - H_n(T_n))^+\right). \quad (18.8)$$

In Proposition 18.2.1 each of the terms $\mathbb{E}(B(T_n)^{-1}(U_n(T_n) - H_n(T_n))^+)$ can be interpreted as a “marginal” exercise value, i.e. as the incremental value of having an exercise at time T_n . The total CLE value is then equal to the sum of the marginal exercise values.

18.3 Monte Carlo Valuation

If a model admits a low-dimensional Markovian representation, then PDE methods are available for valuation of CLEs, and the backward recursion (18.7) is easy to implement, see Section 2.7.4. The situation is more complicated with Monte Carlo based models; fortunately, a range of methods for approximate solutions of optimal exercise problems in Monte Carlo have been developed. The mechanics of the scheme have been broadly outlined in Section 3.5.4, and we now proceed to discuss implementation details for the CLE class.

18.3.1 Regression-Based Valuation of CLEs, Basic Scheme

We start with a basic regression-based method for estimating the value of a CLE. As we recall from the discussion of Section 3.5.4, the regression-based LS (for *Least Squares*) scheme builds on the idea that the expected value of a random variable conditioned on information at time T can be calculated in a Monte Carlo simulation by regressing the random variable against simulated state variables of the model at time T .

To make matters precise, we introduce some notation. Let $\zeta(t) = (\zeta_1(t), \dots, \zeta_q(t))^\top$ be a q -dimensional vector process of *regression variables*, to be defined later. For a given Monte Carlo path ω , let us denote the value of a random variable X on that path by $X(\omega)$. Suppose K paths $\omega_1, \dots, \omega_K$ are generated. For a random variable X , we denote by

$$\mathcal{R}_T(X)$$

the results of regression of the K -dimensional vector $(X(\omega_1), \dots, X(\omega_K))$ on the $K \times q$ matrix of regression variable observations at time T , $(\zeta(T, \omega_1), \dots, \zeta(T, \omega_K))^\top$, i.e.

$$\mathcal{R}_T(X) = \zeta(T)^\top \beta,$$

where the q -dimensional column vector β is obtained by, for instance³, solving the minimization problem

$$\left\| \left(X(\omega_1) - \zeta(T, \omega_1)^\top \beta, \dots, X(\omega_K) - \zeta(T, \omega_K)^\top \beta \right) \right\|^2 \rightarrow \min. \quad (18.9)$$

This least-squares problem can be solved in closed form as explained in Section 3.5.4; to link our discussion here to the results of that section we simply need to set $\zeta_i(t) = \psi_i(x(t))$, $i = 1, \dots, q$, where ψ 's and x 's are defined in Section 3.5.4. For future reference, we denote the solution vector β to (18.9) by

³We discuss the details of the implementation of the regression algorithm later, in Section 18.3.10.

$$\mathcal{C}(\mathcal{R}_T(X)),$$

with the notation meant to be read as “coefficients of the regression of X ”, so that

$$\mathcal{R}_T(X) = \mathcal{C}(\mathcal{R}_T(X))^\top \zeta(T). \quad (18.10)$$

As it turns out, there are several possible LS schemes for CLE valuation. The most basic one is based on the idea of simply replacing the conditional expected value operator E_T in (18.7) with the *regression operator* \mathcal{R}_T introduced above, i.e. to write

$$\tilde{H}_{n-1}(T_{n-1}) = \mathcal{R}_{T_{n-1}} \left(\frac{B(T_{n-1})}{B(T_n)} \max \left(\tilde{H}_n(T_n), U_n(T_n) \right) \right), \quad (18.11)$$

for $n = N - 1, \dots, 1$, where \tilde{H}_n is an approximation to the true hold value H_n . This approach was originally suggested in Carrière [1996] and Tsitsiklis and Roy [2001]. While we shall later describe better LS schemes than (18.11), let us nevertheless take some care in documenting all the steps necessary to apply it. Some of the steps are shared with the standard (non-callable) Monte Carlo valuation algorithm, but we list them anyway for completeness.

1. Choose and calibrate a term structure model (such as the LM model).
2. Decide on what to use for the regression variables process $\zeta(t)$ (we will have more to say about this later).
3. Simulate K paths $\omega_1, \dots, \omega_K$. For the LM model in particular, each ω_k represents one simulated path of all core Libor rates.
4. For each path ω_k calculate simulated values of the numeraire $B(T_n, \omega_k)$, $n = 1, \dots, N - 1$.
5. For each path ω_k , calculate⁴ the value $U_n(T_n, \omega_k)$ of the underlying exotic swap on all exercise dates $n = 1, \dots, N - 1$.
6. For each path ω_k , calculate the values of the q -dimensional regression variables process ζ on the exercise dates, $\zeta(T_n, \omega_k)$, $n = 1, \dots, N - 1$.
7. Set $\tilde{H}_{N-1} \equiv 0$.
8. For each $n = N - 1, \dots, 1$
 - a) Form a K -dimensional vector $V_n = (V_n(\omega_1), \dots, V_n(\omega_K))^\top$,

$$V_n(\omega_k) = \frac{B(T_{n-1}, \omega_k)}{B(T_n, \omega_k)} \max \left(\tilde{H}_n(T_n, \omega_k), U_n(T_n, \omega_k) \right) \quad (18.12)$$

for $k = 1, \dots, K$.

⁴In this basic scheme we implicitly assume that underlying value U_n at time T_n can be calculated in closed form from the simulated yield curve at time T_n (or, more generally, from the model state variables observed on and *before* time T_n). This is a strong restriction which we relax later.

b) Calculate

$$\tilde{H}_{n-1}(T_{n-1}) = \mathcal{R}_{T_{n-1}}(V_n)$$

by regressing $V_n(\omega_k)$ against the matrix of regression variables observed on date T_{n-1} , using (18.9).

9. At this point we have computed $\tilde{H}_0(T_0)$ which is an estimate of the value of the CLE. Return it.

Note that the last iteration in Step 8 ($n = 1$) involves regression on the values of $\zeta(T_0)$, $\tilde{H}_0(T_0) = \mathcal{R}_{T_0}(V_1)$. As $T_0 = 0$, $\zeta(T_0)$ is not random, so the regression here degenerates into a simple average⁵ of $V_1(\omega_1), \dots, V_1(\omega_K)$.

There are a number of shortcomings in the above scheme that make it poorly suited for industrial-strength pricing of CLEs. We list them below as they shall guide us in systematically building more refined versions of the algorithm.

1. In Step 5, there is an assumption that exercise values $U_n(T_n)$ can be computed in closed form from information available at time T_n . While this is possible for, say, simple Bermudan swaptions, more complicated exotic swaps will generally violate this assumption.
2. The use of “regressions upon regressions”, i.e. the fact that we apply $\mathcal{R}_{T_{n-1}}$ in (18.11) to (a function of) a regressed value $\tilde{H}_n(T_n)$, could lead to significant biases building up as the scheme marches backward in time.
3. In general we cannot state whether \tilde{H}_n ’s are low- or high-biased estimates of H_n ’s.

18.3.2 Regression for Underlying

To improve our basic LS scheme, we start out by examining the assumption of the basic LS scheme that exercise value at time t can be evaluated in closed form from time t state variables (typically the simulated yield curve and the state of stochastic volatility parameters). As mentioned earlier, vanilla fixed-for floating-rate swaps underlying Bermudan swaptions certainly satisfy this assumption, as their values can be obtained by discounting projected cash flows on a simulated yield curve at time t . In principle, a number of other exotic swaps could fit under the assumption as well. For example, for a callable inverse floater, the underlying swap is a collection of simple Libor rate options, the values of which could be calculated from simulated market data at times T_n , $n = 1, \dots, N - 1$, by applying option pricing formulas developed earlier in this book. While such a scheme is indeed possible, it would come at a significant computational cost of having to invoke option valuation formulas multiple times for each simulated path, i.e. easily thousands of

⁵A naive numerical implementation of regression will not have this property, see Section 18.3.10.

times overall. Moreover, closed-form caplet/swaption/CMS option pricing is rarely exact in term structure interest rate models; embedded into a backwards recursion algorithm these errors may potentially build up and skew the pricing results for the CLE.

Fortunately, it turns out that the extension of the basic scheme to arbitrary underlying swaps is simple. Specifically, we can use formula (18.2) which states that the exercise value $U_n(T_n)$ is equal to the conditional expected value of all (net) coupons paid after T_n . As we have already introduced the regression operator as a numerical proxy for conditional expected value, all we need to do now is to approximate $U_n(T_n)$ with the regressed value of all (net) coupons paid after T_n . Of course a coupon paid at time T_{i+1} is always⁶ measurable with respect to $\mathcal{F}_{T_{i+1}}$ or, equivalently, can be calculated from the knowledge of the simulated model state variables up to and including T_{i+1} . So, all coupon values are known once a given path is simulated, and we can extend the basic scheme to arbitrary underlyings with the following two modifications. First, we replace Step 5 with

- 5a. For each path ω_k , calculate the values of all net coupons $X_n(T_n, \omega_k)$, $n = 1, \dots, N - 1$.

Second, we must replace the formula (18.12) in Step 8 of the basic scheme with

$$V_n(\omega_k) = \frac{B(T_{n-1}, \omega_k)}{B(T_n, \omega_k)} \max \left(\tilde{H}_n(T_n, \omega_k), \tilde{U}_n(T_n, \omega_k) \right), \quad k = 1, \dots, K, \quad (18.13)$$

where

$$\tilde{U}_n(T_n) = \mathcal{R}_{T_n} \left(\sum_{i=n}^{N-1} \frac{B(T_n)}{B(T_{i+1})} X_i \right). \quad (18.14)$$

We can write (18.14) in a backward-recursive format,

$$Y_n = \frac{B(T_n)}{B(T_{n+1})} (X_n + Y_{n+1}), \quad \tilde{U}_n(T_n) = \mathcal{R}_{T_n}(Y_n), \quad n = N - 1, \dots, 1, \quad (18.15)$$

where we start from $Y_N \equiv 0$. In this form it fits nicely into the backward recursion of the basic LS scheme.

Interestingly, we can rewrite (18.2) (for $t = T_n$) as

$$U_n(T_n) = \mathbb{E}_{T_n} (B(T_n)B(T_{n+1})^{-1} (X_n + U_{n+1}(T_{n+1}))), \quad (18.16)$$

which gives raise to an *alternative backward scheme* for $\tilde{U}_n(T_n)$,

$$\tilde{U}_n(T_n) = \mathcal{R}_{T_n} \left(B(T_n)B(T_{n+1})^{-1} \left(X_n + \tilde{U}_{n+1}(T_{n+1}) \right) \right), \quad n = N - 1, \dots, 1. \quad (18.17)$$

⁶The value of a coupon must obviously be known at the time it is paid.

While (18.16) is trivially equivalent to (18.2) due to the additivity of the conditional expected values and the iterated conditional expectations property, (18.17) is *not* equivalent to (18.15). In (18.17), on step n we add (discounted) coupon X_n to an already-regressed value of future coupons $\tilde{U}_{n+1}(T_{n+1})$, whereas in (18.15) we sum up values of un-regressed coupons from $N - 1$ to n and then regress them to obtain \tilde{U}_n . The difference between the two schemes originates with the fact that the regression operator does not satisfy an equivalent to the iterated conditional expectations property⁷ (see footnote 3 of Chapter 4). An inquisitive reader may ask whether one scheme is better (in the sense of producing smaller bias for the value of the CLE) than the other. The answer is not straightforward. While (18.15) avoids applying regression to the output of other regressions and thus could be expected to produce lower bias, empirical evidence suggests that no scheme is universally better than the other for all CLEs. We consequently recommend a flexible implementation that can use both schemes.

18.3.3 Valuing CLE as a Cancelable Note

Using regressions for the underlying is not the only approach that extends the basic scheme in Section 18.3.2 to arbitrary underlying swaps. An alternative scheme is based on the idea of representing a CLE as a *cancelable note*. To describe this in more detail, let us denote

$$G_n(t) = H_n(t) - U_n(t),$$

and obtain from (18.7) and (18.16) that

$$\begin{aligned} G_{n-1}(T_{n-1}) &= H_{n-1}(T_{n-1}) - U_{n-1}(T_{n-1}) \\ &= E_{T_{n-1}}(B(T_{n-1})B(T_n)^{-1}(\max(H_n(T_n), U_n(T_n)) - (X_{n-1} + U_n(T_n)))) \\ &= E_{T_{n-1}}\left(B(T_{n-1})B(T_n)^{-1}\left((H_n(T_n) - U_n(T_n))^+ - X_{n-1}\right)\right) \\ &= E_{T_{n-1}}(B(T_{n-1})B(T_n)^{-1}(-X_{n-1} + G_n(T_n)^+)), \end{aligned} \quad (18.18)$$

$n = N, \dots, 1$, where we for uniformity of notation have introduced a “fake” coupon at time zero,

$$X_0 = 0.$$

We see that the value $G_0(0)$ is the value of the swap that pays (net) coupons $-X_n$ on dates T_n , $n = 0, \dots, N - 1$, plus the right to cancel it, at zero cost, on any of the exercise dates T_1, \dots, T_{N-1} . In fact, as explained in Section 5.14, it is this structure (a callable structured note) that a bank usually sells to clients. While the representation of the cancelable note as a CLE plus the

⁷A similar issue arises with the regressed hold values, as discussed in Section 18.3.4.

underlying non-callable swap is often convenient for risk management, for the purposes of valuation it is, in fact, often useful to consider the original format.

The LS version of (18.18) is, naturally, given by

$$V_n(\omega_k) = \frac{B(T_{n-1}, \omega_k)}{B(T_n, \omega_k)} \left(-X_{n-1}(\omega_k) + \tilde{G}_n(T_n, \omega_k)^+ \right), \quad k = 1, \dots, K,$$

followed by

$$\tilde{G}_{n-1}(T_{n-1}) = \mathcal{R}_{T_{n-1}}(V_n), \quad (18.19)$$

for each $n = N - 1, \dots, 1$. The starting point is given by $\tilde{G}_{N-1}(T_{N-1}) \equiv 0$. We trust the reader should have no problem amending the basic scheme to use the cancelable note representation.

Interestingly, by linearity (which *does* hold for the regression operator, unlike the tower property), we can rewrite (18.19) as

$$\begin{aligned} \tilde{G}_{n-1}(T_{n-1}) &= -\mathcal{R}_{T_{n-1}}(B(T_{n-1})B(T_n)^{-1}X_{n-1}) \\ &\quad + \mathcal{R}_{T_{n-1}}\left(B(T_{n-1})B(T_n)^{-1}\tilde{G}_n(T_n)^+\right), \quad n = N - 1, \dots, 1. \end{aligned}$$

In other words, to get the value of the cancelable note at time T_{n-1} , start with its value $\tilde{G}_n(T_n)$ at time T_n , apply the optimal cancelability condition $\tilde{G}_n(T_n) \rightarrow \tilde{G}_n(T_n)^+$, discount to T_{n-1} , regress on $\zeta(T_{n-1})$, and then add the time T_{n-1} value of the $(n-1)$ -th coupon. If X_{n-1} is actually $\mathcal{F}_{T_{n-1}}$ -measurable, as is often the case (exceptions include range accrual coupons and coupons that depend on rates observed in-arrears, i.e. at the end of the observation period rather than the beginning), the scheme simplifies a bit more,

$$\begin{aligned} \tilde{G}_{n-1}(T_{n-1}) &= -P(T_{n-1}, T_n) X_{n-1} \\ &\quad + \mathcal{R}_{T_{n-1}}\left(B(T_{n-1})B(T_n)^{-1}\tilde{G}_n(T_n)^+\right), \quad n = N - 1, \dots, 1. \quad (18.20) \end{aligned}$$

Once the value $\tilde{G}_0(0)$ has been calculated, the estimate of the CLE value $H_0(0)$ can be obtained via $H_0(0) = \tilde{G}_0(0) + U_0(0)$, where U_0 may, if not available in closed form, be calculated via a (standard) Monte Carlo algorithm for (18.2). In estimating U_0 , we would normally reuse the paths ω_k , $k = 1, \dots, K$ that were used for the computation of $\tilde{G}_0(0)$.

18.3.4 Using Regressed Variables for Decision Only

Our second criticism of the basic LS scheme of Section 18.3.1 focused on the issue that the values that are regressed at time T_{n-1} themselves come as the result of a regression at time T_n . Such compounded regression could lead

to substantial biases. Interestingly, a small modification of the algorithm reduces this bias significantly.

Let us work with the cancelable note scheme (18.20), although the same idea can easily be applied to the basic CLE scheme. Going back to the expression (18.6) for the value of the CLE expressed as the sum of all coupons paid post optimal exercise, we note that a similar formula holds for cancelable notes,

$$G_0(0) = -E \left(\sum_{n=0}^{\eta-1} B(T_{n+1})^{-1} X_n \right) \quad (18.21)$$

(recall that $X_0 = 0$). The exercise index η here is the same as in (18.6), since the optimal cancel time for the cancelable note G is the same as the optimal exercise time for the CLE H . We rewrite the formula as

$$G_0(0) = -E \left(\sum_{n=0}^{N-1} B(T_{n+1})^{-1} X_n 1_{\{\eta>n\}} \right),$$

and note that

$$1_{\{\eta>n\}} = \prod_{i=1}^n 1_{\{G_i(T_i)>0\}}. \quad (18.22)$$

Let us define V_n 's recursively,

$$V_n = B(T_n) B(T_{n+1})^{-1} (-X_n + 1_{\{G_{n+1}(T_{n+1})>0\}} V_{n+1}), \quad n = N-1, \dots, 0, \quad (18.23)$$

with $V_N = 0$. Then

$$V_0 = - \sum_{n=0}^{N-1} B(T_{n+1})^{-1} \left(\prod_{i=1}^n 1_{\{G_i(T_i)>0\}} \right) X_n$$

and, computing the expected value and using (18.22), we obtain that

$$G_0(0) = E(V_0).$$

Moreover,

$$G_n(T_n) = E_{T_n}(V_n). \quad (18.24)$$

We note that the recursion for V_n involves the value of the cancelable note for exercise decisions only, through the indicators $1_{\{G_i(T_i)>0\}}$, whereas the coupon values that are added up, X_i 's, are never regressed. Following Longstaff and Schwartz [2001], we can take advantage of this observation by defining a new approximation $\hat{G}_n(T_n)$, $n = 1, \dots, N-1$, to the true value of the cancelable note by

$$\widehat{G}_{n-1}(T_{n-1}) = B(T_{n-1})B(T_n)^{-1} \left(-X_{n-1} + 1_{\{\tilde{G}_n(T_n) > 0\}} \widehat{G}_n(T_n) \right), \quad (18.25)$$

$$\widetilde{G}_{n-1}(T_{n-1}) = \mathcal{R}_{T_{n-1}}(\widehat{G}_{n-1}(T_{n-1})), \quad (18.26)$$

defined backwards for $n = N, \dots, 1$. The first equation here comes from (18.23) and the second from (18.24). We emphasize that regression is only used to establish the exercise indicator functions in (18.25), i.e. whether and where to exercise on each path. In (18.24), $\widehat{G}_0(0) = \widehat{G}_0(0, \omega)$ is the (random) accumulation of discounted coupons up to exercise, and our estimate of the true value of the cancelable note G_0 is given by $\widetilde{G}_0(0)$, i.e. the simple average of the realizations of $\widehat{G}_0(0, \omega_1), \dots, \widehat{G}_0(0, \omega_K)$. Clearly, this is the numerical equivalent of the (unconditional) expected value in (18.21). We find that the scheme (18.25)–(18.26) typically has significantly less bias than the naive scheme (18.20).

While (18.25)–(18.26) is the scheme that we recommend for most applications, Egloff et al. [2007] introduce a “blend” of (18.25)–(18.26) and (18.20) with a tunable parameter that can be optimized over to select a scheme with the lowest bias. To briefly outline this idea, we note that the scheme (18.25)–(18.26) only uses un-regressed values of coupons while (18.20) always uses regressed values. A “blended” scheme uses the first few coupons that are unregressed, while the rest are regressed. For example, since we can write

$$G_0(0) = -E \left(\sum_{n=0}^{m-1} B(T_{n+1})^{-1} X_n 1_{\{n>m\}} + B(T_m)^{-1} G_m(T_m) 1_{\{G_m(T_m)>0\}} \right)$$

for any $m = 1, \dots, N-1$, we can replace (18.25) with

$$\begin{aligned} \widehat{G}_{n-1}(T_{n-1}) &= - \sum_{i=n-1}^{\lceil n+m-1 \rceil - 1} \left(\frac{B(T_{n-1})}{B(T_{i+1})} \left(\prod_{j=n}^i 1_{\{\tilde{G}_j(T_j) > 0\}} \right) X_i \right) \\ &+ \mathcal{R}_{T_{n-1}} \left(\frac{B(T_{n-1})}{B(T_{\lceil n+m-1 \rceil})} \left(\prod_{j=n}^{\lceil n+m-1 \rceil} 1_{\{\tilde{G}_j(T_j) > 0\}} \right) \widetilde{G}_{\lceil n+m-1 \rceil}(T_{\lceil n+m-1 \rceil}) \right) \end{aligned}$$

coupled with (18.26), where we have denoted $\lceil l \rceil \triangleq l \wedge N$. We would then run the regression algorithm for different m 's (reusing the paths, of course) and choose the value of m that would give us the highest value⁸. We recover (18.20) for $m = 1$ and (18.25)–(18.26) for $m = +\infty$.

18.3.5 Regression Valuation with Boundary Optimization

The regression-based method can sometimes be improved with methods similar to those of the parametric boundary optimization discussed in Sec-

⁸Jumping a bit ahead, we note that it is essential to use this optimization in conjunction with an independent post-simulation, see Section 18.3.6.

tion 3.5.2, an idea popularized by Bender et al. [2006]. Starting with, for example, scheme (18.25)–(18.26), we fix a collection of trigger thresholds $\mathbf{h} = (h_1, \dots, h_{N-1})$ and rewrite (18.25)–(18.26) as (note dependence on the triggers in the indicator functions)

$$\begin{aligned}\widehat{G}_{n-1}(T_{n-1}, \mathbf{h}) &= B(T_{n-1})B(T_n)^{-1} \left(-X_{n-1} + 1_{\{\tilde{G}_n(T_n, \mathbf{h}) > h_n\}} \widehat{G}_n(T_n, \mathbf{h}) \right), \\ \tilde{G}_{n-1}(T_{n-1}, \mathbf{h}) &= \mathcal{R}_{T_{n-1}}(\widehat{G}_{n-1}(T_{n-1}, \mathbf{h})),\end{aligned}$$

$n = N-1, \dots, 1$. For a given set \mathbf{h} , the scheme produces an estimate $\tilde{G}_0(0, \mathbf{h})$ of the value of the cancelable note. Then, iterating over \mathbf{h} , we find the optimal (highest) value of \tilde{G}_0 and return this as our improved estimate for the cancelable note. Of course, if our original regression was fundamentally sound, we should see the optimal value of \mathbf{h} being very close to $(0, 0, \dots, 0)$, in which case the trigger iteration adds no value. To the extent, however, that the original LS method produces significantly sub-optimal exercise decisions (e.g. due to a poor choice of regression variables), the trigger iteration may lead to pick-up of substantial additional value.

The search for the optimal value of the triggers can be efficiently organized as a sequence of $N-1$ one-dimensional optimizations, along the same lines as the algorithm in Section 3.5.3. In particular, we see that $\widehat{G}_{n-1}(T_{n-1}, \mathbf{h})$ depends on h_n, \dots, h_{N-1} only. Moreover, if for a given n , h_n maximizes the value of $\tilde{G}_0(0, \mathbf{h})$ then it also maximizes the value of the cancelable note with first $n-1$ exercise dates removed, i.e. the value $\mathcal{R}_0(\widehat{G}_{n-1}(T_{n-1}, \mathbf{h}))$ (recall that $\mathcal{R}_0(X) = K^{-1} \sum X(\omega_k)$, i.e. just the average of path values of X). Hence, we find the optimal value of the n -th trigger h_n^* via

$$h_n^* = \underset{h_n}{\operatorname{argmin}} \mathcal{R}_0(\widehat{G}_{n-1}(T_{n-1}, (h_n, h_{n+1}^*, \dots, h_{N-1}^*))), \quad n = N-1, \dots, 1,$$

where, slightly abusing notation, we use $(h_n, h_{n+1}^*, \dots, h_{N-1}^*)$ to denote a vector \mathbf{h} with the last $N-n-1$ elements fixed to the optimal values found on previous steps (and first $n-1$ elements irrelevant). The optimization problem above is easy to solve, but we remind the reader of our comments from Section 3.5.3, where we noted that for a finite-path simulation, the objective functions in each of optimization problems will not be smooth, so one should avoid the use of a derivatives-based numerical optimizer.

18.3.6 Lower Bound via Regression Scheme

All the variations of the regression scheme developed so far produce an estimate of the value of the CLE H_0 (or, equivalently, the value of the corresponding cancelable note G_0) but the bias of the estimate is generally unknown. On one hand, the exercise decisions used in the schemes are necessarily suboptimal, as we use estimates, rather than actual values, of

hold/underlying variables to define them. This, in isolation, suggests that our estimate should be biased low. But on the other hand, our schemes use the same set of sample paths to estimate the exercise decision as to calculate the values of the security if it is exercised. This could lead to an upward bias in the estimate as some amount of future information can affect our decision to exercise, leading to a “perfect foresight” high-bias, see (3.110).

For risk assessment and for price quotation, it is typically desirable to know the sign of biases in computed security prices. To control the sign of the bias in CLE valuation by regression methods, we can follow the advice of Section 3.5 and use the LS regression scheme to estimate the exercise decision rule only, while using an independent simulation to calculate the value of the CLE *given* that exercise rule. A typical implementation algorithm is outlined below.

1. Run the basic regression-based scheme in Section 18.3.1 or any of the alternatives in Sections 18.3.2–18.3.5.
2. The output of the regression are the regression coefficients for the hold and exercise values at all exercise times, $\mathcal{C}(\tilde{H}_n(T_n))$, $\mathcal{C}(\tilde{U}_n(T_n))$, $n = 1, \dots, N - 1$. See (18.10).
3. Simulate additional K' paths, $\omega'_1, \dots, \omega'_{K'}$, that are independent from the paths used in the regression scheme.
4. For each path ω'_k , calculate the values of the q -dimensional regression variables process ζ on the exercise dates, $\zeta(T_n, \omega'_k)$, $n = 1, \dots, N - 1$.
5. For each path ω'_k , calculate an estimate of the exercise index $\tilde{\eta}$ by

$$\begin{aligned}\tilde{\eta}(\omega'_k) &= \min \left\{ n \geq 1 : \mathcal{C} \left(\tilde{U}_n(T_n) \right)^T \zeta(T_n, \omega'_k) \right. \\ &\quad \left. \geq \mathcal{C} \left(\tilde{H}_n(T_n) \right)^T \zeta(T_n, \omega'_k) \right\} \wedge N.\end{aligned}\quad (18.27)$$

In simple terms the exercise is based on the regression estimates of the exercise and hold values obtained in the basic scheme, applied to the new values of the regression variables ζ .

6. Calculate the CLE value as the Monte Carlo value of a knock-in discrete barrier option based on the exercise rule $\tilde{\eta}$,

$$H_0(0) \approx \frac{1}{K'} \sum_{k=1}^{K'} \left(\sum_{n=\tilde{\eta}(\omega'_k)}^{N-1} B(T_{n+1}, \omega'_k)^{-1} X_n(\omega'_k) \right). \quad (18.28)$$

The guaranteed low bias of this two-stage scheme comes at an additional cost of simulating and evaluating payoffs on extra K' paths. For the record, we often choose $K' \approx 10,000$ to $100,000$ and $K \approx K'/4$ to $K'/2$. So the cost is not inconsiderable, on average slowing down a valuation by a factor of up to 2 or so. Importantly, this additional cost is not incurred in the evaluation of most risk sensitivities, since we can reuse the exercise

boundary in calculations of first-order sensitivities, as we explain in Chapter 24. However, if performance is an issue, it is worth pointing out that we often find the values obtained by the scheme (18.25)–(18.26) (with only a single batch of paths simulated) to be pretty close to those from the more costly two-stage scheme (18.28). In particular, it appears that in many practical situations (18.25)–(18.26) is biased low, even though it is not guaranteed to be so.

It is easy to see the strong connection between the scheme (18.25)–(18.26) and (18.28). If, instead of the independent paths in the second stage of (18.28) we used the same paths as in the regression scheme, i.e. $K' = K$ and $\omega'_k = \omega_k$ for $k = 1, \dots, K$, then the values produced by the two scheme would be exactly the same. We leave the verification of this simple fact to the reader.

18.3.7 Iterative Improvement of Lower Bound

Consider the cancelable note that pays (net) coupons $-X_n$ at T_{n+1} , $n = 0, \dots, N-1$. Suppose we are given some exercise policy, generally suboptimal. An exercise policy is, in essence, a stopping time index α that specifies when the note is canceled. For technical reasons we want this exercise strategy to be specified not just for the original cancelable note but also for cancelable notes with the first k , $k = 1, \dots, N-2$ exercise dates removed. This is most conveniently expressed as a collection of stopping times $\alpha = (\alpha_0, \dots, \alpha_{N-1})$, with $\alpha_0 = \alpha$, that satisfy the following *exercise policy consistency conditions*,

$$n + 1 \leq \alpha_n \leq N, \quad \alpha_{N-1} \equiv N,$$

and

$$\alpha_n > n + 1 \Rightarrow \alpha_n = \alpha_{n+1}, \quad n = 0, \dots, N-2.$$

Let us consider the part of the note that includes coupons $-X_n, \dots, -X_{N-1}$ only. Then the value of this note, exercised per stopping time α_k , $k \geq n$, is given at time T_n by

$$\begin{aligned} G_n^{\alpha_k}(T_n) &= -B(T_n)E_{T_n}\left(\sum_{i=n}^{\alpha_k-1} B(T_{i+1})^{-1} X_i\right) \\ &= -B(T_n)E_{T_n}\left(\sum_{i=n}^{N-1} B(T_{i+1})^{-1} X_i 1_{\{\alpha_k > i\}}\right). \end{aligned}$$

We note that $\alpha_k \geq k + 1$, so the coupons $-X_n, \dots, -X_k$ are always paid. The *optimal* exercise policy $\eta \triangleq (\eta_0, \dots, \eta_{N-1})$, as defined by (18.5), satisfies the consistency conditions, and we of course have

$$G_n^{\eta_n}(T_n) = G_n(T_n), \quad n = 0, \dots, N-1,$$

where on the right-hand side we have the actual values of the remaining parts of the cancelable note. The approximations to the optimal exercise rule

we developed in previous sections, such as (18.27), satisfy the consistency conditions as well.

As pointed out previously, for any exercise policy α we have

$$G_n^{\alpha_n}(T_n) \leq G_n(T_n), \quad n = 0, \dots, N-1,$$

so an exercise policy gives us a way to obtain a lower bound, namely $G_0^{\alpha_0}(0)$, for the time 0 cancelable note value $G_0(0)$. It turns out that the general theory of optimal stopping tells us how to improve a given exercise policy, i.e. how to find another exercise policy that would be better, in the sense of producing a higher lower bound. The improvements could be iterated, eventually (after $N-1$ iterations) converging to the optimal exercise policy irrespective of the starting point. This “policy iteration” method was first applied to the pricing of callable derivatives by Bender et al. [2006] and Kolodko and Schoenmakers [2006]. To demonstrate how policy iteration works, for a given policy α , let us define its improvement $\hat{\alpha} = (\hat{\alpha}_0, \dots, \hat{\alpha}_{N-1})$ by

$$\hat{\alpha}_n = \min \left\{ k > n : \max_{j \geq k} G_k^{\alpha_j}(T_k) \leq 0 \right\} \wedge N, \quad n = 0, \dots, N-1. \quad (18.29)$$

To understand this definition, first note that $G_k^{\alpha_j}(T_k)$ is the time- T_k value of the coupon stream from T_k onwards, exercised per stopping time α_j , i.e. we basically add up the (discounted) coupons $-X_k, \dots, -X_j$ and thereafter the remaining coupons subject to the original exercise rule. Hence, the improved exercise rules states that we should exercise (cancel the note) at the first date T_k for which holding on to the note under the original exercise policy does not make sense, specifically if the maximum of the remaining value of the note over all original exercise rules that go strictly past T_k is non-positive.

The proof that $\hat{\alpha}$ is an improvement over α , i.e. that

$$G_0^{\hat{\alpha}_0}(0) \geq G_0^{\alpha_0}(0)$$

can be found in Kolodko and Schoenmakers [2006]. The improvement could be applied to the policy $\hat{\alpha}$ as well, and this can be iterated multiple times. After $N-1$ iterations the exercise policy converges to the optimal one, as also proven in Kolodko and Schoenmakers [2006]. .

While in theory we can find the optimal exercise strategy by iteration starting from any initial policy — such as a trivial policy $\alpha_n = n+1$ for any $n = 0, \dots, N-1$ — performing more than a few iterations is impractical, as shall be discussed shortly. In practice, a sensible strategy would apply just one round of improvements to an already decent exercise policy, such as (18.27) obtained by the basic scheme of Section 18.3.1 or its variations.

The high numerical cost of the policy iteration stems from the presence of terms like $G_k^{\alpha_j}(T_k)$ in (18.29). These are the time T_k values of the coupon stream exercised per some rule, and are rarely, if ever, available in closed

form. One might guess that we could estimate the required conditional expected values via regressions, as we have done for other quantities. Alas, using regressed values in (18.29) in place of true conditional expected values generally leads to *no* policy improvements. Consider, for example, the scheme (18.25)–(18.26), with the corresponding exercise policy defined by

$$\alpha_n = \min \left\{ k > n : \tilde{G}_k(T_k) \leq 0 \right\} \wedge N, \quad n = 0, \dots, N - 1.$$

Then the iterated exercise policy would base the exercise decision on $\max_{j \geq k} \tilde{G}_k^{\alpha_j}(T_k)$. However, as α_n 's are constructed to be the *optimal* stopping policy for $\tilde{G}_n(T_n)$'s, we always have

$$\tilde{G}_k^{\alpha_j}(T_k) \leq \tilde{G}_k^{\alpha_k}(T_k) = \tilde{G}_k(T_k),$$

and

$$\max_{j \geq k} \tilde{G}_k^{\alpha_j}(T_k) = \tilde{G}_k(T_k)$$

for any k . Thus the “improved” policy $\hat{\alpha}$ computed from (18.29) would coincide with α , the original policy.

As regression methods cannot be used, an alternative is to resort to “brute-force” nested simulation to compute unbiased estimates for the conditional expectations $G_k^{\alpha_j}(T_k)$. The basic idea is simple: if we need to estimate $E_T(X)$ for some random variable X then, for each Monte Carlo path ω , we simulate a number of additional sub-paths with the simulated model state at time T on path ω as a starting point, and then estimate the conditional expected value on path ω by averaging the values of X realized on all sub-paths⁹. For policy iteration, such sub-paths must be launched for each (outer) simulated path ω_k , $k = 1, \dots, K$, for each exercise time T_n , $n = 1, \dots, N - 1$, so the computational expense is quite considerable even with a modest number of sub-paths. Bender et al. [2006] recommend using control variates to speed up valuation, and Beveridge and Joshi [2009] list a number of additional suggestions to improve computational performance. Nevertheless, there is no doubt that to keep Monte Carlo simulation error of the conditional expectation low enough for the policy iteration scheme to be effective, a substantial numerical effort is required. As such, we do not recommend routine application of policy iteration in the pricing of CLEs. This, of course, will require that the regression estimates of the exercise policy¹⁰ are sufficiently accurate as is, something that to a large extent depends on how careful we are in choosing the regression variable vector $\zeta(t)$. We turn to this topic in Section 18.3.9, but first we need to address the fundamental problem of how to test whether a given exercise policy is close to optimal in the first place. One useful approach to this problem is to use the estimated

⁹For more detailed discussion of nested Monte Carlo simulation, see Section 18.3.8.

¹⁰Which may include refinements such as that in Section 18.3.5, of course.

exercise policy to construct an *upper bound* for the option price, which, if close to the lower bound, will give us confidence that our exercise policy is close to optimal. Section 18.3.8 below discusses this technique in detail.

18.3.8 Upper Bound

Computing a lower bound for a CLE price is straightforward: pick some exercise policy and price the CLE by Monte Carlo methods. Assuming our computation of the exercise policy did not “cheat” by using information from the Monte Carlo trials subsequently used for valuation, the resulting price estimates will always have a non-positive bias, as the chosen exercise strategy will almost certainly be suboptimal. The lower bound algorithm of Section 18.3.6 was based on precisely such a strategy. To complement a computed lower bound for the CLE, we are now interested in using the lower bound exercise policy to construct an upper bound for the CLE price. Taken together with the lower bound, the upper bound can be used to construct a valid confidence interval for the CLE price; this, in turn, will allow us to assess the quality of the exercise policy.

18.3.8.1 Basic Ideas

Our strategy to construct an upper bound for CLEs will draw directly on the duality results in Section 1.10.2 and the generic upper bound simulation ideas in Section 3.5.5. To formulate these results in the CLE setting, let us start with the description (18.4)–(18.6):

$$H_0(0) = \sup_{\xi \in \mathcal{T}_0} E(B(T_\xi)^{-1} U_\xi(T_\xi)) = E(B(T_\eta)^{-1} U_\eta(T_\eta)). \quad (18.30)$$

Following Section 1.10.2, let \mathcal{K} denote the space of adapted martingales M for which $\sup_{\xi \in \mathcal{T}_0} E|M(T_\xi)| < \infty$. For any martingale $M \in \mathcal{K}$, (1.71) demonstrates that

$$H_0(0) \leq M(0) + E \left(\max_{n=1,\dots,N-1} \left(\frac{U_n(T_n)}{B(T_n)} - M(T_n) \right) \right). \quad (18.31)$$

Also, we know from the duality result (1.72) that this upper bound will become an *equality* provided that M is chosen to be the martingale component of the (supermartingale) deflated value process of the CLE. To emphasize this result, set

$$V_{\text{CLE}}(t) = B(t) E_t (B(T_n)^{-1} \max(U_n(T_n), H_n(T_n))), \quad t \in (T_{n-1}, T_n], \quad (18.32)$$

and use the Doob-Meyer decomposition of Section 1.10.2 to write $V_{\text{CLE}}(t)/B(t) = m_{\text{CLE}}(t) - A(t)$, where $m_{\text{CLE}}(t)$ is a martingale and $A(t)$ an increasing predictable process with $A(0) = 0$. Then setting $M(t) = m_{\text{CLE}}(t)$ yields

$$m_{\text{CLE}}(0) + \mathbb{E} \left(\max_{n=1,\dots,N-1} \left(\frac{U_n(T_n)}{B(T_n)} - m_{\text{CLE}}(T_n) \right) \right) = H_0(0). \quad (18.33)$$

If the underlying model is driven by a vector-valued \mathbb{Q}^B -Brownian motion $W(t)$, the martingale representation theorem (Theorem 1.1.4) shows that any martingale M in (18.31) must be of the form

$$M(t) = \int_0^t \sigma(s)^\top dW(s), \quad (18.34)$$

for some adapted vector-process $\sigma(t)$ satisfying the usual conditions required for the stochastic integral to be proper martingale. Clearly, however, if $\sigma(t)$ is chosen arbitrarily, the resulting upper bound computed from (18.31) is likely to be very loose, and probably not particularly useful. While (18.33) is of little immediate practical use (since we do not know the process $V_{\text{CLE}}(t)/B(t)$), it does suggest that for a chosen martingale $M(t)$ in (18.31) to produce a tight upper bound, it needs to be “close” to $m_{\text{CLE}}(t)$.

Several strategies have been proposed for constructing a good martingale $M(t)$. When working in a simple model setup on simple payouts, sometimes one can make inspired guesses for what $M(t)$ should be. For instance, in a one-dimensional Black-Scholes model, Rogers [2001] shows that using the numeraire-deflated European put option price (which is analytically known) as a guess for $M(t)$ generates good bounds for a Bermudan put option price. This approach, however, does not easily generalize to the CLE setting with its more complicated model and exercise payouts.

18.3.8.2 Nested Simulation (NS) Algorithm

Andersen and Broadie [2004] propose a general strategy for generating upper bounds, starting from any approximation $\tilde{\eta}$ to the optimal exercise strategy η . Typically, this approximation would originate from an LS regression, e.g. as in Section 18.3.6, or from an optimization of a parametric formulation of the exercise rule, as in Section 3.5.2 or Section 19.6.2. In a nutshell, the algorithm in Andersen and Broadie [2004] uses nested simulation — also known as “simulation within a simulation” and already mentioned in Section 18.3.7 — to construct an estimate of the low-bound value process $\tilde{V}_{\text{CLE}}(t)$ generated from $\tilde{\eta}$. The martingale component of the numeraire-deflated value of this process is then used as $M(t)$ in (18.31).

To outline the basic nested simulation (NS) algorithm in further detail, let us work on the exercise time line $\{T_1, T_2, \dots, T_{N-1}\}$ and define

$$\begin{aligned} M(T_{n+1}) - M(T_n) &= \frac{\tilde{V}_{\text{CLE}}(T_{n+1})}{B(T_{n+1})} - \mathbb{E}_{T_n} \left(\frac{\tilde{V}_{\text{CLE}}(T_{n+1})}{B(T_{n+1})} \right) \\ &= \mathbb{E}_{T_{n+1}} \left(\frac{U_{\tilde{\eta}_n}(T_{\tilde{\eta}_n})}{B(T_{\tilde{\eta}_n})} \right) - \mathbb{E}_{T_n} \left(\frac{U_{\tilde{\eta}_n}(T_{\tilde{\eta}_n})}{B(T_{\tilde{\eta}_n})} \right), \end{aligned} \quad (18.35)$$

with $M(0) = \tilde{V}_{\text{CLE}}(0)$ and $M(T_1) = \tilde{V}_{\text{CLE}}(T_1)/B(T_1)$. In the second equality of (18.35), we use $\tilde{\eta}_n$ to denote the restriction of our approximate exercise policy exercise to the index set $\{n+1, \dots, N\}$, such that $\tilde{\eta}_0 = \tilde{\eta}$. For instance, if we use the algorithm in Section 18.3.6, we have (compare to (18.27))

$$\tilde{\eta}_n = \min \left\{ k \geq n+1 : U_k(T_k) \geq \mathcal{C} \left(\tilde{H}_k(T_k) \right)^{\top} \zeta(T_k) \right\} \wedge N. \quad (18.36)$$

For convenience, let us define an $\mathcal{F}_{T_{n+1}}$ -measurable exercise indicator

$$\iota(T_{n+1}) = 1_{\{\tilde{\eta}_n = n+1\}}$$

which will be one at time T_{n+1} if our exercise policy indicates that the CLE should be exercised, and zero otherwise. We can also define hold values

$$\frac{\tilde{H}_n(T_n)}{B(T_n)} = \mathbb{E}_{T_n} \left(\frac{U(T_{\tilde{\eta}_n})}{B(T_{\tilde{\eta}_n})} \right),$$

in which case (18.35) can be rewritten as

$$\begin{aligned} M(T_{n+1}) - M(T_n) &= \iota(T_{n+1}) \frac{U_{n+1}(T_{n+1})}{B(T_{n+1})} \\ &\quad + (1 - \iota(T_{n+1})) \frac{\tilde{H}_{n+1}(T_{n+1})}{B(T_{n+1})} - \frac{\tilde{H}_n(T_n)}{B(T_n)}. \end{aligned} \quad (18.37)$$

Notice that

$$\begin{aligned} \frac{\tilde{V}_{\text{CLE}}(T_{n+1})}{B(T_{n+1})} - \frac{\tilde{V}_{\text{CLE}}(T_n)}{B(T_n)} &= \frac{\tilde{V}_{\text{CLE}}(T_{n+1})}{B(T_{n+1})} - \mathbb{E}_{T_n} \left(\frac{\tilde{V}_{\text{CLE}}(T_{n+1})}{B(T_{n+1})} \right) \\ &\quad - \left(\frac{\tilde{V}_{\text{CLE}}(T_n)}{B(T_n)} - \mathbb{E}_{T_n} \left(\frac{\tilde{V}_{\text{CLE}}(T_{n+1})}{B(T_{n+1})} \right) \right) \\ &= M(T_{n+1}) - M(T_n) - (A(T_{n+1}) - A(T_n)), \end{aligned}$$

where we have denoted

$$\begin{aligned} A(T_{n+1}) - A(T_n) &= \frac{\tilde{V}_{\text{CLE}}(T_n)}{B(T_n)} - \mathbb{E}_{T_n} \left(\frac{\tilde{V}_{\text{CLE}}(T_{n+1})}{B(T_{n+1})} \right) \\ &= \iota(T_n) \left\{ \frac{U_n(T_n)}{B(T_n)} - \frac{\tilde{H}_n(T_n)}{B(T_n)} \right\} \end{aligned}$$

with $A(0) = 0$. The second equality follows from the fact that $\tilde{H}_n(T_n)/B(T_n) = \mathbb{E}_{T_n}(\tilde{V}_{\text{CLE}}(T_{n+1})/B(T_{n+1}))$. Therefore, we have the following decomposition

$$\frac{\tilde{V}_{\text{CLE}}(T_n)}{B(T_n)} = M(T_n) - A(T_n).$$

Notice that the process A is *not* an increasing process (and \tilde{V}_{CLE} therefore not a supermartingale), since we cannot guarantee that $\tilde{H}_n(T_n) \geq U(T_n)$: whenever an incorrect exercise decision is made, A decreases.

For the purpose of computing an upper bound, the hold values $\tilde{H}_n(T_n)$ and $\tilde{H}_{n+1}(T_{n+1})$ in (18.37) *cannot* be estimated by regression; doing so will introduce unknown biases which will destroy the martingale property of M and, in turn, invalidate the inequality in (18.31). Instead, following Andersen and Broadie [2004] we can launch at times T_n and T_{n+1} Monte Carlo simulations to estimate the two expectations in (18.35). Notice that these “inner” Monte Carlo simulations will be nested inside a main “outer” simulation trial that generate sample paths of U , B , and M , as needed to estimate the expectation on the right-hand side of (18.31).

Now, we insert the martingale defined by (18.37) into the right-hand side of (18.31), which gives rise to a high-biased estimate $H_0^{\text{hi}}(0)$ for the CLE value,

$$H_0^{\text{hi}}(0) = \tilde{H}_0(0) + \Delta \geq H_0(0), \quad (18.38)$$

where the *duality gap* Δ is defined as

$$\Delta = E(D), \quad D \triangleq \max_{n=1,\dots,N-1} \left(\frac{U_n(T_n)}{B(T_n)} - M(T_n) \right). \quad (18.39)$$

The hold value $\tilde{H}_0(0) = \tilde{V}_{\text{CLE}}(0)$ can be estimated bias-free from the given exercise strategy $\tilde{\eta}$ by standard Monte Carlo methods (see Section 18.3.6), so we focus on providing an estimate of the duality gap Δ . The following *NS algorithm* can be used for this purpose.

1. Simulate K_U paths $\omega_1, \dots, \omega_{K_U}$.
2. For each path ω_k calculate simulated values of the numeraire $B(T_n, \omega_k)$, $n = 1, \dots, N - 1$.
3. For each path ω_k , calculate the value $U_n(T_n, \omega_k)$ of the underlying exotic swap on all exercise dates $n = 1, \dots, N - 1$.
4. For each path ω_k and each T_n , $n = 1, \dots, N - 2$, launch K_{nest} independent sub-paths $\{\omega_{k,1}^n, \dots, \omega_{k,K_{\text{nest}}}^n\}$ to time T_N and estimate hold values $\tilde{H}_n(T_n, \omega_k)$ as

$$\tilde{H}_n(T_n, \omega_k) \approx \hat{H}_n(T_n, \omega_k) \triangleq \frac{1}{K_{\text{nest}}} \sum_{j=1}^{K_{\text{nest}}} \frac{U_{\tilde{\gamma}(j,k,n)}(T_{\tilde{\gamma}(j,k,n)})}{B(T_{\tilde{\gamma}(j,k,n)})},$$

where, in slightly labored notation, $\tilde{\gamma}(j, k, n)$ is the first exercise date for sub-path j , date T_n , and “outer” path k :

$$\tilde{\gamma}(j, k, n) = \min \{l > n : \iota(T_l, \omega_{k,j}^n) = 1\} \wedge N.$$

5. For each path ω_k , use $\widehat{H}_n(T_n, \omega_k)$, $n = 1, \dots, N - 1$, to form martingale estimates $\widehat{M}(T_n, \omega_k)$ by substituting $\widehat{H}_n(T_n, \omega_k)$ for $\widetilde{H}_n(T_n, \omega_k)$ in (18.37).
6. For each path ω_k , compute pathwise duality gaps as

$$\widehat{D}(\omega_k) = \max_{n=1, \dots, N-1} \left(\frac{U_n(T_n, \omega_k)}{B(T_n, \omega_k)} - \widehat{M}(T_n, \omega_k) \right).$$

7. Estimate Δ as $\widehat{\Delta}$, where

$$\widehat{\Delta} = \frac{1}{K_U} \sum_{k=1}^{K_U} \widehat{D}(\omega_k).$$

18.3.8.3 Bias and Computational Cost of NS Algorithm

Having outlined the basic NS algorithm above, our first order of business is to establish formally that the estimator for $H_0^{\text{hi}}(0)$ resulting from the NS algorithm is, in fact, biased high. Our primary concern here is the effect of the usage in Step 4 of nested Monte Carlo estimators $\widehat{H}_n(T_n, \omega_k)$ in place of the true hold values $H_n(T_n, \omega_k)$. The key result is the following.

Proposition 18.3.1. *In the NS algorithm, the estimator for the duality gap is biased high, i.e.*

$$\mathbb{E} \left(\frac{1}{K_U} \sum_{k=1}^{K_U} \widehat{D}(\omega_k) \right) \geq \mathbb{E}(D).$$

Proof. We drop ω_k throughout the proof, such that $\widehat{H}_n(T_n) = \widehat{H}_n(T_n, \omega_k)$ and so forth. By construction, $\widehat{H}_n(T_n)$ is an unbiased estimator for $\widetilde{H}_n(T_n)$, i.e.

$$\widehat{H}_n(T_n) = \widetilde{H}_n(T_n) + e_n,$$

where e_n is a pure-noise error term with mean 0 and standard deviation proportional to $1/\sqrt{K_{\text{nest}}}$. It follows from (18.37) that Step 5 in the NS algorithm will compute

$$\begin{aligned} & \widehat{M}(T_{n+1}) - \widehat{M}(T_n) \\ &= \iota(T_{n+1}) \frac{U_{n+1}(T_{n+1})}{B(T_{n+1})} \\ &\quad + (1 - \iota(T_{n+1})) \frac{\widetilde{H}_{n+1}(T_{n+1}) + e_{n+1}}{B(T_{n+1})} - \frac{\widetilde{H}_n(T_n) + e_n}{B(T_n)} \\ &= M(T_{n+1}) - M(T_n) + (1 - \iota(T_{n+1})) \frac{e_{n+1}}{B(T_{n+1})} - \frac{e_n}{B(T_n)}. \end{aligned}$$

By induction, it follows that, for $n = 1, \dots, N - 1$,

$$\widehat{M}(T_n) = M(T_n) + q_n$$

where $q_n = q_n(\omega_k)$ is a random variable with zero mean (the explicit expression for q_n is irrelevant) for all n . Hence,

$$\widehat{D} = \max_{n=1,\dots,N-1} \left(\frac{U_n(T_n)}{B(T_n)} - M(T_n) - q_n \right).$$

On the path ω_k , let $\bar{n}(\omega_k)$ be the date index at which $U_n(T_n, \omega_k)/B(T_n, \omega_k) - M(T_n, \omega_k)$ attains its maximum. We can therefore write (again dropping ω_k 's)

$$\begin{aligned} & \mathbb{E} \left(\max_{n=1,\dots,N-1} \left(\frac{U_n(T_n)}{B(T_n)} - M(T_n) - q_n \right) \right) \\ & \geq \mathbb{E} \left(\frac{U_{\bar{n}}(T_{\bar{n}})}{B(T_{\bar{n}})} - M(T_{\bar{n}}) - q_{\bar{n}} \right) \\ & = \mathbb{E} \left(\frac{U_{\bar{n}}(T_{\bar{n}})}{B(T_{\bar{n}})} - M(T_{\bar{n}}) \right) \\ & = \mathbb{E} \left(\max_{n=1,\dots,N-1} \left(\frac{U_n(T_n)}{B(T_n)} - M(T_n) \right) \right), \end{aligned}$$

where the first equality follows from the zero mean of $q_{\bar{n}}$. \square

Proposition 18.3.1 demonstrates that our estimate of the duality gap is biased high for finite values of the sub-path sample size K_{nest} . As such, our finite sample estimate for $H_0^{\text{hi}}(0)$ will itself be biased high and have a mean that is *above* the true CLE value, as desired. By increasing K_{nest} we can reduce the bias originating from the finite sample size of sub-paths, and thereby tighten the upper bound. Of course, even in the limit $K_{\text{nest}} \rightarrow \infty$ we will still produce an estimator that is biased high; the size of this bias will reflect the quality of our exercise strategy choice, and will only vanish in the (unlikely) event that we manage to use precisely the optimal exercise strategy, i.e. when $\tilde{\eta} = \eta$.

While the basic NS algorithm is quite straightforward, the need for nested simulations makes it numerically expensive: the worst-case workload will be proportional to

$$K_{\text{nest}} \cdot K_U \cdot N^2. \quad (18.40)$$

For comparison, with K' simulation trials the lower bound simulation in Section 18.3.6 has a workload proportional to $K' \cdot N$, plus the work required to estimate the exercise rule in a pre-simulation. In many cases the inner simulations of the NS algorithm can be stopped quickly (due to exercise of the CLE), so in practice the dependence on N in (18.40) is often less than quadratic and sometimes close to linear. Finally, K_{nest} can often be set to a number much smaller than K_U without significantly affecting the quality of the upper bound, and even very small values of K_{nest} (e.g., 50–100 or less) may yield informative results. If one additionally takes advantage of the

algorithm refinements discussed in Section 18.3.8.6, it is typically possible to execute an upper bound computation for a long-dated (say, 30 years) CLE in a few minutes. This is still relatively time-consuming, so upper bound computations are often most useful in practice as a way to test the quality of a postulated exercise strategy. We expand on this topic in the next section.

18.3.8.4 Confidence Intervals and Practical Usage

For concreteness, assume that the algorithm in Section 18.3.6 has been used to compute a K' -path lower bound estimate of $V_{\text{CLE}}(0) = H_0(0)$. Let this estimate be denoted \hat{V}_{CLE} , and let its recorded sample standard deviation be \hat{s}_L . Also assume that an independent simulation of the NS algorithm with K_U outer simulation trials has produced an estimate $\hat{\Delta}$ for the duality gap, with sample standard deviation \hat{s}_U . Asymptotically, a $100(1 - \gamma)\%$ confidence interval for the true price $V_{\text{CLE}}(0)$ must be *tighter* than

$$\left[\hat{V}_{\text{CLE}} - u_{\gamma/2} \frac{\hat{s}_L}{\sqrt{K'}}; \hat{V}_{\text{CLE}} + \hat{\Delta} + u_{\gamma/2} \sqrt{\frac{\hat{s}_U^2}{K_U} + \frac{\hat{s}_L^2}{K'}} \right], \quad (18.41)$$

where $\Phi(u_{\gamma/2}) = 1 - \gamma/2$. As already mentioned in Section 3.5.6, the confidence interval is conservative¹¹ because of the low bias in the sample estimator \hat{V}_{CLE} (i.e., $E(\hat{V}_{\text{CLE}}) \leq H_0(0)$) and the high bias in $\hat{V}_{\text{CLE}} + \hat{\Delta}$, which originates in part from the nature of the upper bound, and in part from the earlier mentioned additional high bias introduced by the finite sample size of the inner simulations (see Proposition 18.3.1).

As noted in Section 3.5.6, it is not uncommon that the upper and lower bounds for the option price often are roughly symmetric around the true value, so in the event that we have computed both bounds, the obvious point estimate

$$\hat{V}_{\text{CLE}} + \frac{1}{2} \hat{\Delta} \quad (18.42)$$

will often give better price estimates than either the upper or lower bound alone.

Upper bound simulation algorithms can typically be expected to be both more involved and/or more expensive than lower bound simulation methods. In many cases, the best use of the upper bound simulation algorithm will therefore be to test whether postulated lower bound exercise strategies are tight or not. Specifically, starting from some guess for the exercise strategy, we can produce confidence intervals using (18.41) to test whether the lower

¹¹In addition to random Monte Carlo error, simulation of some models may also involve a systematic error stemming from the time-discretization of the model dynamics. Such discretization errors are not accounted for in (18.41), nor in any previous argument. In most cases, however, the discretization bias will be negligible relative to the random Monte Carlo error.

bound estimate is of good quality, in which case the confidence interval can be made tight by using large values of K_U and K' (as well as the number of inner simulation trials, K_{nest}). In case the lower bound estimator is deemed unsatisfactory, we can iteratively refine it, by altering the choice of basis functions, say, until the confidence interval is tight. Importantly, such tests can often be done at a high level, covering entire classes of payouts and/or models. Once an exercise strategy has been validated for a particular product or model, day-to-day pricing of callable securities can be done by the lower bound method, with only occasional runs of the upper bound method needed (e.g., if market conditions change markedly). If upper bound methods are predominantly used in this fashion, the fact that they may sometimes be computationally intensive¹² becomes less punitive.

18.3.8.5 Non-Analytic Exercise Values

The observant reader might have noticed that the NS algorithm outlined in Section 18.3.8.2 assumes (in Steps 3, 4, and 6) that exercise values are directly computable at each maturity and each state of the world. A similar assumption was made in the basic LS regression algorithm of Section 18.3.1, but relaxed later, in Section 18.3.2. To handle CLEs that involve exotic swap underlyings for which the values are not easily computable, we can modify the upper bound simulation algorithm in at least two ways.

Our first approach is straightforward, and based on the representations

$$\frac{U_n(T_n)}{B(T_n)} = \mathbb{E}_{T_n} \left(\sum_{i=n}^{N-1} \frac{X_i}{B(T_{i+1})} \right),$$

and

$$\begin{aligned} \frac{\tilde{H}_n(T_n)}{B(T_n)} &= \mathbb{E}_{T_n} \left(\frac{U_{\tilde{\eta}_n}(T_{\tilde{\eta}_n})}{B(T_{\tilde{\eta}_n})} \right) = \mathbb{E}_{T_n} \left(\frac{\mathbb{E}_{T_{\tilde{\eta}_n}} \left(\sum_{i=\tilde{\eta}_n}^{N-1} \frac{X_i}{B(T_{i+1})} \right)}{B(T_{\tilde{\eta}_n})} \right) \\ &= \mathbb{E}_{T_n} \left(\frac{\sum_{i=\tilde{\eta}_n}^{N-1} \frac{X_i}{B(T_{i+1})}}{B(T_{\tilde{\eta}_n})} \right), \end{aligned} \quad (18.43)$$

where the last equality follows from the optional sampling theorem. The relevant expectations can be computed bias-free by launching a nested simulation at time T_n , generating sample paths from time T_n to T_{N-1} . For instance, for the path ω_k and date T_n , we would write

¹²Note that in testing the viability of a class of exercise rules through an upper bound simulation, it is often acceptable to work with a reduced set of exercise opportunities — e.g. change a quarterly exercise schedule to an annual one, say — in order to save computation time (see (18.40)).

$$U_n(T_n, \omega_k) \approx \widehat{U}_n(T_n, \omega_k) = \frac{1}{K_{\text{nest}}} \sum_{j=1}^{K_{\text{nest}}} \sum_{i=n}^{N-1} \frac{X_i(\omega_{k,j}^n)}{B(T_{i+1}, \omega_{k,j}^n)},$$

where $\{\omega_{k,1}^n, \dots, \omega_{k,K_{\text{nest}}}^n\}$ is the set of “inner” sub-paths spawned at time T_n for the “outer” path ω_k . An estimator $\widehat{H}_n(T_n, \omega_k)$ for $\widetilde{H}_n(T_n, \omega_k)$ may be computed from (18.43) the same way, except that only net coupons $X_i(\omega_{k,j}^n)$ *after* exercise will be counted in the accumulation of deflated coupons. With these estimators used in Steps 4 and 6, the NS algorithm in Section 18.3.8.2 may proceed as before. As the nested simulation will produce bias-free (but noisy) estimates for the true exercise and hold values, it follows from Proposition 18.3.1 that the resulting upper bound estimator will still be guaranteed to be biased high.

While using nested simulation in the manner described above does not add to the order of the computational complexity of the upper bound algorithm (it remains as in (18.40)), the need to construct exercise values by simulating coupon streams will obviously add additional noise to the basic algorithm, which in turn will increase the finite sample bias and also widen the confidence interval (18.41) for a given computational budget.

In our second approach to dealing with non-analytic exercise values, we follow the alternative route also taken in Section 18.3.3, and focus on pricing the cancelable note $G_0(0)$, where

$$G_0(0) = H_0(0) - U_0(0).$$

As a starting point, consider the expression (18.21), which we may rewrite as

$$\begin{aligned} G_0(0) &= \sup_{\xi \in \mathcal{T}_0} E \left(- \sum_{n=0}^{\xi-1} B(T_{n+1})^{-1} X_n \right) \\ &= \sup_{\xi \in \mathcal{T}_0} E(B(T_\xi)^{-1} J(T_\xi)), \quad J(T_\xi) \triangleq - \sum_{n=0}^{\xi-1} \frac{B(T_\xi)}{B(T_{n+1})} X_n. \end{aligned} \quad (18.44)$$

Financially, the quantity $J(T_\xi)$ can be interpreted as the payout at the exercise date T_ξ from re-investing all pre-exercise coupons into the numeraire B . Effectively, this formulation removes all pre-exercise cash flows from the expectation for $G_0(0)$, making (18.44) structurally identical to (18.30), but with the exercise value $U_\xi(T_\xi)$ replaced by $J(T_\xi)$. Of course, while $U_\xi(T_\xi)$ may be difficult to compute at time T_ξ in a Monte Carlo simulation, $J(T_\xi)$ is not.

To construct an upper bound for $G_0(0)$, we follow the same principles that lead to (18.35) and construct a martingale M as

$$M(T_{n+1}) - M(T_n) = E_{T_{n+1}} \left(\frac{J(T_{\tilde{\eta}_n})}{B(T_{\tilde{\eta}_n})} \right) - E_{T_n} \left(\frac{J(T_{\tilde{\eta}_n})}{B(T_{\tilde{\eta}_n})} \right), \quad (18.45)$$

where $\tilde{\eta}$ is a given exercise strategy. With $\tilde{G}_0(0)$ being the lower bound value for $G_0(0)$ computed by using $\tilde{\eta}$ as the exercise strategy, an upper bound for $G_0(0)$ is (compare to (18.38))

$$G_0^{\text{hi}}(0) = \tilde{G}_0(0) + \Delta_G \geq G_0(0),$$

where the duality gap Δ_G is defined as

$$\Delta_G = E(D_G), \quad D_G \triangleq \max_{n=1,\dots,N-1} \left(\frac{J(T_n)}{B(T_n)} - M(T_n) \right).$$

When using Monte Carlo simulation to estimate Δ_G , we can use nested simulation to establish a bias-free estimator for the martingale in (18.45), in the same way as was done for the NS algorithm. By the arguments in Proposition 18.3.1, the resulting estimator for Δ_G will be biased high, i.e. our upper bound is valid. Confidence intervals for $G_0(0)$ (and for $H_0(0) = G_0(0) + U_0(0)$) can be constructed using the principles of Section 18.3.8.4.

18.3.8.6 Improvements to NS Algorithm

In Broadie and Cao [2008], the authors outline a number of improvements to both upper and lower bound simulations. As some of the proposed techniques are fairly involved, we cannot give this paper full justice, but contend ourselves with listing one relatively straightforward trick from Broadie and Cao [2008]. The reader interested in additional techniques should consult Broadie and Cao [2008] directly.

Let us return to the setting of Section 18.3.8.2, and assume that an easily computable lower limit $\underline{H}(T_n)$ exists for the hold value at time T_n ,

$$H_n(T_n) \geq \underline{H}(T_n), \quad n = 1, \dots, N-1. \quad (18.46)$$

We comment on how to choose $\underline{H}(\cdot)$ later. Let us also assume that the inequality (18.46) is honored by the given exercise policy $\tilde{\eta}$, i.e. we assume that exercise will never take place when $U_n(T_n) < \underline{H}(T_n)$. When using a regression approach, we can ensure that this assumption is true by simply writing (compare to (18.36))

$$\tilde{\eta}_n = \min \left\{ k \geq n+1 : U_k(T_k) \geq \max \left(\underline{H}(T_k), C \left(\tilde{H}_k(T_k) \right)^T \zeta(T_k) \right) \right\} \wedge N. \quad (18.47)$$

Modifying an exercise policy to accommodate bounds such as (18.46) is sometimes known as *policy fixing*, see Broadie and Glasserman [2004].

Before stating our next result, let us make the additional assumption that $\underline{H}(t)/B(t)$ is a submartingale in measure Q^B for $t \leq T_N$; this is, for instance, the case when $\underline{H}(t)$ is chosen to be either zero or the price of an asset that pays no cash flows before T_N . With this assumption, we have

$$\frac{\tilde{H}_n(T_n)}{B(T_n)} = \mathbb{E}_{T_n} \left(\frac{U_{\tilde{\eta}_n}(T_{\tilde{\eta}_n})}{B(T_{\tilde{\eta}_n})} \right) \geq \mathbb{E}_{T_n} \left(\frac{\underline{H}(T_{\tilde{\eta}_n})}{B(T_{\tilde{\eta}_n})} \right) \geq \frac{\underline{H}(T_n)}{B(T_n)}. \quad (18.48)$$

We use this result to show the following proposition, adapted from Broadie and Cao [2008].

Proposition 18.3.2. *Let the martingale $M(t)$ be defined as in (18.35) and (18.37), and assume that $\underline{H}(t)/B(t)$ is a submartingale. Assume also that $U_k(T_k) < \underline{H}(T_k)$ for all $k = l, \dots, n$, with $1 \leq l \leq n$. Then we have,*

$$M(T_n) = M(T_{l-1}) - \tilde{H}_{l-1}(T_{l-1})/B(T_{l-1}) + \tilde{H}_n(T_n)/B(T_n)$$

and

$$U_n(T_n)/B(T_n) - M(T_n) < \tilde{H}_{l-1}(T_{l-1})/B(T_{l-1}) - M(T_{l-1}). \quad (18.49)$$

In particular, if $l = 1$ then

$$M(T_n) = \tilde{H}_n(T_n)/B(T_n), \quad U_n(T_n)/B(T_n) < M(T_n). \quad (18.50)$$

Proof. First, we notice that exercise will never take place on the interval $[T_l, T_n]$ such that, from (18.37),

$$\begin{aligned} M(T_n) &= M(T_{n-1}) + \frac{\tilde{H}_n(T_n)}{B(T_n)} - \frac{\tilde{H}_{n-1}(T_{n-1})}{B(T_{n-1})} \\ &= M(T_{l-1}) + \sum_{j=l}^n \left(\frac{\tilde{H}_j(T_j)}{B(T_j)} - \frac{\tilde{H}_{j-1}(T_{j-1})}{B(T_{j-1})} \right) \\ &= M(T_{l-1}) - \frac{\tilde{H}_{l-1}(T_{l-1})}{B(T_{l-1})} + \frac{\tilde{H}_n(T_n)}{B(T_n)}. \end{aligned}$$

Using (18.48), we also have

$$\frac{\tilde{H}_n(T_n)}{B(T_n)} \geq \frac{\underline{H}(T_n)}{B(T_n)} > \frac{U_n(T_n)}{B(T_n)}, \quad (18.51)$$

which proves (18.49). The results for the special case where $l = 1$ follows from the fact that $M(0) = \tilde{H}_0(T_0)$. \square

Recalling (18.39), it is clear from (18.50) that there is no contribution to the upper bound increment D before the underlying process enters the region where $U_n(T_n) \geq \underline{H}(T_n)$. Moreover, whenever the option is inside a region where the exercise value is less than the lower limit \underline{H} , M does not depend on the actual path taken inside the region. This allows for the following straightforward modification to the NS algorithm: whenever at time T_n we observe that $U_n(T_n) < \underline{H}(T_n)$, we simply skip the nested simulations and proceed to the next date T_{n+1} ; otherwise we launch K_{nest} sub-simulations

and update $M(T_n)$ according to Proposition 18.3.2. When computing \widehat{D} in Step 6, we can ignore all dates where $U_n(T_n) < \underline{H}(T_n)$. For options that are deeply out-of-the money, a substantial number of nested simulations can be avoided in this fashion, leading to significant improvements in computational performance.

The refinement in the Broadie and Cao [2008] algorithm hinges upon our ability to establish a meaningful lower limit submartingale \underline{H} . While this must generally be done on a case-by-case basis, one obvious possibility in the standard CLE setting is to set the lower limit to zero. This is a relatively coarse bound, and it may be tempting to sharpen it by observing that, for any $t \geq 0$ and $n, m \geq 0$,

$$H_n(t) \geq H_{n+m}(t) \geq U_{n+m+1}(t),$$

which, assuming we can calculate $U_{n+m+1}(t)$'s analytically (admittedly a strong assumption), shows that a lower limit can be computed as

$$\underline{H}(T_n) = \max_{k=n+1, \dots, N-1} (U_k(T_n)).$$

However, while this choice of the lower limit¹³ $\underline{H}(\cdot)$ can certainly be used in policy fixing to improve the lower bound, it is generally not a submartingale and therefore *cannot* be used in upper bound computations.

18.3.8.7 Other Upper Bound Algorithms

To avoid the need for nested simulation, it is tempting to return to the representation in (18.34) and contemplate whether one can estimate the optimal choice of $\sigma(t)$ by extracting it from the empirical volatility of (the martingale component of) $\tilde{V}_{\text{CLE}}(t)/B(t)$. One advantage to this approach is that any errors in the estimation of $\sigma(t)$ will not affect the martingale property of $M(t)$ in (18.34). Starting again from a postulated exercise strategy $\tilde{\eta}$, Belomestny et al. [2007] use this observation to construct a regression on a set of basis functions to produce an estimate for the function $\sigma(t)$. By applying regression techniques this way, the authors are able to construct a true martingale process $M(t)$, which can be turned into a valid upper bound through (18.31). While the resulting algorithm involves no nested simulation, it requires considerable care in its implementation, in part because the optimal integrand $\sigma(t)$ can be expected to be considerably less regular than the optimal martingale $M(t)$ itself. This, in turn, requires additional thought in the selection of appropriate basis functions for the regression. One possibility advocated in Belomestny et al. [2007] is to include, whenever available, exact or approximate expressions for the diffusion term

¹³This limit is essentially equivalent to imposing a so-called “carry” restriction on exercise, a notion that we explore in some detail in Section 18.3.10.2 and Proposition 19.7.1.

in dynamics of several still-alive European options underlying the Bermudan option. This strategy is akin to that of Rogers [2001], and its feasibility depends on the pricing problem at hand. In cases where it does apply, the authors of Belomestny et al. [2007] demonstrate that their method gives good results, with the upper bound often being nearly as tight as that of the nested algorithm in Andersen and Broadie [2004]. They also show how to use their technique to develop a variance-reduced version of the algorithm in Andersen and Broadie [2004].

Additional techniques for computing upper bounds can be found in Broadie and Glasserman [1997] and Glasserman and Yu [2002]. The latter is based purely on regression, but requires strong conditions on regression basis functions, that may be hard to check in practice.

18.3.9 Regression Variable Choice

The single most critical determinant of the performance of the regression-based valuation methods is the choice of the regression variables¹⁴ $\zeta(t) = (\zeta_1(t), \dots, \zeta_q(t))^\top$. Recall that the values of these variables at time T partly¹⁵ serve to approximate the information contained in the sigma-algebra \mathcal{F}_T . Some of this information is relevant to the valuation of a given security and some is not. The closer $\zeta(T)$ approximates the information in \mathcal{F}_T that is relevant for the security, the better the regression method performs, i.e. the smaller is the bias of lower bound estimates of the security value.

18.3.9.1 State Variables Approach

For Markovian models, all information observed at (but not before) time T is encoded in the Markovian state variables, say some (d -dimensional) vector $x(T)$. As discussed earlier in Section 3.5.4, for such models it is often natural take the $\zeta(T)$'s to be deterministic functions of the state variables; these functions should approximately span the set of all functions of the state variables. A good choice is to use monomials of $x(T)$, i.e. functions of the type $\prod_{i=1}^d x_i(T)^{p_i}$, which corresponds to using a polynomial basis in the LS regression. One caveat applies: the filtration $\sigma(x(T))$ generated by the state variables at $x(T)$ is clearly smaller than \mathcal{F}_T , since the former consists of model information observed at time T only, while the latter contains information for *all* times from 0 to T . For some derivatives the reduction from \mathcal{F}_T to $\sigma(x(T))$ is irrelevant — this includes all callable securities whose coupons are not path-dependent. On the other hand, for securities such

¹⁴While for convenience we write ζ as a function of time t , the regression variables only need to be defined at times when we perform regression, i.e. at exercise dates T_1, \dots, T_{N-1} .

¹⁵Beyond representing information available at time T , $\zeta(T)$ also serves to define the function space that is obtainable by least-squares regression.

as callable snowballs (Section 5.14.4) whose coupons depend on the whole history of interest rate evolution, the regression variables must be augmented with some that carry information from the past. To do this in a product-independent way, we can, for instance, include state variables from previous times in $\zeta(T)$.

Markovian interest rate models that would allow for such essentially product-independent choice of regression variables include quasi-Gaussian models and quadratic Gaussian models, at least if the number of Markov state variables is not too high. While Libor market (LM) models are Markovian in the set of all Libor forward rates (and any stochastic volatility variables), the dimensionality of the state vector is typically so large that the regression will suffer from numerical problems, an issue we shall discuss further later. One general way to reduce the dimension of the state variable vector is to perform a principal components analysis. To demonstrate this idea in a LM model setting, suppose that we wish to synthesize d “state variables” by forming the first d principal components of the vector of Libor rates. At a given time T_n we have a vector of (centered) still-alive forward Libor rates

$$A = (L_n(T_n) - \mathbb{E}(L_n(T_n)), \dots, L_{N-1}(T_n) - \mathbb{E}(L_{N-1}(T_n)))^\top.$$

We assume $N - n \geq d$. The term covariance matrix of the rates,

$$c = \mathbb{E}(AA^\top),$$

can be estimated using methods from Chapter 14. Then, by principal components analysis (see Section 3.1.3) we can find an $(N - n) \times d$ matrix D such that DD^\top is the closest (in the Frobenius norm) rank- d approximation to c . Then

$$A \approx Dx$$

for some d -dimensional vector $x = x(T_n)$, which we recover from the least-squares (regression) problem

$$\|A - Dx\|^2 \rightarrow \min,$$

with the solution

$$x = (D^\top D)^{-1} D^\top A.$$

Doing this for each T_n gives us a set of d approximate state variables for regression on each exercise date.

18.3.9.2 Explanatory Variables

While the state variable approach (with or without principal components dimension reduction) is appealing in its independence of security specifics, it has several shortcomings. First, as mentioned earlier, situations may arise where the information carried in the state variables is inadequate for the

security in question (e.g. snowballs or other path-dependent CLEs). Second, there may be cases where there is *too much* information in the state variables, leading to an overabundance of regression variables. For instance, a standard Bermudan swaption turns out (see Section 19.6.1) to primarily depend on the overall level of interest rates, so if, say, the number d of PCA state variables used in an LM model is much larger than 1, many of the regression variables will have little or no explanatory power.

Having too many regression variables not only adds more work to the numerical scheme, it ultimately tends to reduce the quality of the CLE price estimate. Indeed, for a finite budget of simulated paths used in regression, having too many regression variables¹⁶ will induce errors in the regression coefficient estimates and, therefore, in the exercise rule and, ultimately, in the value of the callable security. Moreover, if some of the regression variables add no explanatory power to the regression, their inclusion in the regression may lead to spurious noise and further issues with the estimation quality of the exercise rule¹⁷.

In light of the comments above, in practice it is hard to avoid a careful analysis of each security type, to ensure that neither too little or too much information is contained within the set of regression variables. In such a “product-specific” approach to choosing regression variables, we would aim to choose regression variables in a way that maximizes their explanatory power for the particular security in question.

Similar to the state variables approach, we find it convenient to specify regression variables $\zeta(t)$ as simple functions of so-called *explanatory variables*, which can be thought of as product-specific analogs to state variables of the model. Let us use $x(t) = (x_1(t), \dots, x_d(t))^\top$ to denote these explanatory variables; the fact that we recycle the notation for x and d from Section 18.3.9.1 should not lead to any confusion. When constructing regression variables from explanatory variables, we typically choose monomials of explanatory variables

$$\prod_{i=1}^d x_i(t)^{p_i}, \quad (18.52)$$

up to a low order r , say 3 to 4, so that $\sum_{i=1}^d p_i \leq r$. This, of course, amounts to fitting exercise and hold values with polynomials of degree r in explanatory variables. It is worth noting that while higher-order polynomials could give a closer fit, they could also lead to unexpected behavior outside of the range of the points used in fitting. As exercise boundaries may lie far away

¹⁶Glasserman and Yu [2004] determine that the number of paths should grow exponentially in the number of regression variables. Under different assumptions, Moni [2005] shows that the number of paths should grow polynomially. Polynomial rate of growth is also derived in Egloff et al. [2007].

¹⁷Section 18.3.10 lists a number of regularization approaches that may help guard against some of the issues associated with spurious regression variables.

from typical simulation scenarios, this could lead to poor estimation of the exercise rule.

When selecting explanatory variables, we should look for variables that have high explanatory power in the regression of hold and exercise values of a given CLE. Inevitably, this process is trade-type specific and combines elements of both science and art. While trial and error is always needed, we can make a few recommendations. Generally, we always prefer using financially meaningful explanatory variables such as various market rates or values of (simple) market instruments, and find it convenient to distinguish three classes of potential explanatory variables. The first class contains the market variables that directly drive the values of coupons. For example, if a coupon is linked to a CMS spread between two rates, we should include the spread as a potential explanatory variable. The second class contains those variables that describe the past; these variables are only relevant for CLEs with path-dependent coupons. For instance, recall that in a snowball the value of a current coupon depends on the values of the previous coupons, so we would recommend including the current coupon in the set of potential explanatory variables. Finally, the last class contains those market variables that are thought to drive the exercise decision. Inspired by the case of a simple Bermudan swaption, we would often include a variable responsible for the overall level of the yield curve. In particular, for the regression on date T_n we would typically include a swap rate that fixes at T_n and spans all periods to the last exercise (a so-called “core” swap rate). We might also include a variable that reflects the slope of the yield curve on the exercise date; a front Libor rate, i.e. a Libor rate that fixes on the exercise date, is a good choice for this.

After collecting potential explanatory variables from all three classes, we would typically proceed to analyze the list and try to reduce it to a manageable number of variables, such as 2 to 4, that we would then use in our regression method. While nothing replaces careful analysis, the selection process can be automated somewhat and done during the actual regression — we discuss this in more details later on, in Section 18.3.10.1.

Whatever variables we choose, we should be careful to always choose variables for regression on T_n that are \mathcal{F}_{T_n} -measurable. In plain speak, variables should not be “future-looking”, but should be computable by using only the state of the model as observed up to and at time T_n . This requirement, while seemingly technical, is critical to the success of the regression algorithm, as one must not be allowed to “see into the future” when making decisions about exercise.

In selecting explanatory variables we should be thinking about qualitative impact of various variables on exercise/hold values but, fortunately, we do not need to be quantitatively exact in capturing the effects. As long as the general influence is accounted for correctly, the fitting of parametric functions will generally take care of choosing the best scalings and/or linear

combinations of explanatory variables¹⁸. For example, one may decide (as we often recommend) to include in the set of explanatory variables the level of the yield curve, as measured by a core swap rate, as well as the slope of the yield curve, as measured by the difference between the swap rate and some short-tenor Libor rate. For the latter variable, one does not need to explicitly include the difference of the rates as an explanatory variable, since the short-tenor Libor rate will work just as well — the properly weighted difference of the two rates will implicitly be used in the fitting of the polynomials.

18.3.9.3 Explanatory Variables with Convexity

Sometimes the quality of exercise boundary estimation is improved by such seemingly minor changes as using core swap *values* in the set of explanatory variables, rather than core swap rates. In practical applications, we always recommend trying both and seeing which one is better. The reasons for this effect are not always entirely clear, but originate with subtle convexity differences between the two choices. It turns out that convexity, i.e. non-linear dependence between exercise values and the explanatory variables, can be quite important. In particular, instead of using simple rates such as a core swap rate or a front Libor rate, we sometimes find it useful to use *functions* of them as explanatory variables. While this may seem superfluous — after all we will ultimately be applying polynomial functions to explanatory variables to get our regression variables — it turns out that using functional mappings more finely tuned to the features of the trade is often beneficial. The effect is due to the fact that for many CLEs the underlying swap consists of coupons that are options on rates. By matching, roughly, the resulting convexity in value of coupons relative to underlying rates, we can often improve the regression fit. Sometimes this can be accomplished by simply using $(S(T_n) - K)^+$, instead of some rate $S(T_n)$, as an explanatory variable for a CLE whose coupons are strike- K options on the rate $S(T_n)$ (or some other, relatively similar, rate). For a more refined approach we could try to roughly estimate the option value of the remaining coupons as of the exercise date T_n . Even if coupons do not have any optionality, such as for Bermudan swaptions, we may wish to use European options on the core swap rates as explanatory variables to better fit *hold* values, which are always convex due to callability.

If one uses a model that allows for closed-form valuation of European options on rates at any future time, then, for suitably simple CLEs, it may be possible to use the exact exercise value as an explanatory variable. This, however, only works for a subset of models, and the approach can significantly impair the speed of valuation, since one needs to compute European option prices a large number of times (one per path per exercise time per underlying

¹⁸As discussed in Section 18.3.10, a simple pre-normalization of variables before the regression may still be useful in preventing numerical problems.

option). Fortunately, it is not necessary to be particularly precise in matching the curvature of the exercise values with the explanatory variables, and a rough estimate will typically do just fine. For example one can use the Black formula with an approximate volatility to value the options in the underlying, even in non-Black models. Also, if the underlying is a strip of options, one does not need to value all options in the underlying, but just one, e.g. the one in the middle of the strip.

To give an example of an implementation using option-like explanatory variables, consider a callable capped floater. The structured coupon at time T_n is here given by

$$C_n = \min(L_n(T_n) + s, c),$$

received against a Libor rate payment,

$$X_n = \tau_n \times (C_n - L_n(T_n)).$$

We note that the net coupon X_n can be written

$$\begin{aligned} X_n &= \tau_n \times (\min(L_n(T_n) + s, c) - L_n(T_n)) \\ &= \tau_n \times s - \tau_n \times (L_n(T_n) - (c - s))^+, \end{aligned}$$

so the payment is a combination of a fixed rate payment with rate s and a call option on the Libor rate with strike $c - s$.

For the first explanatory variable we may use an approximate value of the exercise value on each exercise date

$$x_1(T_n) = \sum_{i=n}^{N-1} \tau_i \times P(T_n, T_{i+1}) \times (s - c_B(T_n, L_i(T_n); T_i, c - s; \sigma_{n,i})),$$

where the Black formula $c_B(\dots)$ is defined in Remark 7.2.8. The volatility $\sigma_{n,i}$ is the log-normal volatility of the Libor rate $L_i(t)$ over the interval $[T_n, T_i]$ as given by the model. To reiterate, very crude approximations can be used here. For the displaced log-normal Libor market model (see (14.12)) one may use $\sigma_{n,i} = \bar{\lambda}$, where $\bar{\lambda}$ is some average of relative forward Libor model volatilities over time and Libor rates, and for the stochastic volatility version (see (14.16)) one may use

$$\sigma_{n,i} = \bar{\lambda} \sqrt{z(T_n)}. \quad (18.53)$$

For the second explanatory variable, we can use a front Libor rate fixing on an exercise date, as we recommended before. Alternatively, we can use just the approximate value of the first coupon payment,

$$x_2(T_n) = \tau_n \times P(T_n, T_{n+1}) \times (s - c_B(T_n, L_n(T_n), c - s; T_n; \sigma_{n,n}))$$

or the last one,

$$x_2(T_n) = \tau_n \times P(T_n, T_N) \times (s - c_B(T_n, L_{N-1}(T_n); T_{N-1}, c - s; \sigma_{n, N-1})).$$

In conjunction with the first explanatory variable, either choice of $x_2(T_n)$ will capture the effect of changes in the interest rate curve slope.

In the example above, we proposed using an explanatory variable dependent on the stochastic volatility process (see (18.53)). It turns out that this is quite important to do in general, even for CLEs that do not have optionality in the underlying, such as Bermudan swaptions. The stochastic volatility process $z(t)$ can either be incorporated into an explanatory variable in the way of (18.53), or used as a separate explanatory variable.

18.3.10 Regression Implementation

While the selection of the regression variables is the primary key to the success of a regression-based method for CLE valuation, the details of implementation of the numerical algorithm for performing regressions are important as well. The basic regression algorithm does not involve much more than matrix inversion, see Section 3.5.4, but there are a number of ways in which the algorithm could be made more robust. We discuss them in this section.

For future reference, let us quickly fix notations. Throughout, we shall consider a particular date T_n and assume that paths $\omega_1, \dots, \omega_K$ have been simulated, allowing us to collect a $K \times q$ matrix Z of simulated values of the regression variables $\zeta(T_n)$, with

$$Z_{k,j} = \zeta_j(T_n, \omega_k), \quad k = 1, \dots, K, \quad j = 1, \dots, q.$$

Recall that each $\zeta_j(T_n)$ is typically obtained as a particular monomial (18.52) applied to the vector of explanatory variables (or state variables, per Section 18.3.9.1) observed at time T_n . We also assume that we have available a K -dimensional vector of simulated values $Y = (Y_1, \dots, Y_K)^\top$ that we would like to regress, for example the vector $Y_k = V_n(\omega_k)$, $k = 1, \dots, K$, from (18.12) in the basic regression scheme of Section 18.3.1. The goal is to find a q -dimensional vector β such that $Z\beta$ approximates Y in some sense, e.g. in the least-squares sense

$$\|Y - Z\beta\|^2 \rightarrow \min. \tag{18.54}$$

A (naive) solution to this problem is, of course, given by

$$\beta = (Z^\top Z)^{-1} Z^\top Y. \tag{18.55}$$

18.3.10.1 Automated Explanatory Variable Selection

As we explained in Section 18.3.9.2, we often find ourselves in a situation where we have many potential candidates for the explanatory variables but want to prune the set to keep only the most relevant variables for regression.

One approach here is to analyze potential candidates for explanatory variables and, based on experience, choose the ones that subjectively appear to be the most relevant. Another approach is to try to extract a subset of variables based on numerical criteria of regression fit. It turns out that the latter approach is quite common in econometrics circles, allowing us to draw on known techniques. In the problem of *automatic econometric model selection*, one considers a given time series of data — say investment returns of a hedge fund — and tries to choose macro variables, such as equity returns, oil prices, etc., that contribute the most to explaining the time series data. Conceptually, the solution to the problem is simple: one tries regressing the time series on subsets of potential regression variables and observes which subset of variables provides the best fit. The details of how these trials are conducted are, however, quite important, as a brute-force test of all variable combinations would often be impractical: even with just 10 variables to choose from, the number of potential subsets of variables to check is $2^{10} - 1$. The literature on the subject of automatic econometric model selection is quite extensive¹⁹ so we do not go into much detail here, but merely scratch the surface by demonstrating a simple algorithm.

In an algorithm for automatic pruning of explanatory variables, the first question that needs to be answered is how to measure the quality of fit of a given set of regression variables. In the context of linear regression (18.54), an often-used measure²⁰ is R^2 (“R-squared”) which measures the variance of the residual $Y - Z\beta$ relative to the variance of Y

$$R^2 = 1 - \frac{\|Y - Z\beta\|^2}{\|Y\|^2}.$$

With the solution (18.55), R^2 is equal to

$$R^2 = \frac{Y^\top Z\beta}{Y^\top Y}.$$

The higher R^2 , the better the fit is deemed to be. So, in our algorithm, we might choose different subsets of explanatory variables and examine the resulting values of the regression R^2 . The results could, say, be used to choose a subset of variables that would give us the highest R^2 , subject to a constraint on the maximum size d of the variable vector.

Many search paths through the collection of subsets of explanatory variables are possible. For demonstration, let us consider a special case of the *General-to-Specific* (GETS) approach²¹. Here we start with all potential

¹⁹The interested reader could start with Campos et al. [2005].

²⁰A potentially more accurate measure is “modified R-squared”, see Campos et al. [2005]. For our applications where the dimension of the dataset (K) is very large, modified R-squared and ordinary R-squared are almost identical, however.

²¹See e.g. www.pcgive.com/pcgets.

explanatory variables included in the regression and calculate the baseline R^2 . Then we remove each explanatory variable in turn (always returning previously-removed variables to the set after each regression) and calculate the R^2 . Having tried removing all variables, we then select the one that gave us the smallest reduction in R^2 and throw it away. With the potential set of explanatory variables reduced by one, we repeat the procedure and again exclude a variable that gives us the smallest effect on the R^2 . We continue until we have either reached a pre-determined number of explanatory variables, or until we have reached some maximum allowed reduction in R^2 . With each regression typically being pretty quick to execute, this approach does not affect the overall valuation time much, yet generally gives good results. As one can imagine, however, substantially more sophisticated approaches for variable pruning are possible; we refer the reader to Campos et al. [2005] as a starting point.

While we motivated our discussion of explanatory variable selection with a problem in econometrics, we note that our regression problem is not exactly the same as that typically faced in econometric analysis. While econometricians often try to explain virtually all changes in their time series through changes in explanatory variables — i.e. they seek values of R^2 that are close to one — the problems we are interested in here would typically be characterized by values of R^2 that are much lower than one. This is a consequence of the fact that even the full information set available at time T cannot explain all changes in hold/exercise values after time T . In fact, if in the regression we obtain high values of R^2 , this would indicate that there is something wrong with our choice of explanatory variables and they most likely are “future-looking”.

18.3.10.2 Suboptimal Point Exclusion

The main point of performing a regression on exercise/hold values is to determine an exercise rule. If for a given path, at a given point in time, we can prove that the exercise can never be optimal then, arguably, we should exclude this point from the regression that defines the exercise rule.

Interestingly, it turns out that in some general cases we can indeed establish situations where exercise is never optimal. Consider a cancelable note on a coupon stream $-X_1, \dots, -X_{N-1}$, see Section 18.3.3. From (18.18),

$$\begin{aligned} G_{n-1}(T_{n-1}) &= E_{T_{n-1}} \left(B(T_{n-1})B(T_n)^{-1} \left(-X_{n-1} + (G_n(T_n))^+ \right) \right) \\ &\geq -E_{T_{n-1}} (B(T_{n-1})B(T_n)^{-1} X_{n-1}) \end{aligned} \quad (18.56)$$

so if, for a given path ω , $-(E_{T_{n-1}}(B(T_{n-1})B(T_n)^{-1} X_{n-1}))(\omega)$ is positive, then $G_{n-1}(T_{n-1}, \omega)$ is positive and it cannot be optimal to cancel the note at T_{n-1} on that path²². Fortunately, in many cases, X_{n-1} is $\mathcal{F}_{T_{n-1}}$ -measurable and

²²A more general result is derived in Section 19.7.2.

$$\mathbb{E}_{T_{n-1}}(B(T_{n-1})B(T_n)^{-1}X_{n-1}) = P(T_{n-1}, T_n)X_{n-1},$$

a quantity easily computable at time T_{n-1} . Since we know the scenarios where it is never optimal to exercise, we can do two things. First, as a simple application of the policy fixing rule (18.47), we can outright forbid exercise in the suboptimal scenarios, even if our regression-based rule instructs us to exercise. Second, as suggested by Beveridge and Joshi [2008], we can exclude those paths ω_k for which $X_{n-1}(\omega_k)$ is negative from participating in the regression fit at time T_{n-1} , potentially improving the fit in the region that matters for the exercise rule.

The idea of excluding “uninteresting” paths could be taken further, albeit based on more practical than theoretical considerations. For example, we may decide to exclude (or at least de-emphasize) paths that are very deeply in or out of the money, since we want most precision around the exercise frontier itself. We formalize this idea in the next section.

18.3.10.3 Two Step Regression

As mentioned in our discussion on convexity in Section 18.3.9.3, simple polynomial functions may not capture to sufficient precision the functional dependence of regressed hold and exercise variables on regression variables. The functional mapping of Section 18.3.9.3 gives us one way of addressing this shortcoming. Another idea is to fit polynomials separately in different regions of values of explanatory variables. For example, with a single explanatory variable, we can split the range of explanatory variable values into a few intervals, and fit polynomials separately in each interval. With many explanatory variables, simple subdivision of the space into intervals quickly becomes impractical. In this case, however, Beveridge and Joshi [2009] suggest a somewhat similar idea of space subdivision based on the moneyness of the derivative in question. They point out that global regression fit of polynomials often works well for values of the regressed variables “in the wings”, while for points near the “at-the-money” of the CLE (whatever that might mean), the richer functional structure of the hold and exercise values is often not well-approximated by globally fit polynomials. Beveridge and Joshi [2009] propose a two-step scheme, that we describe in the context of cancelable notes of Section 18.3.3 and, in particular, for the scheme (18.25)–(18.26). First, for a given T_n , we regress the variables $Y_k = \widehat{G}_n(T_n, \omega_k)$ to obtain values $\zeta(T_n, \omega_k)^\top \beta$, an approximation to $\widetilde{G}_n(T_n, \omega_k)$, $k = 1, \dots, K$. Subsequently we can choose those ω_k for which $\widetilde{G}_n(T_n, \omega_k)$ is close to zero, which we consider the definition of “at-the-money” (ATM) for a cancelable note. Specifically, for some threshold value ϵ we set

$$\mathcal{W}^\epsilon = \{\omega_k : |\zeta(T_n, \omega_k)^\top \beta| \leq \epsilon\}.$$

Excluding paths outside of \mathcal{W}^ϵ , we now perform a separate fit of $\widehat{G}_n(T_n, \omega)$, $\omega \in \mathcal{W}^\epsilon$, with the same polynomials; since the set of regression points is

different, a new regression coefficient vector β^ϵ is obtained. Finally, we set

$$\tilde{G}_n(T_n, \omega_k) = \begin{cases} \zeta(T_n, \omega_k)^\top \beta, & \omega_k \notin \mathcal{W}^\epsilon, \\ \zeta(T_n, \omega_k)^\top \beta^\epsilon, & \omega_k \in \mathcal{W}^\epsilon, \end{cases}$$

i.e. we use the values of the original regression for non-ATM points $\omega_k \notin \mathcal{W}^\epsilon$, and a new regression for ATM points $\omega_k \in \mathcal{W}^\epsilon$ only. The value of ϵ could be set to be some (small) fraction of the notional of the derivative, say 5%; ultimately some experimentation here may be required.

Rather than a binary division of regression points into “near” and “far” groups, one could also imagine using some kind of smooth kernel that weighs points in the regression differently, depending on how close one is perceived to be to the exercise point, e.g. by how far away from zero $\tilde{G}_n(T_n, \omega_k)$ is. One can also refine it through iterations, where the procedure is repeated multiple times.

18.3.10.4 Robust Implementation of Regression Algorithm

In Section 18.3.1 we assumed that the regression operator \mathcal{R}_T at time $T = 0$ reduces to averaging of the values of the random variable it is applied to, because the values of the regression variables $\zeta(0)$ are the same for all paths. This is indeed the case for the true solution of the problem (18.54). If we assume that all the rows of Z are the same and equal to a vector ζ^\top , then the objective function in (18.54) is equal to

$$\sum_{k=1}^K (Y_k - \zeta^\top \beta)^2.$$

Differentiating with respect to β_j and setting the derivative to zero we obtain

$$\sum_{k=1}^K (Y_k - \zeta^\top \beta) \zeta_j = 0,$$

giving us a regressed solution

$$\zeta^\top \beta = \frac{1}{K} \sum_{k=1}^K Y_k,$$

i.e. the average of Y_k 's as advertised. The solution, however, is *not* given by (18.55) as the inverse of $Z^\top Z = K\zeta\zeta^\top$ here will not exist²³.

This simple example highlights the danger of using the textbook regression solution (18.55) to the problem (18.54). Beyond the degenerate case

²³On a related note, Rasmussen [2005] suggests starting simulation in the past to ensure sufficient variability in the state space for small times, and shows that this step improves estimation of the exercise boundary.

described above, similar issues will arise whenever the matrix $Z^\top Z$ is ill-conditioned, i.e. close to singular. This can happen either due to outright user error (e.g. an inexperienced user accidentally entering the same explanatory variable twice) or due to subtle near-linear dependencies between explanatory variables. In such cases, the regression problem becomes ill-posed and the numerical solution of the regression problem will be unstable. To counteract this, the user will normally have to add additional structure to the regression problem in order for a robust solution with desirable properties to exist.

To stabilize an ill-posed regression problem, we should first contemplate what would constitute desirable properties for the vector β provided that the regression data matrix Z imposes insufficient constraints on its behavior. A standard approach is to give preference to solutions for β with smaller norms, a choice we can motivate by the observation that if some of the regression coefficients are not constrained by the data, they should be set to zero to ensure that the corresponding regression variables (monomials) do not contribute to the fit. With this in mind, we choose a scalar regularization weight $w_{\text{reg}} > 0$ and replace (18.54) with

$$\|Y - Z\beta\|^2 + w_{\text{reg}} \|\beta\|^2 \rightarrow \min. \quad (18.57)$$

Intuitively, we should choose the weight w_{reg} small enough so that the extra term in the objective function does not interfere with the regression objective, yet sufficiently large that the regularization term performs its function. Numerous methods for data-driven selection of w_{reg} have been published in the literature, including the *L-curve method* in Hansen [1992], the *discrepancy principle* in Morozov [1966], and *generalized cross-validation* in Craven and Wahba [1979]. Andersen [2005] contains a review of many of these methods in the setting of yield curve construction. Whatever method is used, the size of w_{reg} should obviously reflect the relative scale of the numbers used in the regression. To examine the scaling issue a bit further, notice that the quadratic (in β) term in the objective function is given by

$$\beta^\top (Z^\top Z + w_{\text{reg}} I) \beta,$$

where I is the identity $q \times q$ matrix. The sum of squares of all elements in I and $Z^\top Z$ equals q and

$$\text{tr} \left((Z^\top Z)^\top (Z^\top Z) \right),$$

respectively. Consequently it is natural to write

$$w_{\text{reg}} = \epsilon \left(\frac{1}{q} \text{tr} (Z^\top Z Z^\top Z) \right)^{1/2}, \quad (18.58)$$

where ϵ is a new scale-free constant to be determined. While ideally we should rely on one of the data-driven approaches above, in a pinch we can always try to set ϵ equal to a small number, such as 10^{-4} .

The formal solution to (18.57) is given by

$$\beta = (Z^\top Z + w_{\text{reg}} I)^{-1} Z^\top Y. \quad (18.59)$$

Note that the matrix $Z^\top Z + w_{\text{reg}} I$ is of full rank even if $Z^\top Z$ is not, so the matrix inversion in (18.59) is always well-defined, even when Z is ill-conditioned. The resulting method (which we have used before, in Section 6.4.3) is often called *Tikhonov regularization* or *ridge regression*.

Tikhonov regularization is attractive because it retains a fair amount of intuition as to what happens to the regression coefficients as a result of regularization. We should, however, mention other regularization alternatives, in particular the *pseudo-inverse* or truncated singular value decomposition (TSVD) method. To briefly outline this approach, let us rewrite (18.55) as a system of linear equations on β ,

$$M\beta = Z^\top Y, \quad (18.60)$$

where $M = Z^\top Z$. The SVD method, see Press et al. [1992], allows us to decompose the $q \times q$ matrix M into a product of three matrices,

$$M = U\Sigma V^\top,$$

where U (not to be confused with the exercise value notation from earlier in the chapter) and V are $q \times q$ orthogonal matrices (i.e. $U^\top = U^{-1}$, $V^\top = V^{-1}$) and Σ is a diagonal $q \times q$ matrix. The diagonal elements of Σ are called *singular values* and are ordered by their absolute value (highest first). The decomposition applies even to singular matrices M . In particular, if M is of rank r , $r \leq q$, then only the first r diagonal elements of Σ will be non-zero.

The pseudo-inverse of the matrix M is defined by

$$M^+ = V\Sigma^+U^\top,$$

where Σ^+ is a diagonal matrix with elements

$$\Sigma_{i,i}^+ = \begin{cases} 1/\Sigma_{i,i}, & \Sigma_{i,i} \neq 0, \\ 0, & \Sigma_{i,i} = 0. \end{cases}$$

Then we have

$$MM^+ = M^+M = \text{diag}((1, \dots, 1, 0, \dots, 0)^\top),$$

where there are as many 1's on the diagonal as the rank of the matrix M .

The pseudo-inverse allows us to define a solution to (18.60) that always exists. For numerical stability, it is common to modify the solution slightly by choosing a truncation cut-off value $\epsilon > 0$ and defining a diagonal matrix Σ^ϵ by

$$\Sigma_{i,i}^\epsilon = \begin{cases} 1/\Sigma_{i,i}, & |\Sigma_{i,i}| \geq \epsilon|\Sigma_{1,1}|, \\ 0, & |\Sigma_{i,i}| < \epsilon|\Sigma_{1,1}|. \end{cases}$$

Then the solution to (18.60) and, ultimately, the regression problem (18.54), is given by

$$\beta = V \Sigma^\epsilon U^\top Z^\top Y.$$

A possible choice for ϵ is $\epsilon = 10^{-6}$.

To understand better the intuition behind TSVD, let us highlight an interesting connection between Tikhonov regularization and TSVD. Specifically, it can be shown that the Tikhonov solution (18.59) to the regression problem can be written as

$$\beta = V \Sigma^{\text{Tikhonov}} U^\top Z^\top Y.$$

where Σ^{Tikhonov} is a diagonal $q \times q$ matrix computed from the singular value matrix Σ as

$$\Sigma_{i,i}^{\text{Tikhonov}} = \frac{\Sigma_{i,i}}{\Sigma_{i,i}^2 + w_{\text{reg}}}.$$

We recognize this as a smoothed version of the cut-off matrix Σ^ϵ above, with the Tikhonov factor w_{reg} determining how much small singular values get damped out.

Singular values of widely different magnitudes in the matrix $Z^\top Z$ generally cause numerical problems in inversion of the matrix, something that the Tikhonov and TSVD method can help rectify. However, widely different scales of singular values do not necessarily arise only when explanatory variables are poorly chosen (e.g. highly dependent among themselves), but can also emerge if explanatory variables used are themselves of different scales. For example, if one variable is a swap rate measured in the units of a few percent and another is the value of the current coupon measured in the unit of millions of dollars, such scale discrepancy could lead to numerical problems in the regression. The problem is exacerbated by our choice of polynomials as basis functions, as one million to the power of, say, 4 is obviously quite different from one percentage point to the 4th power.

Fortunately, such scaling issues are easy to rectify, as we only need to rescale all variables to the same base before applying the regression. So, instead of the matrix Z we would use \tilde{Z} , whose elements are given by

$$\tilde{Z}_{i,j} = \frac{Z_{i,j} - \bar{Z}_{\cdot,j}}{\left(K^{-1} \sum_{k=1}^K (Z_{k,j} - \bar{Z}_{\cdot,j})^2 \right)^{1/2}}, \quad i = 1, \dots, K, \quad j = 1, \dots, q, \tag{18.61}$$

where $\bar{Z}_{\cdot,j} = K^{-1} \sum_{k=1}^K Z_{k,j}$. This transformation sets all columns in the matrix Z to have zero mean and unit (empirical) standard deviation. The only caveat with (18.61) is that, if applied to a column of Z with constant numbers — a column that is always present as we typically include a constant function in our regression — a division by 0 would occur. One obvious workaround is to simply avoid scaling the constant columns.

Once the standard deviation, as required by (18.61), of each column of Z is calculated, then we can apply another simple method to increase the robustness: we can watch out for points that are well outside a “reasonable” range — say outside of 10 standard deviations for a given column — and exclude them from our regression.

18.4 Valuation with Low-Dimensional Models

Libor market models are often our preferred choice for valuation and risk management of callable Libor exotics. For some CLEs, however, we can use simpler and faster models without sacrificing the benefits of proper calibration and good model dynamics. The trick here is to calibrate a simpler model in a special way, an approach we call the *local projection method*.

18.4.1 Single-Rate Callable Libor Exotics

The local projection method builds on the calibration discussion of Section 18.1 by calibrating a “local”, low-dimensional model to the volatility information that we identified as important to the CLE valuation. Information obtained from the market is used directly, and the rest is extracted from a “global”, fully-calibrated model such as the Libor market model. The success of the method depends on our ability to identify the relevant volatility information, and how well the local model can calibrate to this information. As a low-dimensional model has only a limited number of parameters, it can only be successfully calibrated for a CLE that depends on a relatively small subset of all available market information.

Callable Libor exotics most amendable to the local projection method are those that have, for each n , $n = 1, \dots, N - 1$, coupon C_n that is a function of at most a single market rate. We denote such structures *single-rate* CLEs. Examples include Bermudan swaptions, callable inverse floaters, callable CMS capped floaters and fixed-rate callable range accruals. Excluded are CLEs whose coupons depend on spreads between CMS rates, floating-rate callable range accruals, and similar.

The main attraction of using a low-dimensional model is the ability to apply PDE methods for valuation. We have already briefly discussed a relevant pricing scheme, see Section 2.7.4, and the mechanics of the valuation algorithm typically present no special difficulties unless the underlying CLE has path-dependent features. Section 18.4.5 discusses certain PDE techniques that can be used if the path-dependency is sufficiently weak.

18.4.2 Calibration Targets for the Local Projection Method

To start, let us focus on the heart of the local projection method, namely the choice of calibration targets for the local model. Let $\{S_n^1(t)\}_{n=1}^{N-1}$ be the

strip of swap rates that define coupons of a CLE, so that each C_n depends on $S_n^1(T_n)$, $n = 1, \dots, N - 1$. We assume that the swap rate S_n^1 has $\mu(n)$ periods, so that $S_n^1(t) = S_{n,\mu(n)}(t)$, where $S_{n,m}(t)$ is the standard notation for a swap rate fixing on T_n and covering m periods, see e.g. (4.10). For example, for a callable inverse floater we have $S_n(t) = L_n(t)$ (a Libor rate) and $\mu(n) = 1$; and for a callable CMS capped floater $\mu(n) \equiv k$, where k is the number of periods for the underlying CMS rate. In addition, we define a second strip of swap rates $\{S_n^2(t)\}_{n=1}^{N-1}$ to be the core, or coterminal, swap rate strip, i.e. $S_n^2(t) = S_{n,N-n}(t)$, $n = 1, \dots, N - 1$.

The underlying of the CLE, an exotic swap, can be expressed as a strip of options where the n -th option is written on the rate $S_n^1(T_n)$. Thus, for a model to reprice the underlying correctly, it should be calibrated to the market (spot) volatilities of the first swap rate strip, i.e. the implied volatilities of the European swaptions defined by $\{S_n^1(T_n)\}_{n=1}^{N-1}$. The ability to match the underlying exotic swap (i.e. the non-callable CLE) is certainly a prerequisite for any reasonable model for a callable CLE, but, as we have already seen in Section 18.1.1, we also need to consider that the volatilities and inter-temporal correlations of core swap rates $\{S_n^2(\cdot)\}_{n=1}^{N-1}$ will affect the value of the callability feature of the CLE. In light of this, as a starting point (to be refined later, see Section 18.4.4) we suggest that any low-dimensional model be calibrated to the following targets:

- The *underlying volatilities*, or swap rate volatilities for $\{S_n^1(T_n)\}_{n=1}^{N-1}$ that correspond to strikes relevant for the coupons C_n or, in a pinch, to at-the-money strikes.
- The *core volatilities*, or swap rate volatilities for $\{S_n^2(T_n)\}_{n=1}^{N-1}$. The choice of swaption strikes used to define core volatilities is often not straightforward, but at-the-money strikes is a common choice. In some cases more advanced methods for strike selection are available, see e.g. Section 19.3.
- The *core correlations*, or inter-temporal correlations for $\{S_n^2(T_n)\}_{n=1}^{N-1}$.

While volatilities of swap rates are directly observable from the market, the inter-temporal correlations are not. This is where we can draw on the LM (or similar) global model; once it has been calibrated to the market as a whole, we can calculate the required correlations from the global model. In a nutshell, the role of the global model is to serve as our “correlation extractor”. The important point here is that by including dynamic information such as inter-temporal correlations as calibration targets, the local model not only captures the static information about interest rate volatilities at valuation time, but also the transition densities and dynamics of the volatility structure, as seen by a global, fully calibrated and, presumably, realistic model.

18.4.3 Review of Suitable Local Models

The one-dimensional quasi-Gaussian (qG) model developed in Chapter 13 is a natural candidate to consider for the role of the local model in the local projection method for single-rate CLEs. A simple, yet useful, special case setup is based on the version of the qG model with linear local volatility, see Section 13.1.6. The volatility structure of such a qG model is controlled by several time-dependent functions, including the volatility function, the skew function and the mean reversion function. If convexity in the volatility smile is deemed important, the model could be upgraded to the stochastic volatility version in Section 13.2.

Let us first look at the volatility structure specification; we will consider skew and stochastic volatility parameter selection later on. With the volatility function and the mean reversion function discretized over the tenor structure $\{T_n\}_{n=0}^N$, the qG model has $2(N - 1)$ independent parameters for volatility calibration. As discussed in Section 13.1.7, the volatility parameters can be used to calibrate the model to term volatilities for one of the swap rate strips. The mean reversion can be used to either match the term volatilities for the second swap rate strip (Section 13.1.8.2), or the inter-temporal correlations (Section 13.1.8.3). As such, the one-factor qG model is not large enough to match all three sets of calibration targets identified in Section 18.4.2 above. In some situations, however, this might be acceptable as some securities may turn out to depend only weakly on one of the three calibration targets. For example, for shorter-dated CLEs, inter-temporal correlations may not affect the CLEs value all that much. Likewise, if the underlying has options on the rates $\{S_n^1(T_n)\}_{n=1}^{N-1}$ that are deep in or out of the money, this set of calibration targets can potentially be dropped. Of course, all such decisions must be supported by extensive testing, which fortunately is easy to do as we always have the global (LM or similar) model to benchmark against.

For derivatives where all three sets of calibration targets are important, a one-factor qG model will not suffice²⁴, and we ideally need to move to models with more stochastic factors. One particularly simple choice is here the two-factor Gaussian model, see Section 12.1.4, which has enough degrees of freedom to match all volatility targets. In this model, some of the time-dependent parameters can be chosen to be constant to make the dynamics of the volatility structure implied by the model more realistic, see Lemma 12.1.11.

The disadvantage of the two-factor Gaussian model is, of course, its lack of control over the volatility smile, so calibration to the volatility targets will require us to identify a single strike per swaption. Improved smile fits can be accomplished by using two-factor versions of either the quasi-Gaussian model (Chapter 13), the affine model, or the quadratic Gaussian model (Chapter

²⁴ Although we can always try to increase the range of applicability of the qG model with some of the techniques of Chapter 21.

12). While all these models are different in some regards, the underlying philosophy and calibration methods will be quite similar.

For models that are sufficiently rich to incorporate volatility skew/smile effects (such as local volatility or stochastic volatility qG models), we also need to select the market information to which we wish to calibrate skew and smile parameters. Normally we would extend one of the swaption strips, $\{S_n^1\}$ or $\{S_n^2\}$, to multiple strikes for this purpose (for mechanics of calibration see e.g. Section 13.2.3). The choice of the strip is typically driven by an analysis of relative importance of the two sets of smiles to the value of the CLE. It is difficult to state any firm general guidelines here, but we can observe that it is often fairly easy to match the underlying exotic swap value by a judicious choice of a single strike per maturity in the swaption strip $\{S_n^1\}$. On the other hand, it is often difficult to establish which strikes are the most relevant for the “callability” value. Given this, it is often reasonable to use whatever skew/smile parameters we have at our disposal to improve the broad fit of implied core swaption volatilities (the strip $\{S_n^2\}$) at multiple strikes per maturity. If we only have skew, but not smile, parameters, we can use these to match *two* volatilities at each maturity, or to match the slope of the volatility smile at a given strike. The latter could be important if the underlying structured coupons are not simple European options but, for example, of digital or range-accrual type, in which case it is the slope of the volatility smile, and not the overall level, that drives the underlying value. In this case we might, in fact, want to use skew to calibrate the underlying, rather than the callability, value.

18.4.4 Defining a Suitable Analog for Core Swap Rates

When we in Section 18.4.2 looked for the elements of the volatility structure that are relevant for a callable Libor exotic security, we argued that the callability value is driven by the volatilities of core swap rates $S_n^2(t) = S_{n,N-n}(t)$, since a CLE is related to a standard Bermudan swaption. This argument, clearly, has limitations of its applicability. For instance, in Section 19.4 we study Bermudan swaptions on amortizing swaps and show that the most relevant European swaptions in this case are not the standard core European swaptions, but swaptions with tenors based on the durations of the underlying amortizing swap.

In light of this, let us try to refine the selection of the volatility targets relevant for the callability option of a CLE. As a starting point we can use the idea that the local model should match the values of *European options on exercise values* $U_n(T_n)$, $n = 1, \dots, N - 1$. While market values of options on $U_n(T_n)$ could be hard to come by, we can linearize the underlying $U_n(T_n)$ of the CLE and use the resulting rate as a replacement for the core swap rate that should be used in volatility and correlation calibration.

Using a LM model as a backdrop for our analysis, the exercise values $U_n(T_n)$, for each $n = 1, \dots, N - 1$, are functions of the vector of primary

Libor rates

$$\mathbf{L}(T_n) = (L_n(T_n), \dots, L_{N-1}(T_n))^{\top}$$

observed on the date T_n , i.e.,

$$U_n(T_n) = f_n(\mathbf{L}(T_n)), \quad n = 1, \dots, N - 1.$$

Linearizing this expression, we obtain

$$U_n(T_n) \approx f_n(\mathbf{L}(0)) + \nabla f_n(\mathbf{L}(0)) (\mathbf{L}(T_n) - \mathbf{L}(0)),$$

where $\nabla f_n(x)$ is the (row vector) gradient of $f_n(x)$. Hence, the value of the European option on the underlying,

$$E \left(B(T_n)^{-1} (U_n(T_n))^+ \right),$$

can be approximated with

$$E \left(B(T_n)^{-1} \left(\nabla f_n(\mathbf{L}(0)) \mathbf{L}(T_n) - \left(\nabla f_n(\mathbf{L}(0)) \mathbf{L}(0) - f_n(\mathbf{L}(0)) \right) \right)^+ \right). \quad (18.62)$$

We can therefore argue that the most relevant “interest rate” is

$$R_n(T_n) = \nabla f_n(\mathbf{L}(0)) \mathbf{L}(T_n) = \sum_j w_{n,j} L_j(T_n), \quad w_{n,j} \triangleq \frac{\partial f_n}{\partial L_j}(\mathbf{L}(0)).$$

Being a linear combination of Libor rates, $R_n(T_n)$ is not, strictly speaking, a market swap rate. However, the volatility of the rate $R_n(T_n)$ can be approximated in a Libor market model (along the same lines as in Section 14.4.2), as well as in local models we may wish to use. Therefore, we can easily use the term volatilities of $R_n(T_n)$, $n = 1, \dots, N - 1$, as volatility targets in place of core swap rates.

The underlying $U_n(t)$ typically consists of options on market rates. The derivatives $\partial f_n / \partial L_j$ can then be computed quite easily with, say, Black-type approximations to option values. Volatilities that should be used in these calculations are the forward (as observed at time T_n) volatilities. Needless to say, high degree of precision is not necessary in these calculations.

To see how consistent the method defined above with our recommendations in Section 18.4.2, let us apply it to a standard Bermudan swaption. For a payer swap with coupon K we have

$$\begin{aligned} U_n(T_n) &= \sum_{i=n}^{N-1} \tau_i (L_i(T_n) - K) P(T_n, T_{i+1}) \\ &= \sum_{i=n}^{N-1} \tau_i (L_i(T_n) - K) \left(\prod_{k=n}^i \frac{1}{1 + \tau_k L_k(T_n)} \right). \end{aligned}$$

Hence,

$$f_n(x) = \sum_{i=n}^{N-1} \tau_i (x_i - K) \left(\prod_{k=n}^i \frac{1}{1 + \tau_k x_k} \right),$$

and

$$\frac{\partial f_n}{\partial L_j}(x) = \tau_j \prod_{k=n}^j \frac{1}{1 + \tau_k x_k} - \frac{\tau_j}{1 + \tau_j x_j} \sum_{i=j}^{N-1} \tau_i (x_i - K) \left(\prod_{k=n}^i \frac{1}{1 + \tau_k x_k} \right).$$

Thus

$$w_{n,j} = \tau_j P(0, T_n, T_{j+1}) (1 - U_j(0)/P(0, T_j)).$$

If we compare these weights with those obtained by decomposition of the swap rate into a sum of Libor rates via the “freezing” techniques in Section 14.4.2, we see that they are roughly the same, up to a constant scaling. Thus

$$R_n(T_n) = \sum_{j=n}^{N-1} w_{n,j} L_j(T_n)$$

is quite close to the (scaled) core swap rate $S_{n,N-n}(T_n)$. Therefore, for standard Bermudan swaptions, using $R_n(T_n)$ to define volatility calibration targets should be approximately consistent with the standard method of using core swap rates. When applied to non-standard (e.g., amortizing) Bermudan swaptions, this method produces results that are similar to what we propose later in Section 19.4.

It should be clear that the choice of calibration targets has carries significant impact on the value of a CLE in a local model. Equally important, it also defines the “basis” for vegas (volatility sensitivities), i.e. the set of swaption volatilities to which the CLE is sensitive; hedging of volatility exposure in local model would therefore, as a practical matter, only be done with the swaptions included in the calibration strips. Using a particular swaption for calibration implies the dependence of the CLE value to the volatility of that swaption; conversely, omitting a swaption from the calibration set makes the CLE value (in the local model) insensitive to its volatility. This, of course, is not wholly realistic as even simple CLEs (including Bermudan swaptions) would, when priced in a global model, typically show sensitivity to volatilities of *all* swaptions whose total maturity is no greater than the CLE maturity²⁵. In some sense, vegas from a local model (to a subset of swaptions) can be thought of as an aggregation of vegas from the global model. Some traders may in fact prefer such an aggregated view as it (seemingly) simplifies the job of vega hedging.

²⁵For more detail on this topic, see Chapter 26.

18.4.5 PDE Methods for Path-Dependent CLEs

As mentioned earlier, one attraction of the local projection method is the fact that the resulting model state can often be represented by a low-dimensional state vector $x(t)$. If the dimension of $x(t)$ is less than 3 or 4, this will often allow us to state the CLE value as the solution to a PDE, a problem that can be attacked by the finite difference methods in Chapter 2. While callability is easy to handle (see Section 2.7.4), most path-dependent CLEs are outside the scope of finite difference methods. Exceptions do exist, however, if the path-dependency is sufficiently mild. We show some examples of this below. Of course, even in those cases where a PDE solution is technically possible, one should contemplate whether a local projection model is fundamentally suitable for the path-dependent derivative in question. In particular, the basic single-rate CLE calibration strategies may need adjustment to better capture the path-dependent feature of the payout.

18.4.5.1 CLEs Accreting at Coupon Rate

One particular class of path-dependent CLEs that is amenable to PDE methods has its path-dependency confined to the CLE *notional* only, see Piterbarg [2002]. A prime example of such a CLE is a callable Libor exotic accreting at a coupon rate, see Section 5.14.5, which is the example we consider here. Recall that a CLE is defined by its structured coupon C_n that is fixed at time T_n and paid at time T_{n+1} , $n = 1, \dots, N - 1$. In the standard CLE, the notional of the coupon is constant, or at least deterministic, and has been factored out from the definition in Section 18.2.1. For a coupon-accreting CLE, the notional to which the coupon rate and the Libor rate are applied at time T_{n+1} is equal to the notional at time T_n times an accretion factor that depends on C_n .

Formally, we replace (18.1) with

$$X_n = D_n \tau_n (C_n - L_n(T_n)),$$

where $D_1 = 1$ and

$$D_n = D_{n-1} \times (1 + \tau_{n-1} C_{n-1}), \quad n = 2, \dots, N - 1. \quad (18.63)$$

A coupon-accreting CLE is defined as a Bermudan-style option to enter, on date T_n , the remaining part of the underlying, i.e. an exotic swap with the value (at time $t \leq T_n$),

$$U_n(t) = B(t) \sum_{i=n}^{N-1} E_t (B(T_{i+1})^{-1} X_i).$$

The backward-induction scheme (18.7), (18.16) is trivially extended,

$$\begin{aligned} U_n(T_n) &= \tau_n D_n (C_n - L_n(T_n)) P(T_n, T_{n+1}) \\ &\quad + B(T_n) E_{T_n} (B(T_{n+1})^{-1} U'_{n+1}(T_{n+1})) , \end{aligned} \quad (18.64)$$

$$H_n(T_n) = B(T_n) E_{T_n} (B(T_{n+1})^{-1} \max(U'_{n+1}(T_{n+1}), H'_{n+1}(T_{n+1}))) . \quad (18.65)$$

As written, the scheme cannot be directly implemented in a PDE solver because of path-dependency under the expected value operator in (18.64). However, by employing the method of *similarity reduction*, the scheme can be rewritten in a way amendable to a PDE representation.

Dividing both sides of (18.64)–(18.65) by D_n , and using the fact that D_n is \mathcal{F}_{T_n} -measurable, we get

$$\begin{aligned} U'_n(T_n) &= \tau_n (C_n - L_n(T_n)) P(T_n, T_{n+1}) \\ &\quad + B(T_n) E_{T_n} \left(B(T_{n+1})^{-1} \frac{D_{n+1}}{D_n} U'_{n+1}(T_{n+1}) \right) , \\ H'_n(T_n) &= B(T_n) E_{T_n} \left(B(T_{n+1})^{-1} \frac{D_{n+1}}{D_n} \max(U'_{n+1}(T_{n+1}), H'_{n+1}(T_{n+1})) \right) , \end{aligned}$$

where we have defined

$$U'_n(T_n) = \frac{U_n(T_n)}{D_n}, \quad H'_n(T_n) = \frac{H_n(T_n)}{D_n}, \quad n = 1, \dots, N-1.$$

From (18.63),

$$\frac{D_{n+1}}{D_n} = 1 + \tau_n C_n,$$

where $1 + \tau_n C_n$ is \mathcal{F}_{T_n} -measurable. Therefore, the factor D_{n+1}/D_n can be pulled out from inside the expected value operator, to give us

$$\begin{aligned} U'_n(T_n) &= \tau_n (C_n - L_n(T_n)) P(T_n, T_{n+1}) \\ &\quad + (1 + \tau_n C_n) B(T_n) E_{T_n} (B(T_{n+1})^{-1} U'_{n+1}(T_{n+1})) , \end{aligned} \quad (18.66)$$

$$\begin{aligned} H'_n(T_n) &= (1 + \tau_n C_n) B(T_n) \\ &\quad \times E_{T_n} (B(T_{n+1})^{-1} \max(U'_{n+1}(T_{n+1}), H'_{n+1}(T_{n+1}))) . \end{aligned} \quad (18.67)$$

These equations are used for $n = N-1, \dots, 0$, with $U'_N(T_N) = H'_N(T_N) = 0$. We have the following result.

Proposition 18.4.1. *For a coupon-accreting CLE, $U'_n(T_n)$ and $H'_n(T_n)$ can, for each $n = 0, \dots, N$, be written as deterministic functions of the model state variables model at time T_n .*

Proof. The proof is by induction. The statement is trivially true for $n = N$. To prove the induction step, we assume it is true for $n+1$. We note that

$$B(T_n) E_{T_n} (B(T_{n+1})^{-1} U'_{n+1}(T_{n+1}))$$

and

$$B(T_n)E_{T_n} \left(B(T_{n+1})^{-1} \max (U'_{n+1}(T_{n+1}), H'_{n+1}(T_{n+1})) \right)$$

can be written as functions of the state variables at time T_n by applying the PDE rollback scheme to

$$U'_{n+1}(T_{n+1}), \quad \max (U'_{n+1}(T_{n+1}), H'_{n+1}(T_{n+1}))$$

(which are functions of the state variables at time T_{n+1} by the induction hypothesis). The accreting factor $(1 + \tau_n C_n)$ and the marginal coupon $\tau_n(C_n - L_n(T_n))P(T_n, T_{n+1})$ are functions of the state variables at time T_n as well. The proposition is proved. \square

The new scheme (18.66)–(18.67) in fact looks just like (18.64)–(18.65) for a unit-notional CLE, with one modification: on each backward induction step, the values of U' and H' are rescaled by the “marginal” accreting notional $(1 + \tau_n C_n)$. The key fact here is that this factor is known at time T_n .

18.4.5.2 Snowballs

Certain other path-dependent CLEs can be valued by PDE methods by introducing extra state variables, along the lines of Section 2.7.5. Among the more popular CLEs for which this method is applicable are the snowball swaps and callables, see Chapter 5. In a snowball, the structured coupon at time T_n is a function of the structured coupon at time T_{n-1} (and rates at time T_n). As an example, recall the basic structure from Section 5.13.5 with the coupon at time T_n given by

$$C_n = (C_{n-1} + s_n - g_n \times L_n(T_n))^+, \quad n = 1, \dots, N-1,$$

with C_0 being a fixed initial coupon.

We can value snowball swaps and callables by PDE method after the introduction of an extra state variable $I(t)$ defined to be the current coupon,

$$I(t) = \sum_{n=0}^{N-1} 1_{\{t \in [T_n, T_{n+1})\}} C_n.$$

The backward recursion for the exercise value at time T_n then reads,

$$U_n(T_n) = \tau_n (I(T_n) - L_n(T_n)) P(T_n, T_{n+1}) + E_{T_n} \left(\frac{B(T_n)}{B(T_{n+1})} U_{n+1}(T_{n+1}) \right). \quad (18.68)$$

Hence we obtain the following continuity condition for a given value of $I(T_n) = I$ (where T_n- is a time immediately prior to T_n),

$$\begin{aligned} U_n(T_n-, I) &= \tau_n (I - L_n(T_n)) P(T_n, T_{n+1}) \\ &\quad + U_n \left(T_n, (I + s_n - g_n \times L_n(T_n))^+ \right), \end{aligned}$$

combined with the following one-period rollback scheme,

$$U_n(T_n, I) = \mathbb{E}_{T_n} \left(\frac{B(T_n)}{B(T_{n+1})} U_{n+1}(T_{n+1}-, I) \right), \quad (18.69)$$

$n = N - 1, \dots, 1$. For the hold value, the continuity condition is

$$H_n(T_n-, I) = H_n \left(T_n, (I + s_n - g_n \times L_n(T_n))^+ \right), \quad (18.70)$$

where

$$H_n(T_n, I) = \mathbb{E}_{T_n} \left(\frac{B(T_n)}{B(T_{n+1})} \max(U_{n+1}(T_{n+1}-, I), H_{n+1}(T_{n+1}-, I)) \right). \quad (18.71)$$

The PDE scheme may be implemented by discretizing I over an appropriate range, solving PDEs (18.69)–(18.71) on each I -plane, and interfacing the solutions between I -planes at times T_n , $n = 1, \dots, N - 1$ using (18.68)–(18.70). The details of implementation follow the general plan of Section 2.7.5, and we do not repeat them here.

Bermudan Swaptions

After our general discussion of callable Libor exotics in the previous chapter, we now turn our attention to an important subset of the generic CLE class, the Bermudan swaptions. Bermudan swaptions are among the most liquid exotic interest rate derivatives, and the demands they place on accuracy, fidelity and performance of term structure models have driven many advances in interest rate modeling. While the ideas and methods from the previous chapter all apply to Bermudan swaptions, the simpler structure of Bermudan swaptions compared to general CLEs allows us to considerably deepen our analysis of valuation and risk management methods.

19.1 Definitions

As defined in Section 5.12, a Bermudan swaption is a callable Libor exotic with the coupon paying a fixed rate, $C_n = k$, $n = 1, \dots, N - 1$. Alternatively, we can consider it a Bermudan-style option to enter a simple fixed-for-floating swap. The fixed rate k is often referred to as the *strike* of the Bermudan swaption. Exercise dates of a Bermudan swaption are typically¹ a subset of a tenor structure $\{T_n\}_{n=0}^N$ that defines the underlying swap. In a standard structure, exercise is restricted to the dates $\{T_n\}_{n=s}^{N-1}$ where $s \geq 1$; as we explained in Section 5.12, the period up to T_s is known as the lockout or no-call period. Recall that a Bermudan swaption on, say, a 10 year swap with a 2 year lockout period (at inception) is known as a “10 no-call 2”, or “10nc2”, Bermudan swaption. For convenience (and without loss of generality) we assume in most of this chapter that all $\{T_n\}_{n=1}^{N-1}$ are, in fact, exercise dates. If the Bermudan swaption is exercised at time T_n , the exercise value, for a payer swap, is given by

$$U_n(t) = \sum_{i=n}^{N-1} \tau_i P(t, T_{i+1}) (L_i(t) - k), \quad (19.1)$$

¹But see Sections 19.4.7 and 19.4.8 below for exceptions.

where k is the fixed rate. We note that $U_n(t)$ can be written as

$$U_n(T_n) = A_n(T_n)(S_n(T_n) - k),$$

where $A_n(t) \triangleq A_{n,N-n}(t)$ is the annuity, and $S_n(t) \triangleq S_{n,N-n}(t)$ is the swap rate for the swap into which one can exercise at time T_n (see notations (4.8), (4.10)). The definition of hold values carries over unchanged from Chapter 18.

19.2 Local Projection Method

As Bermudan swaptions are liquid and their volume is relatively high, the performance advantages of PDE methods over Monte Carlo simulation lead many market participants to value Bermudan swaptions in low-factor Markovian models, using either finite difference grids or trees. A sound framework for the usage of low-dimensional Markovian models is provided by the local projection method for single-rate CLEs that we discussed in Section 18.4. The method takes a particularly simple form for Bermudan swaptions, as the underlying swaps have no optionality and only the volatility parameters of core swap rates $\{S_n(\cdot)\}_{n=1}^{N-1}$ are relevant. As we discussed before (in Section 13.1.8.1), we can view a Bermudan swaption as the option to choose the “best” among a collection of swap values observed on different dates. This implies that a Bermudan swaption value is driven by *core volatilities*, or volatilities of the core swap rates $\{S_n(T_n)\}_{n=1}^{N-1}$, and *core correlations*, or inter-temporal correlations of the core swap rates $\{S_n(T_n)\}_{n=1}^{N-1}$. Alternatively, of course, we can think of forward volatilities in place of inter-temporal correlations as the source of “exotic” risk in Bermudan swaptions, see Section 18.1.1.

The relative simplicity of the dependence of Bermudan swaptions on the volatility structure allows us to use models as simple as the one-factor Gaussian model (Section 10.1.2) for valuation and risk management. The time-dependent volatility is typically calibrated to core swaption volatilities, while the mean reversion is calibrated to inter-temporal correlations of core swap rates (see Section 13.1.8.3); these correlations could, for instance, be extracted from an LM model. In practice, it is not unusual to skip the last step — since Bermudan swaptions have been traded well before LM models (or other practical multi-factor models) were invented, a market practice has developed whereby the mean reversion of a Gaussian (or similar) model is used essentially as a free parameter, rather than implied from a global model. This practice continues today, with mean reversions often set to match the “market” prices of Bermudan swaptions that are sometimes observable, or quasi-market prices such as independently-produced averages of dealer-submitted prices of a few typical structures². Another fairly popular

²At the time of writing this is done by Markit, see www.markit.com.

choice would set mean reversions to match caplet volatilities, although using caplets for mean reversion calibration is rather arbitrary and can sometimes lead to odd mean reversion curves (see for instance the discussion in Section 10.1.2.3). The practice can, however, perhaps be justified if caplets are used as hedging instruments for inter-temporal correlation or forward volatility; technical details are available in Section 13.1.8.2.

Turning to the issue of volatility smile, we recall that the Gaussian model basically has no control over it, and the model can only be calibrated to one³ volatility per expiry T_n , $n = 1, \dots, N - 1$. A one-factor quasi-Gaussian (qG) model with local volatility (Section 13.1) would constitute an improvement, as it will also allow to capture the slopes of volatility smiles of core swaptions, in addition to volatilities at specific strikes. Finally, the stochastic volatility version of the one-dimensional qG model (Section 13.2) would essentially allow for (best-fit) calibration to all core swaption volatility smiles across all strikes. On balance, the local volatility qG model is probably sufficient for effective risk management of Bermudan swaptions, although we would of course choose the SV version if available computing power permits. Finally, a two-factor quadratic Gaussian model of Section 12.3 is also a viable choice for Bermudan swaption pricing.

While we are on the subject of the model choice, let us briefly comment on the discussion around what number of factors is appropriate for a Bermudan swaption model. While the usage of single-factor, or essentially single-factor models such as the qG model, for Bermudan swaption valuation is widespread, some argue that single-factor models significantly underprice Bermudan swaptions. The basic claim is that higher de-correlation in rates has a positive impact on Bermudan swaption prices (as a Bermudan swaption is a “best-of” option on swap rates) and two- and multi-factor models intrinsically are able to de-correlate rates more than a single-factor model (where the instantaneous correlations between moves in all forward rates is always one). There are a number of flaws in this argument, starting with the fact that the correlations relevant for Bermudan swaptions are the *inter-temporal* correlations, which can be easily manipulated in a one-factor model through the choice of mean reversion. In addition, when comparing one- and multi-factor rates models, it is obviously important that calibration to European swaptions is unaffected by changes in the number of factors. Careful analysis in Andersen and Andreasen [2001] of Bermudan swaption pricing in a one- and a two-factor Gaussian models shows that, if the models are calibrated in consistent fashion to core European swaptions, the two-factor model price is in fact slightly *lower* than the one-factor price. Experiments with LM models with different numbers of Brownian motions, but all calibrated to the full swaption grid, confirm this analysis.

³For more information on *which* volatility to calibrate to, see Section 19.3 below.

The reason for the slight decrease in Bermudan swaption price as a function of the number of factors may seem puzzling at first and is, indeed, a rather subtle effect. While there are numerous factors in play (see Andersen and Andreasen [2001] for a full analysis), one important observation is that *forward* volatility in a low-factor model generally is higher than in a multi-factor model, as long as both models are in calibration with the European swaption market. A technical explanation for this phenomenon can be found in Appendix 19.A of this chapter; loosely speaking the effect stems from the fact that a one-factor model will imply a lower time 0 instantaneous forward volatility term structure than will a multi-factor model, a relationship that must be reversed as time progresses to preserve overall variance.

19.3 Smile Calibration

Let us now discuss the issues of smile calibration in more detail. For concreteness, we consider the linear local volatility version of the quasi-Gaussian model from Section 13.1. As we have seen before, the model has enough flexibility to match the level and the slope of the volatility smile for each of the core swap rates. The market volatility smile is, of course, not close to linear, so we often seek to match the volatility of a *particular strike* exactly, or as closely as possible, while roughly matching the overall slope of the volatility smile.

The simplest approach to choosing the strikes that define core volatilities for calibration involves using at-the-money (ATM) strikes for each core swap rate. This rather crude approach is still in use for (we assume) historical reasons, as Bermudan swaptions started trading well before pronounced market smiles developed in interest rate markets, and probably even before the non-ATM points of the swaption volatility cube became liquid enough to keep track of them. Proponents of this approach sometimes rely on hedging arguments, as volatility exposure of a Bermudan swaption is often vega hedged using the most liquid European swaptions — which happen to be at-the-money. Yet another possible justification for the ATM strike choice notes that using the same volatilities for Bermudan swaptions of different strikes (seemingly) ensures consistency of valuation across Bermudan swaptions of different strikes. In reality, however, the ATM strike choice leads to *inconsistent* valuation between European and Bermudan swaptions: if one uses a Bermudan swaption model calibrated to ATM swaptions, and applies it to a Bermudan swaption with a non-ATM strike and just a *single* exercise date, the value is going to be different from the value of the same derivative priced as a European swaption. Clearly, this is a strongly undesirable feature of the ATM strike calibration idea.

To ensure consistency between the Bermudan swaption and its underlying core European swaptions, it suffices to set the calibration strike equal to that of the Bermudan swaption itself. That is, if the fixed rate of the

Bermudan swaption is k , then for each expiry T_n one uses the volatility of the appropriate core European swaption that corresponds to the strike k . This method automatically ensures that a European swaption valued as a single-exercise Bermudan swaption has exactly the same value in the model as in the market. The fact that all swaptions we can exercise into are priced exactly is intuitively appealing, and also ensures that certain rational bounds for the Bermudan swaption price will not be violated. Indeed, if $V_{\text{swaption},n}(0; k)$ is the price of a k -strike, T_n -expiry swaption on a swap that matures at time T_N , then clearly⁴ the k -strike Bermudan swaption price $V_{\text{Berm}}(0; k)$ at time 0 must satisfy (compare to (18.3))

$$V_{\text{Berm}}(0; k) \geq \max_{n=1, \dots, N-1} V_{\text{swaption},n}(0; k), \quad (19.2)$$

where we as mentioned earlier have assumed that exercise can take place at all T_n , $n \geq 1$. If our model fundamentally matches all swaption prices inside the max-operator on the right-hand side of (19.2), then pricing the Bermudan swaption in, say, a finite difference grid will always return a Bermudan swaption that satisfies (19.2). We notice as an aside that the (non-negative) difference between the left- and right-hand sides of (19.2) is sometimes known in trader jargon as the *Bermudanality* of the Bermudan swaption.

While enforcing consistency with European swaptions is useful, the idea of at-the-strike calibration is not a panacea, as Bermudan and European swaptions can behave quite differently. For instance, Bermudan swaptions have other “interesting”, e.g. high-convexity, points in the swap rate dimension than just the underlying strike k , the most important of which is the exercise boundary, i.e. the swap rate level (for each time T_n) at which the decision to exercise the swaption switches to the decision to hold. The importance of this point is clearly seen from the marginal exercise value decomposition (18.8), as it corresponds to “strikes” of European options in the representation of the value of a Bermudan swaption as a sum of European options. Hence, a third calibration option available is to use the swap rate volatilities that correspond to the exercise boundary on each of the exercise dates. As was the case for the Bermudan strike method, this method makes valuation of single-exercise Bermudan swaptions consistent with the valuation of (equivalent) European swaptions, since the exercise boundary for a European swaption coincides with its strike. For the same reason, any weighted average of the strike and the exercise boundary would also produce a consistent scheme. We do point out, however, that using calibration strikes other than that of the Bermudan swaption may lead to violations of (19.2).

⁴The optimal exercise strategy for a Bermudan swaption must be as least as good as simply picking at time 0 one of the exercise dates and never changing one's mind.

To provide a bit more detail on the idea of using the exercise boundary to select calibration strikes, let us first observe that for any given model, one can determine the exercise boundary as a function of the state variables. To be able to calibrate to European swaption volatilities with the strike at the exercise boundary, one has to be able to translate this “model” exercise boundary into a value of the corresponding core swap rate. Strictly speaking, this can be done unambiguously only in single-factor models, such as the Gaussian model. For the one-factor qG model with its two state variables the boundary is, in fact, represented by a line in a two-dimensional plane of possible values of the x and y state variables; each of the points on this line corresponds, potentially, to a different value of the core swap rate. However, the dependence of the exercise boundary on y is rather mild, a fact that should come as no surprise if one recalls the “auxiliary” nature of the y state variable, see Chapter 13. Hence, for the qG model, we can use the expected value of $y(T_n)$ when converting the exercise boundary for the state variable $x(T_n)$ into the strike for the swap rate $S_n(T_n)$.

The choice of exercise boundary for calibration is, unfortunately, rather inconvenient from the implementation point of view because the exercise boundary information is not available until *after* the valuation algorithm has been run. One can try a recursive scheme where one uses some (e.g., strike-calibrated) volatilities for an initial calibration, values a Bermudan swaption, calculates the exercise boundary, looks up the core volatilities for calibration at this boundary, calibrates the model again, values the Bermudan swaption, and so on. Such a procedure can in fact diverge; thus, one is forced to limit the number of iterations artificially, potentially resulting in unstable risk sensitivities and other problems. Moreover, this scheme in general consumes more computational resources due to the multiple valuations required. For all these reasons the at-the-boundary calibration is of limited use and the at-the-strike volatility calibration method is probably the most reasonable in practice, combining ease of implementation and consistency with European swaptions.

19.4 Amortizing, Accreting, and Other Non-Standard Bermudan Swaptions

A standard (or *vanilla* or *bullet*) Bermudan swaption is characterized by the fact that the notional on which coupons are paid are all identical, as in (19.1) (where the notional of all coupons are 1). A relatively popular extension involves making the notional of a Bermudan swaption time-dependent and deterministic, with the exercise value given by

$$\tilde{U}_n(t) = \sum_{i=n}^{N-1} R_i \tau_i P(t, T_{i+1}) (L_i(t) - k) \quad (19.3)$$

(compare to (19.1)). Here R_i is the notional for the i -th coupon, $i = 1, \dots, N - 1$.

If the notional increases with the coupon index, the Bermudan swaption is said to be *accreting*; if it decreases, it is said to be *amortizing*. Other profiles are possible but are much less common. Amortizing Bermudan swaptions are often used as hedges for pools of mortgages, with the amortization feature mimicking prepayments on the pool. Accreting Bermudan swaptions, on the other hand, often appear as a result of issuing “zero coupon” structured notes, i.e. notes with the repayment notional growing over time but paying no coupons during the life of the note, as explained in Section 19.4.6.

Since the notionals in Bermudan swaptions of type (19.3) are still deterministic (even if time-varying), their valuation in a properly calibrated model does not present any particular technical difficulties⁵ — but what constitutes “properly calibrated”, however, is not always obvious. Of course, in models with global calibration (e.g. LM models), calibration for non-standard Bermudan swaptions is no different from calibration for standard Bermudan swaptions, as model calibration is product-independent by definition. On the other hand, for models requiring local calibration, such as a one-factor Gaussian or a quasi-Gaussian model, calibration for non-standard Bermudan swaptions will require additional analysis. We consider this problem in the next few sections, but it is worth pointing out that, in the opinion of some, making notionals time-dependent pushes Bermudan swaptions across the boundary that separates those securities for which local models are acceptable to use from those for which globally-calibrated multi-factor models are required.

Before commencing on a more detailed analysis, let us first briefly try to understand the basic complications involved in local model calibration for non-vanilla Bermudan swaptions. As established previously, a locally-calibrated model should, as a minimum, be calibrated to the volatilities of core swap rates. For an amortizing Bermudan swaption, say, a core swap rate would correspond to an amortizing swap. Volatilities of amortizing swaps can be extracted from amortizing European swaptions, but the liquidity of such swaptions is significantly poorer than for vanilla European swaptions — in fact, amortizing European swaptions are about as illiquid as Bermudan swaptions themselves. In practice, if one wishes to calibrate a local model to core (amortizing) swaptions, one may need to use a “pre-processing” step to extract amortizing European swaption prices from a model calibrated to liquid vanilla European swaptions, as we do later in Section 19.4.4. Alternatively, one needs to choose a different set of calibration targets in the first place, see Section 19.4.3.

⁵Although this is not true for the family of Markov-functional models, see Appendix 11.A in Chapter 11; in fact the difficulty of handling amortizing/accreting Bermudan swaptions is often cited as one of the problems with such models.

19.4.1 Relationship Between Non-Standard and Standard Swap Rates

Regardless of the calibration method ultimately used, it is useful to understand the relationship between non-standard and standard swaps (and, hence, swap rates). To be consistent with (19.3), consider a swap that starts at T_n , ends at T_N , and has a notional schedule $\{R_i\}$. The time t value of such a swap is given by $\tilde{U}_n(t)$ in (19.3). To make some of the formulas below simpler, let us extend the notional schedule by one period and set $R_N = 0$. We denote the annuity and the swap rate that correspond to this non-standard swap by $\tilde{S}_n(\cdot)$ and $\tilde{A}_n(\cdot)$, so that

$$\tilde{A}_n(t) = \sum_{i=n}^{N-1} R_i \tau_i P(t, T_{i+1}), \quad \tilde{S}_n(t) = \tilde{A}_n(t)^{-1} \sum_{i=n}^{N-1} R_i \tau_i P(t, T_{i+1}) L_i(t).$$

We would like to decompose the non-standard swap into a linear combination of standard swaps. Such a decomposition is, however, not unique and could be done in a multitude of ways, potentially using any of the standard swaps with starting date on or after T_n , and final payment date on or before T_N . To narrow down the problem, we note that we here are ultimately interested in establishing the volatility of the non-standard rate $\tilde{S}_n(\cdot)$ over the period $[0, T_n]$. Since the values of standard European swaptions provide us with the information on the volatilities of swap rates *only* over the period from time 0 to their start dates, we should focus only on standard swaps that start on T_n ; this choice makes the decomposition unique.

Let us denote the value of a standard swap starting at T_n and covering m periods by $V_{n,m}(t)$, and the corresponding annuity and swap rate by $S_{n,m}(t)$ and $A_{n,m}(t)$ (see (4.8), (4.10)). In light of the discussion above, we want to find weights $\{v_{n,m}\}$, $m = 1, \dots, N - n$, such that

$$\tilde{U}_n(T_n) = \sum_{m=1}^{N-n} v_{n,m} V_{n,m}(T_n).$$

Note that only swaps starting at time T_n are used in the right-hand side of this expression. Matching terms to (19.3) we obtain

$$v_{n,m} = R_{n+m-1} - R_{n+m}, \quad m = 1, \dots, N - n, \quad (19.4)$$

so that

$$\tilde{U}_n(T_n) = \sum_{m=1}^{N-n} (R_{n+m-1} - R_{n+m}) V_{n,m}(T_n).$$

After some algebraic manipulations, we obtain the following relationship for the swap rates,

$$\tilde{S}_n(T_n) = \sum_{m=1}^{N-n} w_{n,m}(T_n) S_{n,m}(T_n), \quad (19.5)$$

where (recall that we set $R_N = 0$)

$$w_{n,m}(T_n) = (R_{n+m-1} - R_{n+m}) \frac{A_{n,m}(T_n)}{\tilde{A}_n(T_n)}.$$

While the swap weights $v_{n,m}$ are deterministic, the swap *rate* weights $w_{n,m}(T_n)$ are not. For a qualitative discussion, however, we note that the weights $w_{n,m}(T_n)$ can be approximated reasonably well by their values at time 0,

$$\tilde{S}_n(T_n) \approx \sum_{m=1}^{N-n} w_{n,m}(0) S_{n,m}(T_n), \quad w_{n,m}(0) = (R_{n+m-1} - R_{n+m}) \frac{A_{n,m}(0)}{\tilde{A}_n(0)}. \quad (19.6)$$

From the expression (19.6) it follows that the volatility of the non-standard swap rate is a function of volatilities of all standard swap rates with a given expiry (T_n in our case), and of their correlations. Putting correlations aside for a moment, observe that to price a non-standard Bermudan swaption, in principle one needs to calibrate the model to volatilities of standard rates with all expiries (T_1, \dots, T_{N-1}) and all maturities, something a low-dimensional local model will virtually never be able to do. Below, we discuss two possible strategies for going forward.

19.4.2 Same-Tenor Approach

Perhaps the simplest calibration approach for non-standard Bermudan swaptions is to simply pretend that they are standard Bermudan swaptions and set up the model calibration accordingly. So, for expiry T_n , one would choose a European swaption on the $(N - n)$ -period swap as the calibration instrument. While easy to implement, the merits of this approach are obviously somewhat wanting, and we do not recommend it. Nevertheless, it is instructive to investigate the issues that would come up if we adopted this scheme. As an example, consider an amortizing Bermudan swaption and a one-factor model calibrated to standard swaptions of the same tenor as the core amortizing swaptions. As a thought experiment, suppose that we increase mean reversion in the model, while keeping it calibrated to our calibration swaption set. In this case, the core amortizing swaption prices would increase, a simple consequence of our decomposition (19.6) and the fact that shorter-tenor (standard) swap rate volatilities increase as a function of mean reversion when volatility of a longer-tenor swap rate is kept fixed, see the discussion in Section 13.1.8.1. As a consequence, mean reversion would affect not only the inter-temporal correlations that are important for Bermudan swaptions, but would also affect the volatilities of core swap rates. In the context of a local projection method, we would then face a dilemma as to which targets to calibrate the mean reversion to: the inter-temporal correlations or the prices of amortizing European swaptions (the latter, just

like the former, would be available from a global model). Of course it would be highly unlikely that both calibration targets would imply the same mean reversion.

Volatility smile calibration presents another challenge for the same-tenor approach. For instance, if one chooses a particular strike of the non-standard swaption to calibrate to (e.g. the fixed rate of the non-standard swap), which strike for the *standard* swaption would that correspond to?

19.4.3 Representative Swaption Approach

The idea of the representative swaption approach is to choose a standard swap that approximates the non-standard swap in some reasonable sense, and then to calibrate the Bermudan swaption model to the market-implied volatilities of swaptions on these standard swaps, one per exercise date.

One can define a “representative” swap in many ways. For example, a fairly simple *PVBP matching* method chooses the standard swap whose PVBP (Present Value of a Basis Point, see Section 5.5) matches the PVBP of the non-standard swap most closely. In this case, for expiry T_n , we would choose the tenor μ_n of the standard swap by

$$\mu_n = \operatorname{argmin}_m \left\{ \left| R_n A_{n,m}(0) - \tilde{A}_n(0) \right| \right\} \quad (19.7)$$

(note that we use the notional of the first period of the swap \tilde{U}_n to scale the PVBP of the standard swap). While somewhat simplistic, the method actually turns out to be reasonably robust for some non-standard Bermudan swaptions. We proceed to improve it and make more rigorous, which will also help us identify not just the right tenor, but also the right strike for the standard calibration swaptions.

We work in the context of a one-factor Gaussian model to demonstrate the main idea, although the method is not tied to a particular model. Let us fix a start date T_n and note that the value of (any) swap at time T_n is a function of the Gaussian short rate state $x = x(T_n)$ on that date. Let $\tilde{U}_n(x)$ be the value of the non-standard swap $\tilde{U}_n(T_n)$, as a function of the short rate state. Note that $\tilde{U}_n(x)$ depends on the mean reversion, but not the volatility parameter of the model, as follows from bond reconstruction formulas (Proposition 10.1.7). Define $V(x; R, q, m)$ to be the value of a standard swap starting on T_n as a function of x , with constant notional R , fixed rate q , and covering m periods. In departure from our normal conventions, we here allow m to be any real number and not necessarily an integer; we interpret a value of, say, $m = 5.3$ as 5 full periods plus three tenths of the sixth period. The rationale for allowing for fractional periods will become clear shortly.

Now, we have three parameters that define the standard swap, R , q , and m . In the *payoff matching* method, we choose the three parameters by

matching the level, slope and curvature of the swap payoffs as functions of the state variable,

$$V(x_0; R, q, m) = \tilde{U}_n(x_0), \quad (19.8)$$

$$\frac{\partial}{\partial x} V(x_0; R, q, m) = \frac{\partial}{\partial x} \tilde{U}_n(x_0), \quad (19.9)$$

$$\frac{\partial^2}{\partial x^2} V(x_0; R, q, m) = \frac{\partial^2}{\partial x^2} \tilde{U}_n(x_0), \quad (19.10)$$

where the expansion point x_0 is the expected value of $x(T_n)$ (or close to it). In the parameterization of Section 10.1.2, it suffices to set $x_0 = 0$. The system of equations (19.8)–(19.10) is easy to solve (numerically) in the one-factor Gaussian model; let us denote the solution by R_n^*, q_n^*, m_n^* .

Even though we fix the parameters of the standard swap by local conditions around x_0 , numerical experiments show that the swap that solves (19.8)–(19.10) tends to match that of the non-standard swap across a large range of state values $x(T_n)$, suggesting considerable robustness. In addition, even though the functions V and \tilde{U}_n depend on mean reversion, numerical experiments show that the best-fit parameters R_n^*, q_n^*, m_n^* are only mildly sensitive to mean reversion.

Incorporating the payoff matching method into a volatility calibration routine is quite easy, since the choice of the best-fitting standard swaps is independent of volatility, which allows us to identify the calibration targets *before* we commence on the volatility calibration. Moreover, the strike of the calibration swaptions are produced automatically as part of the payoff matching routine, facilitating calibration to a particular point of the observed volatility smile.

Before discussing application of the representative swaption idea to accreters and amortizers, let us briefly motivate our usage of fractional swap tenors. If tenors are restricted to an integer number of periods, then a perturbation of the market data, e.g. when shifting a yield curve to calculate an interest rate delta, could potentially alter the solved-for number of periods m by plus or minus one period. Hence, restricting m to be an integer would potentially introduce a discontinuity in the calibrated model parameter — and therefore in the value of the Bermudan swaption — as a function of market data. Such a discontinuity would be purely artificial and, as explained at length in Chapter 23, highly undesirable for stability of risk sensitivities. By allowing fractional tenors, we eliminate these problems. Of course, to complete the volatility calibration, we need to know implied volatilities of swaptions with fractional tenors, but these could be obtained by (smoothly!) interpolating implied volatilities of swaptions with integer-valued tenors.

Now, let us see how the representative swaption method works for an amortizing Bermudan swaption, i.e. a Bermudan swaption with decreasing notional $R_1 \geq R_2 \geq \dots \geq R_{N-1}$. We note that, according to (19.4), all weights $v_{n,m}$ in the decomposition of the amortizing swap into standard

swaps are positive, $v_{n,m} \geq 0$, $m = 1, \dots, N - n$. It follows that both the PVBP matching and payoff matching methods produce a standard swap whose tenor is some average of tenors of the standard swaps in the basket; in particular, the resulting standard swap have a final maturity that is *shorter* than the amortizing swap. This is an intuitive result, and leads to an amortizing Bermudan swaption being sensitive to interest rate volatilities of standard swaps V_{n,m_n^*} , $n = 1, \dots, N - 1$, with $n + m_n^* \leq N$ for any $n = 1, \dots, N - 1$.

The situation is different for accreting Bermudan swaptions, i.e. when $R_1 < R_2 < \dots < R_{N-1}$. According to (19.4),

$$\begin{aligned} v_{n,m} &< 0, & m &= 1, \dots, N - n - 1, \\ v_{n,m} &> 0, & m &= N - n. \end{aligned}$$

So, an accreting swap decomposes into a standard swap of the maximum tenor $N - n$ and maximum notional R_{N-1} , *minus* a basket of swaps of smaller tenors. The PVBP $\tilde{A}_n(0)$ of an accreting swap is larger than the PVBP of a standard swap of the matching tenor (times starting notional R_n) so the PVBP matching method would calculate an optimal tenor m^* that is *longer* than the tenor of the amortizing swap, $m^* > N - n$. The same would be true of the payoff matching method as well. This is, of course, rather problematic, as our calibration method would suggest that an accreting Bermudan swaption is sensitive to volatilities of swaptions with total length (the sum of expiry and tenor) exceeding the final maturity of the Bermudan swaption. This would be in direct contradiction to what, say, a globally-calibrated LM model would suggest, as in the latter the price of any derivative is fully determined by the volatility structure of Libor rates that fix before the final maturity of the derivative, and this volatility structure, in turn, is fully determined by the volatilities of swaptions with total length less than the final maturity.

The reader may ask why we are getting reasonable results for amortizing swaptions and unreasonable ones for accreting swaptions. A bit of reflection reveals that the discrepancy originates with the single-factor assumption that we made implicitly in the PVBP matching method, and explicitly in the payoff matching method. For the amortizing swap, the decomposition resulted in a basket of standard swaps with positive weights, a basket that can be reasonably well-hedged with a single swap of average tenor. In the accreting case, our decomposition resulted in a *spread* position in standard swaps: long a long-tenor swap and short a basket of short-tenor swaps. Our one-factor methods suggest hedging this spread position with a single (very) long-dated swap — perfectly reasonable in a one-factor world, but not in actual reality.

While sensibly hedging the spread position in an accreting swap with a single standard swap is not possible, things improve markedly if we allow usage of *two* standard swaps in the hedge. In particular, we may then take

as one of our hedges a long position in the T_N -maturity swap with maximum (R_{N-1}) notional, and construct the second swap hedge by PVBP (or payoff) matching the remaining short basket of swaps $\sum_{m=1}^{N-n-1} v_{n,m} V_{n,m}$. As all weights in the short basket are of the same (negative) sign, a single standard swap would often provide a good hedge. Thus, to get a reasonable calibration scheme for an accreting Bermudan swaption, we would need to calibrate to two standard European swaptions per expiry, both of which would have their final payment date on or before the final payment date of the Bermudan swaption. Of course, we would find it difficult to accurately calibrate a one-factor Markovian model to two swaptions per expiry, and would likely need to move on to more elaborate models with additional factors.

In conclusion, we note that the two-swaps approach works universally for Bermudan swaptions with arbitrary notional schedules. To apply it, for each n we would combine swaps with positive weights $v_{n,m}$ into one basket, and swaps with negative weights $v_{n,m}$ into another basket. Then we would represent each basket by one standard swap by the procedures discussed above, yielding the calibration swaption targets for that expiry.

19.4.4 Basket Approach

The discussion above suggests that the pricing of at least some Bermudan swaptions with non-standard notional schedules is best done in multi-factor models. Still, one-factor Markovian models are highly popular due to their performance advantages, and the desire to use them even in situations where they might be overstretched is often considerable. Consequently, rather creative ways of using one-factor models for non-standard Bermudan swaptions have been developed, resulting in a family of approaches that we here all categorize as *basket methods*.

The basket methods generally split the valuation of a non-standard Bermudan swaption into two stages. During the first stage, some model is used to calculate values of core non-standard European swaptions. During the second stage, a one-factor model is calibrated to these values of non-standard core European swaptions, and subsequently used to compute the value of the non-standard Bermudan swaption. Various method differ in how the values of non-standard European swaptions are calculated. One perfectly sound method uses a globally calibrated model such as the LM model for the task, resulting in a *local projection method for non-standard Bermudan swaptions*. We have discussed the local projection method in various flavors often enough, so we trust the reader with filling in remaining details. Instead, we review some alternatives for how to execute the first stage of the basket method.

For concreteness, let us focus on the first non-standard European swaption underlying the Bermudan swaption, i.e. the option expiring at T_1 on a swap that covers $N - 1$ periods. As follows from (19.4), this non-vanilla European swaption can be interpreted as an option on a basket of standard swaps,

all starting on T_1 but with different maturities. Hence, to compute the price of the non-standard swaption, one can use a model calibrated to the volatilities of options on such swaps, as well as relevant swap rate correlations. Notice that the standard swaptions involved here all form a “row” of the swaption grid (see Section 5.10), as they all share the same expiry but have different tenors. A one-factor mean-reverting Gaussian (or quasi-Gaussian or quadratic Gaussian) model can be calibrated to this set of European swaptions, although the calibration will be different from our standard procedure. In particular, the prices of all swaption targets depend on the model volatility function over the *same* interval $[0, T_1]$ and, thus, the short rate volatility function (e.g. $\sigma_r(\cdot)$ in the notation of Proposition 10.1.7) cannot be used if we want to match each swaption volatility exactly. Upon reflection, it should be clear that we instead can use the time-dependent mean reversion function ($\kappa(\cdot)$ in the notation of Proposition 10.1.7) as our main calibration “knob”, since the pricing of a T_1 -expiry swaption on a swap that covers m periods will depend on the mean reversion function over the period $[T_1, T_{1+m}]$. Hence, a sequential mean reversion calibration is possible: after calibrating to the first m standard swaptions, the $(m+1)$ -th is matched by changing the mean reversion function⁶ over the time interval $t \in [T_{m+1}, T_{m+2}]$, for $m = 0, \dots, N-2$. The (constant) level of volatility over the first period could be set arbitrarily, as its scaling effect would be compensated by the mean reversion calibration. However, it is advisable to keep it at a “typical” value of, say, 1% so the calibrated mean reverisions would also remain in a “typical” region, as the numerical implementation of the model might not cope well with extreme values of mean reversion.

To summarize, the basket method for a mean-reverting one-factor short rate model works like this. First, for each row of the swaption grid that corresponds to an exercise date T_n , $n = 1, \dots, N-1$, of the non-vanilla Bermudan swaption, we fit separate instances of the one-factor model by sequentially calibrating the mean reversion function to all relevant swaptions in the row. For each T_n , the relevant instance of the model is then used to compute the price of the T_n -expiry non-vanilla European swaption that the Bermudan swaption can be exercised into. Finally, we calibrate the model once more by setting its short rate volatility function to match the prices of the non-standard European swaptions established in the previous step. In the final calibration, we would typically keep the mean reversion fixed, either at a user-specified level or (ideally) at a level that makes inter-temporal correlations of core swap rates match those coming from a global model; see Section 19.2 for additional discussion.

So far we have side-stepped the issue of what strikes to choose for various calibrations in the basket scheme above. It is fair to say that this choice is a non-trivial problem. Various ad-hoc schemes could be imagined, such as using

⁶It is probably advisable to impose smoothness constraints on the time-dependent mean reversion function while performing such calibration.

the standard swaption of the same relative moneyness as the non-standard one, but they are rarely entirely satisfactory. Using a quasi-Gaussian or quadratic Gaussian model (or, even better, a stochastic volatility extension of these models) is obviously preferable to, say, using a Gaussian model, as the former models will allow us to calibrate to the volatility smile at more than a single strike, thereby alleviating somewhat the strike selection problem.

Another issue that we should touch on concerns correlations. By using a one-factor model for establishing non-standard European swaption values, we are implicitly assuming high⁷ correlations between standard swaps in the portfolio that replicates the exercise value of the swaption. This is not necessarily as constraining as it might appear to be, as in reality these swaps are indeed highly correlated. Still, we may want to contemplate methods to somehow incorporate into our procedure observations about correlation extracted, say, from historical analysis. One possible route for this would be to apply approaches inspired by basket valuation methods from equities modeling. For example, we could (rather crudely) value a non-standard European swaption by the Black formula on the (non-standard) swap rate whose volatility is obtained by moment matching⁸. This method would need to approximate the weights in the decomposition (19.5) as being deterministic (although they are not). We can also use more advanced approaches, such as the copula methods, or even SV methods, that we developed for multi-rate derivatives in Chapter 17, allowing us to incorporate volatility smile information into the basket valuation.

A few final comments on the method developed in this section are in order. First, it is worth pointing out that most of the basket methods are consistent with the way standard Bermudan swaptions are valued. Specifically, if we apply these methods to standard Bermudan swaptions, we obtain the same price as if we had used the “standard” valuation method of Section 19.2. Second, notice that the method produces volatility sensitivities for non-standard Bermudan (and European, for that matter) swaptions that tend to be intuitive and reasonable. In particular, each underlying European non-standard swaption will show sensitivities only to the correct row of the swaption grid and to ordinary swaptions on swaps with maturities that do not exceed the maturity of the swap in the corresponding non-standard swaption. For an accreting European swaption in particular with, say, an expiry T_n and swap maturity T_N , the sensitivities will be negative for all standard swaptions with expiry T_n and swap maturities T_i , $i = n + 1, \dots, N - 1$, and will be positive for a swaption with expiry T_n and swap maturity T_N , thus faithfully representing the accreting swap rate as a “spread”.

⁷Term correlations between swap rates in one-factor models are not exactly 100% due to time-dependence in parameters and presence of mean reversion.

⁸Appendix 19.B gives a quick tour of the classical moment matching ideas.

19.4.5 Super-Replication for Non-Standard Bermudan Swaptions

The replication method of Proposition 8.4.13 links the value of a European option with a non-standard payoff to that of a portfolio of standard European options. Not only does that give us a way to value a non-standard derivative, it also allows us to fully hedge it in a model-independent way. Such static replication results are extremely convenient, and much research in derivatives pricing theory have been directed towards the search for static hedges for exotic derivatives, see e.g. Andersen et al. [2002]. Unfortunately, the availability of truly model-free static replication methods for non-European options is an exceedingly rare phenomenon and no such results are known to exist for Bermudan swaptions. Interestingly, however, Bermudan swaptions with non-standard notional schedules can be *super-replicated* in a model-independent way, in the sense that for any given non-standard Bermudan swaption we can find a portfolio of standard Bermudan swaptions that would dominate the value of the non-standard Bermudan swaption in all states of the world. Moreover, in some cases the difference in value between the non-standard Bermudan swaption and its super-replicating portfolio can be quite small. While not as convenient as a replicating portfolio, the super-replicating portfolio has several practical uses. First, the value of the portfolio provides a hard no-arbitrage bound for the value of the non-standard Bermudan swaption, and any modeling procedure (including calibration rules, choice of strikes, etc.) can be checked against this bound. Second, if the upper bound provided by the portfolio is known to be relatively tight (as is the case for, e.g., amortizing Bermudan swaptions, as we shall see shortly), then the super-replicating portfolio can be used directly for valuation purposes, perhaps amended with a small ad-hoc adjustment. More importantly, the super-replicating portfolio can be used as a robust hedge that requires little rebalancing over time.

The easiest way to demonstrate the construction of a super-replicating portfolio is by example. Consider first an amortizing Bermudan swaption. For concreteness, assume it is a 10 year (10y) Bermudan swaption with exercises every year, starting in year 1. Suppose the initial notional is 10 and it decreases by 1 every year. The super-replicating portfolio then consists of nine standard Bermudan swaptions, all with unit notional: a 10 no-call 1, a 9 no-call 1, . . . , a 3 no-call 1, and a 2 no-call 1. To see that this is indeed a super-replicating portfolio, suppose the amortizing Bermudan swaption is exercised at year 5. Then the option holder receives an amortizing 5 year swap, with a starting notional of 5 that decreases by 1 every year; the value of this swap is equal to the sum of standard swaps of tenors 5y, 4y, . . . , 1y, each with notional 1. But clearly the value of each of these standard swaps is dominated by the value of a corresponding standard Bermudan swaption in the super-replicating basket.

Another way to test that the super-replicating strategy dominates the amortizing Bermudan swaption is to impose a particular exercise strategy on

the portfolio of standard Bermudan swaptions. Specifically, we simply require that all still-alive (i.e. non-expired) standard Bermudan swaptions shall be exercised at the same time as when the amortizing Bermudan swaption is exercised. A little thought shows that the exercise value obtained from the super-replicating basket is then exactly the same as from the amortizing Bermudan swaption. Hence, the value of the portfolio of standard Bermudan swaptions with this specific exercise rule enforced must precisely equal the value of the amortizing Bermudan swaption. However, as the chosen exercise strategy will generally be sub-optimal for each of the standard Bermudan swaptions in the portfolio, the true value (i.e. the value obtained with optimal exercise) of the basket of standard Bermudan swaptions will be higher than the amortizing Bermudan swaption, and will dominate its value in all states of the world.

To give another example, consider a 10 year accreting Bermudan swaption with an initial notional of 1 that increases by 1 every year. Assuming that exercise can take place annually starting in year 1, the super-replicating portfolio will now consist of a collection of nine 10 year standard unit notional Bermudan swaptions exercisable annually, with lockout periods of 1, 2, . . . , 9 years, respectively. To check that this hedge works, let us again assume the accreter is exercised in year 5. Then the holder would receive a swap that can be decomposed into a spot-starting 5 year standard swap, a 4 year standard swap starting in 1y, a 3 year standard swap starting in 2 years and so on, all with notional 1. Again, the value of each of these swaps is dominated by a corresponding standard Bermudan swaption in the super-replicating portfolio.

Super-replicating portfolios for Bermudan swaptions with arbitrary notional schedule always exist, but are rather tedious to write down explicitly (for example, see a related algorithm in Evers and Jamshidian [2005]). The basic idea, however, is quite simple: for any exercise opportunity one needs to ensure that each of the standard swaps into which the exercise value can be decomposed (see (19.4)) is matched by a standard Bermudan swaption in the super-replicating portfolio.

Let us show some numerical results as a way to examine the tightness of the upper value bound produced by the super-replicating portfolio. For the numerical experiments, we throughout use a one-factor quasi-Gaussian model with some reasonable, representative parameters and a yield curve flat at 6%. In Table 19.1 we show values for Bermudan (European) amortizing swaptions of different maturities (with 1y lockout) against the values of corresponding super-replicating portfolios of standard Bermudan (European) swaptions; all contracts are receivers (options on receive-fixed swaps) with 6% strike. The notionals of amortizing swaps decrease linearly from the initial notional indicated in the table by 1 every year. We notice that the upper bounds produced by the super-replicating portfolio are here quite tight, something that appears to hold generally for amortizing Bermudan swaptions.

Maturity	10y	10y	30y	30y
Initial notional	10	10	30	30
Bermudan/European	E	B	E	B
Amortizer value	0.606	0.630	2.180	2.700
Portfolio value	0.614	0.650	2.230	2.830

Table 19.1. The value of an amortizing Bermudan or European swaption vs. the value of a super-replicating portfolio of standard Bermudan/European swaptions for different maturities and contract types.

As a second test we look at a particular amortizing Bermudan swaption across a range of strikes, see Table 19.2. We consider a 30 year amortizing Bermudan swaption with the initial notional of 30, and compare it to the super-replicating portfolio. For reference, the vega (change in Bermudan swaption value to 1% change of Black volatilities of European swaptions) for the 6% Bermudan swaption is about 0.1, with European swaption implied volatilities around 12%. Again, the results from the super-replicating portfolio are quite close to the real option values.

Strike	3%	4%	5%	6%	7%	8%
Amortizer value	0.267	0.632	1.372	2.700	4.640	7.017
Portfolio value	0.355	0.766	1.534	2.830	4.710	7.052

Table 19.2. The value of a 30 year amortizing Bermudan swaption vs. the value of a super-replicating portfolio of standard Bermudan swaptions for different strikes.

Next, we look at the results for accreting swaptions. In Table 19.3 we have the results for accreting receivers with 6% strike with notional accreting at 6% relative rate, across different contract types and maturities. Clearly, the super-replicating portfolio here is substantially more expensive than the accreting Bermudan swaption.

Maturity	10y	10y	30y	30y
Initial notional	1	1	1	1
Bermudan/European	E	B	E	B
Accreter value	0.096	0.113	0.121	0.255
Portfolio value	0.123	0.143	0.280	0.380

Table 19.3. The value of an accreting Bermudan or European swaption vs. the value of a super-replicating portfolio of standard Bermudan/European swaptions for different maturities and contract types.

Finally, we look at a particular 30 year accreting Bermudan swaption with initial notional of 1 across different strikes, see Table 19.4. For each contract, the notional compounds at the rate given by the fixed rate.

Strike	3%	4%	5%	6%	7%	8%
Accreter value	0.021	0.053	0.121	0.255	0.495	0.880
Portfolio value	0.021	0.063	0.164	0.380	0.773	1.400

Table 19.4. The value of a 30 year accreting Bermudan swaption vs. the value of a super-replicating portfolio of standard Bermudan swaptions for different strikes.

Judging by Tables 19.1–19.4, it appears that the super-replicating portfolio tends to produce a tighter bound for amortizers than for accreters. This observation, as it turns out, is not tied to the particular structures and market data used in the tables, but is true in general. To understand why this is the case, recall that the super-replicating portfolio for a non-standard Bermudan swaption will have the exact same value as the non-standard Bermudan swaption if each standard Bermudan swaption in the portfolio is exercised at the same time as the non-standard Bermudan swaption. In other words, the tightness of the upper bound for amortizers therefore suggests that it is optimal to exercise all standard Bermudan swaptions in the super-replicating portfolio at about the same time (for an amortizing Bermudan swaption, we recall that this portfolio consists of Bermudan swaptions with identical lockout periods, but different maturities). However, according to Proposition 19.7.1 proven later in this chapter, arbitrage arguments can be used to show that, at any exercise date, if one does not exercise the standard Bermudan swaption with the shortest tenor among remaining in the basket (i.e. with a remaining 1 year swap in the example above), then one should never exercise *any* of the remaining standard Bermudan swaptions (with tenors 2,3,... years). In light of this result, the tightness of the super-replication bound is therefore not surprising. The same argument does *not* hold for the accreters, because for these structures the super-replicating portfolio consists of standard Bermudan swaptions with different lockout periods, rather than different underlying swap tenors. As such, it is obviously not reasonable to assume that the standard Bermudans will be optimally exercised at the same time.

Finally, let us briefly comment on the *lower* bound for the value of a non-standard Bermudan swaption. While we know of no general results, the carry argument of Section 19.7.2 allows us to show that an amortizing Bermudan swaption with a final notional (i.e. the notional for the final exercise date) of 1 is bounded from below by the standard Bermudan swaption of notional 1, as long as we keep the exercise schedule, strike, etc. unchanged. This holds

in a model-independent way. Unfortunately, this lower bound is typically quite loose, except for Bermudan swaptions that amortize slowly.

19.4.6 Zero-Coupon Bermudan Swaptions

Let us momentarily turn to the question of where accreting Bermudan swaptions come from in the first place. Consider an accreting (receiver) Bermudan swaption with the notional defined by

$$R_i = \prod_{j=0}^{i-1} (1 + \tau_j k), \quad i = 1, \dots, N-1. \quad (19.11)$$

According to (19.3), its n -th exercise value at time T_n is given by

$$\begin{aligned} \tilde{U}_n(T_n) &= \sum_{i=n}^{N-1} R_i \tau_i P(T_n, T_{i+1}) (k - L_i(T_n)) \\ &= \sum_{i=n}^{N-1} \left(\prod_{j=0}^{i-1} (1 + \tau_j k) \right) (P(T_n, T_{i+1}) (1 + \tau_i k) - P(T_n, T_i)) \\ &= \sum_{i=n}^{N-1} (R_{i+1} P(T_n, T_{i+1}) - R_i P(T_n, T_i)) \\ &= R_N P(T_n, T_N) - R_n, \end{aligned} \quad (19.12)$$

where we have used the defining relation $L_i(T_n) = (P(T_n, T_i) - P(T_n, T_{i+1}) / (\tau_i P(T_n, T_{i+1}))$.

Now consider a contract in which an investor gives⁹ the dealer 1 at time T_0 , and the dealer promises to pay the investor the amount of $R_N = \prod_{j=0}^{N-1} (1 + \tau_j k)$ at T_N . The payment of R_N at T_N can be seen as the value of the original investment compounded at the fixed rate k over the time period $[T_0, T_N]$. The contract is essentially a zero-coupon bond with a discretely compounding rate of k . Suppose now that the dealer is granted a Bermudan-style option to cancel the zero-coupon note at any time T_n , $n = 1, \dots, N-1$, in return for paying the investor the accumulated amount to that date, i.e. R_n . For reasons that should be obvious, the embedded Bermudan option is called a *zero-coupon Bermudan swaption*. The payoff to the dealer upon exercise at time T_n of this option is evidently equal to i) an immediate outflow of R_n ; and ii) release from the obligation to pay R_N at time T_N , the value of which is $P(T_n, T_N) R_N$. In other words, the exercise value of a zero-coupon Bermudan swaption, $U_{ZC,n}(T_n)$, is equal to

$$U_{ZC,n}(T_n) = -R_n + R_N P(T_n, T_N),$$

⁹Or, equivalently, pays 1 at time T_N in addition to running Libor coupons on a unit notional.

which allows us to deduce from (19.12) that

$$\tilde{U}_n(T_n) = U_{ZC,n}(T_n), \quad n = 1, \dots, N - 1.$$

Therefore, the value of a zero-coupon Bermudan swaption with rate k is equal to the value of an accreting Bermudan swaption with the notional accreting at rate k , as in (19.11). Our earlier discussion of accreting Bermudan swaptions therefore holds unchanged for zero-coupon Bermudan swaptions.

19.4.7 American Swaptions

Having a time-varying notional schedule is not the only non-standard feature that can be attached to Bermudan swaptions. One relatively common deviation from the standard contract permits the option holder to exercise *on any business date*, after a lockout period. Not surprisingly, such swaptions are called *American swaptions*. These are fairly popular in the US as hedges for mortgage bonds, presumably because the American exercise feature might be considered a better hedge for the prepayment behavior of mortgage borrowers.

The coupon-paying nature of the underlying swap makes the definition of an American swaption somewhat complicated. If exercise takes place during a period $[T_n, T_{n+1}]$, then the option holder receives the swap starting at T_{n+1} , i.e. $U_{n+1}(\cdot)$, as well as an “exercise fee” equal to the difference between the Libor rate effective for the period $[T_n, T_{n+1}]$ (i.e. $L_n(T_n)$) and the fixed rate, times the notional and times the remaining time to T_{n+1} (in the appropriate day count convention). In mathematical notations, the exercise value per unit notional at t , $t \in (T_n, T_{n+1}]$, is given by

$$U_n^A(t) = (L_n(T_n) - k)(T_{n+1} - t) + U_{n+1}(t).$$

We emphasize that, as a rule, the time t fee is set to the “accrued current coupon” $(L_n(T_n) - k)(T_{n+1} - t)$, *not* its discounted value from t to the payment date T_{n+1} . This choice, as is true of many others related to contract specification, is made by those who write term sheets (documents outlining details of derivatives contracts) rather than by those responsible for valuation algorithms, and implies that the exercise value will be discontinuous in time,

$$U_n^A(T_n+) \neq U_n^A(T_n) = U_n(T_n). \quad (19.13)$$

Odd as it is, this discontinuity is not the main issue with American swaptions. From a valuation standpoint the biggest problem with American swaptions is the fact that the exercise value at t (if not equal to one of the coupon dates) depends on the value of the Libor rate at $T_n < t$ and is thus *path-dependent*.

For Monte Carlo based valuation methods, the path dependence is not a problem, as the value of the Libor rate is known on each path when estimating the exercise boundary or using it in valuation. In a PDE setting, matters

are more complicated. Before describing possible methods, let us comment on the somewhat prevalent view that one can approximate an American swaption with a Bermudan swaption with high frequency of exercise dates.

19.4.7.1 American Swaptions vs. High-Frequency Bermudan Swaptions

Let us choose a particular period $[T_n, T_{n+1}]$ and consider it subdivided into M periods

$$T_n = t_0 < t_1 < \dots < t_M = T_{n+1}.$$

Then, consider a Bermudan swaption with exercises at $\{t_i\}_{i=1}^M$, versus an American swaption that can be exercised on the same dates. Also, for simplicity assume that the “exercise fee” for an American is in fact, properly discounted, such that the discontinuity in (19.13) is removed. Concentrating only on the exercise value contributions $u_n^A(t)$, $u_n^B(t)$ paid in the period $[T_n, T_{n+1}]$, so that

$$U_n^A(t_m) = u_n^A(t_m) + U_{n+1}(t_m), \quad U_n^B(t_m) = u_n^B(t_m) + U_{n+1}(t_m),$$

a standard Bermudan swaption exercised at t_m gives the holder an exercise value

$$u_n^B(t_m) = \sum_{i=m}^{M-1} (t_{i+1} - t_i) P(t_m, t_{i+1}) \left(\frac{P(t_m, t_i) - P(t_m, t_{i+1})}{(t_{i+1} - t_i) P(t_m, t_{i+1})} - k \right).$$

Here

$$\frac{P(t_m, t_i) - P(t_m, t_{i+1})}{(t_{i+1} - t_i) P(t_m, t_{i+1})}$$

is just the time- t_m forward Libor rate for the period $[t_i, t_{i+1}]$. Using the standard rearrangement of terms, we obtain

$$\begin{aligned} u_n^B(t_m) &= L(t_m, t_m, T_{n+1}) (T_{n+1} - t_m) P(t_m, T_{n+1}) \\ &\quad - k \sum_{i=m}^{M-1} (t_{i+1} - t_i) P(t_m, t_{i+1}), \end{aligned} \quad (19.14)$$

with

$$L(t_m, t_m, T_{n+1}) = \frac{1 - P(t_m, T_{n+1})}{(T_{n+1} - t_m) P(t_m, T_{n+1})},$$

where we have used the full notation (4.2) for forward Libor rates. The first term represents payment of a Libor rate covering the period $[t_m, T_{n+1}]$, and the second is a collection of fixed-rate payments on a schedule $\{t_m, \dots, t_M = T_{n+1}\}$.

Let us contrast this with the exercise value of a true American swaption exercised at t_m ,

$$u_n^A(t_m) = L(T_n, T_n, T_{n+1}) (T_{n+1} - t_m) P(t_m, T_{n+1}) - k(T_{n+1} - t_m) P(t_m, T_{n+1}), \quad (19.15)$$

with

$$L(T_n, T_n, T_{n+1}) = \frac{1 - P(T_n, T_{n+1})}{(T_{n+1} - T_n) P(T_n, T_{n+1})}.$$

Clearly, there are two differences between u_n^B and u_n^A . The first one is the difference in the fixed leg, the payment of the annuity $k \sum_{i=m}^{M-1} (t_{i+1} - t_i) P(t_m, t_{i+1})$ versus a single bullet payment $k(T_{n+1} - t_m) P(t_m, T_{n+1})$. The difference is in the timing of discounting and is normally small, as $T_{n+1} - T_n$ would often be equal to 3 months in the US. Even if one deems this to be an issue, any perceptible difference could be eliminated almost fully by imposing appropriately chosen deterministic exercise fees. The second difference, on the other hand, is much greater, and concerns the difference of which Libor rate is applied to the period $[t_m, T_{n+1}]$. For the Bermudan swaption, it is $L(t_m, t_m, T_{n+1})$ and for the American, it is $L(T_n, T_n, T_{n+1})$. Not only will the two Libor rates have different forward values, their different fixing dates (t_m vs. T_n) will affect the amount of volatility each Libor rate experiences over its lifetime. These effects can yield quite significant valuation differences, especially for steeper yield curves and shorter maturities.

19.4.7.2 The Proxy Libor Rate Method

Besides highlighting the fallacy of using a high-frequency Bermudan swaption as a proxy for an American, the analysis above also hints at proper remediation. Indeed, it should be clear that the approximation of an American swaption with a Bermudan swaption suffers the most not from the mismatched exercise frequency, but from the difference in the exercise values in the two contracts, as the Bermudan approximation uses a Libor rate that has the wrong forward value and the wrong volatility. The idea behind the *proxy Libor rate method* involves correcting the forward value/volatility as appropriate, while removing the path-dependence of exercise that hinders the application of backward-induction methods.

Continuing with the notations of the previous section, we define the proxy Libor rate $\tilde{L}(t_m, t_m, T_{n+1})$ by

$$\begin{aligned} \tilde{L}(t_m, t_m, T_{n+1}) &= L(0, T_n, T_{n+1}) \\ &+ \frac{\text{Stdev}(L(T_n, T_n, T_{n+1}))}{\text{Stdev}(L(t_m, t_m, T_{n+1}))} (L(t_m, t_m, T_{n+1}) - L(0, t_m, T_{n+1})). \end{aligned} \quad (19.16)$$

Here $\text{St.dev}(X)$ is defined as the Normal (or basis-point, see Remark 7.2.9) term volatility of the rate X , which we may compute from any particular model used for American swaption valuation. The proxy Libor rate enjoys the following properties.

- Its expected value under the T_{n+1} -forward measure is equal to the forward $L(0, T_n, T_{n+1})$, i.e. the forward of the rate used in the “real” American swaption.
- Its (term) volatility is equal to that of $L(T_n, T_n, T_{n+1})$.
- It is a function of the yield curve at time t_m and, unlike for $L(T_n, T_n, T_{n+1})$, its value is available in a backward-induction scheme at t_m .

Having defined the proxy Libor rate, we define a Bermudan swaption whose exercise value approximates the exercise value of an American swaption (compare to (19.14)),

$$u_n^B(t_m) \approx \tilde{L}(t_m, t_m, T_{n+1})(T_{n+1} - t_m) P(t_m, T_{n+1}) \\ - k \sum_{i=m}^{M-1} (t_{i+1} - t_i) P(t_m, t_{i+1}).$$

Assuming that the underlying model is Markovian and low-dimensional, a Bermudan swaption with these exercise values can easily be evaluated in a finite difference lattice, and is a close proxy for the “real” American swaption.

19.4.7.3 The Libor-as-Extra-State Method

While having an accurate approximation is a good step forward, it may still be of use to have an exact lattice-based valuation method for American swaptions, especially when assessing the accuracy of various approximations. One approach for this is the “extra state variable” method, where path-dependence is dealt with by back-propagating values of the security in all possible states of the path-dependent state variable, and then applying update conditions between different “slices”, see Sections 2.7.5 and 18.4.5. For our purposes here, we define the extra state variable by

$$I(t) = \sum_{n=0}^{N-1} 1_{\{t \in [T_n, T_{n+1})\}} L(T_n, T_n, T_{n+1}).$$

Clearly, the payoff of an American swaption at time t could be expressed as a function of the yield curve observed at t , and the state variable $I(t)$. To elaborate, let $H_n(t, I)$ be the value of an American swaption with exercise dates after T_n at time t given $I(t) = I$. Focusing on the period $[T_n, T_{n+1})$ and assuming the exercise dates fall on the grid $\{t_m\}$ as defined in Section 19.4.7.1, we have the following recursion,

$$H_n(t_m, I) \\ = E_{t_m} \left(\frac{B(t_m)}{B(t_{m+1})} \max(u_n^A(t_{m+1}, I) + U_{n+1}(t_{m+1}), H_n(t_{m+1}, I)) \right)$$

for $m = M - 1, \dots, 0$, where we have defined (see (19.15))

$$u_n^A(t_m, I) = (I - k)(T_{n+1} - t_m) P(t_m, T_{n+1}).$$

As $I(T_n) = L(T_n, T_n, T_{n+1})$ is a function of the state variables of the model at T_n , we have that

$$H_n(T_n) = H_n(t_0, L(T_n, T_n, T_{n+1}))$$

is independent of I and only a function of the state variables of the model at T_n . We then use $H_n(T_n)$ to start a recursion for the $(n - 1)$ period. The full American value is given by $H_0(0)$.

19.4.8 Mid-Coupon Exercise

Half-way between standard Bermudan swaptions and American swaptions lie Bermudan swaptions that allow exercises on the standard tenor dates $\{T_n\}$ plus a select few — often just one — extra dates per coupon period. At this point, we should note that while we so far for convenience have assumed that fixed and floating payments take place on the same schedule, in reality swaps in many currencies (including USD and EUR) pay floating coupons more frequently than fixed coupons. Taking the US as an example, standard conventions specify that floating rate payments occur every three months (and are linked to three-month Libor rates), while fixed rate payments are made every six months. Exercise dates are often chosen to coincide with floating rate fixing dates, i.e. are spaced three months apart. Having an exercise take place in the middle of a fixed-rate coupon period is not a problem, however, as the value of the remaining part of the fixed-rate coupon is trivial to estimate on the exercise date¹⁰. Less common are exercise dates in the middle of a floating-rate coupon period. For such contracts we face the same issue as with American swaptions: the exercise value on the exercise date is path-dependent, as it is linked to a fixing of the Libor rate that occurs prior to the exercise date.

PDE pricing of structures with exercise taking place inside a floating rate period involves the same issues as those discussed in Section 19.4.7 for American swaptions, and remediation follows the same path. As an approximation we can use an expression like (19.16) to replace the Libor rate with a proxy Libor rate setting on the exercise date, with the proxy rate constructed to have the same forward value and volatility as the real rate. For exact valuation, the state-variable approach of Section 19.4.7 can be used.

¹⁰We note that exercise in the middle of a fixed-rate period is most often accompanied by an *exercise fee*, a deterministic amount of money payable upon exercise that is agreed upon in advance. The fee is typically calculated to reflect the value of the part of the fixed-rate coupon accrued from the beginning of the period to the exercise date.

19.5 Flexi-Swaps

A Bermudan swaption can be interpreted as a fixed-floating swap with zero notional and a (single) option to increase the notional to a given level on any of the exercise dates. Likewise, a cancelable swap can be seen as a swap of full notional with an option to decrease the notional to zero on any of the exercise dates; once the option is exercised, the right goes away. More flexibility in choosing swap notionals is afforded in a so-called *flexi-swap* (also known as a *chooser swap* or a *band swap*), a swap with multiple options to change the notional on a given set of exercise dates, subject to certain constraints. Flexi-swaps are related to the flexi-caps discussed in Section 2.7.6 and are most often used as hedges for so-called *balance-guarantee swaps*, i.e. swaps with a notional linked to a pool of mortgages. The ability to gradually decrease the notional in a flexi-swap on each exercise date allows the option holder to mimic (random) prepayments in the mortgage pool. We briefly consider flexi-swap valuation in this section.

With a tenor structure $\{T_n\}_{n=0}^N$ and a collection of net coupons X_n with unit notional, fixing at T_n and paying at T_{n+1} , $n = 0, \dots, N - 1$, we define a flexi-swap to be a contract that pays a net coupon $X_n R_n$ at time T_{n+1} , where the starting notional R_0 is fixed up-front, and time- T_n notional R_n is chosen by the holder of the option at time T_n (so that R_n is \mathcal{F}_{T_n} -measurable) for each $n = 1, \dots, N - 1$, subject to some constraints. The *constraint set* for the decision at time T_n may include

- Global deterministic bounds, e.g. $R_n \in [g_n^{\text{lo}}, g_n^{\text{hi}}]$.
- Local bounds that are functions of the current notional, e.g. $R_n \in [l_n^{\text{lo}}(R_{n-1}), l_n^{\text{hi}}(R_{n-1})]$.
- Bounds that are function of market data x_n (such as Libor and swap rates) at time T_n , e.g. $R_n \in [m_n^{\text{lo}}(x_n), m_n^{\text{hi}}(x_n)]$.

In a general flexi-swap the constraint set for time T_n is the intersection of the global, local and market constraint sets; let us denote this set $\mathcal{C}_n(R_{n-1}, x_n)$, so that $R_n \in \mathcal{C}_n(R_{n-1}, x_n) \subset \mathbb{R}$. It is common to require that the notional may only decrease, so that $\mathcal{C}_n(R_{n-1}, x_n) \subset [0, R_{n-1}]$. Of course, the larger the constraint set is for each date (i.e. the fewer the constraints enforced), the more expensive the flexi-swap will be.

The valuation of a flexi-swap may proceed by backward induction, while keeping track of “current notional”. To demonstrate, let $V_n(t, R)$ be the time t of the part of the flexi-swap paying strictly after T_n , given that $R_n = R$. At time T_{n-1} , the flexi-swap value must be equal to the discounted expected value of the maximum value at time T_n , with the maximum taken over all possible choices of the notional. This observation allows us to write down the backward recursion equation,

$$V_{n-1}(T_{n-1}, R) = P(T_{n-1}, T_n) X_{n-1} R + B(T_{n-1}) \mathbb{E}_{T_{n-1}} \left(B(T_n)^{-1} \max_{R' \in C_n(R, x_n)} \{V_n(T_n, R')\} \right) \quad (19.17)$$

for $n = N, \dots, 1$, with the terminal condition $V_N(T_N, R) \equiv 0$. The time 0 actual value of the flexi-swap is given by $V_0(T_0, R_0)$. The recursion (19.17) can be implemented in a PDE model by introducing an extra state variable to keep track of the current notional R , along the lines of Sections 2.7.5 and 19.4.7.3.

19.5.1 Purely Global Bounds

Using ideas similar to those from Section 19.4.5, Evers and Jamshidian [2005] demonstrate that a flexi-swap with purely global deterministic bounds can be decomposed exactly into a portfolio of Bermudan swaptions. While theoretically interesting, the replication is sometimes awkward in practice as a typical flexi-swap will decompose into hundreds of Bermudan swaptions; valuing them all one by one is rarely more efficient than just applying the recursion (19.17). On the other hand, when using a local model such as a one-factor qG model, valuing the Bermudan swaptions one by one would allow us to tailor calibration to each individual Bermudan swaption, leading to increased precision in the value of the flexi-swap. The choice of the valuation method will be dictated by the trade-off between calibration accuracy and performance.

19.5.2 Purely Local Bounds

Flexi-swaps that involve non-global constraints generally cannot be replicated with portfolios of Bermudan swaptions. However, in the practically relevant special case of *purely local bounds of scaling type* we may obtain a more efficient valuation formula involving no state variables beyond those driving the yield curve. Specifically, let us assume that only local constraints are enforced and that these are given by the current notional multiplied by lower and upper multipliers, i.e.

$$l_n^{\text{lo}}(R_{n-1}) = \lambda_n^{\text{lo}} R_{n-1}, \quad l_n^{\text{hi}}(R_{n-1}) = \lambda_n^{\text{hi}} R_{n-1}, \quad 0 \leq \lambda_n^{\text{lo}} \leq \lambda_n^{\text{hi}},$$

for $n = 1, \dots, N - 1$. To simplify the valuation method (19.17), we make the critical observation that the value of the flexi-swap scales linearly in notional, i.e.

$$V_n(T_n, R) = R V_n(T_n, 1) \quad (19.18)$$

for any T_n, R . This follows from the fact that all coupons scale linearly with notional R , as do all constraints. In particular, as there are no global constraints, our exercise decision at any time T_n is independent of the absolute size of the notional.

An important corollary to (19.18) is that on any step n , the optimal notional choice is of the “all or nothing” type (known in control theory as “bang-bang”). This follows from the fact that

$$\begin{aligned} \max_{R' \in [\lambda_n^{\text{lo}} R, \lambda_n^{\text{hi}} R]} \{V_n(T_n, R')\} &= \max_{x \in [\lambda_n^{\text{lo}}, \lambda_n^{\text{hi}}]} \{V_n(T_n, xR)\} \\ &= \max_{x \in [\lambda_n^{\text{lo}}, \lambda_n^{\text{hi}}]} \{V_n(T_n, R)x\} \end{aligned}$$

whereby the function being maximized, $V_n(T_n, R)x$, is linear in the maximization variable x . Hence the maximum is attained at the boundary of the interval,

$$\max_{R' \in [\lambda_n^{\text{lo}} R, \lambda_n^{\text{hi}} R]} \{V_n(T_n, R')\} = \max(V_n(T_n, \lambda_n^{\text{lo}} R), V_n(T_n, \lambda_n^{\text{hi}} R)). \quad (19.19)$$

We can use the two observations above to simplify the valuation algorithm. Rewriting (19.17) with the help of (19.19) we get

$$\begin{aligned} V_{n-1}(T_{n-1}, R) &= P(T_{n-1}, T_n)X_{n-1}R \\ &\quad + B(T_{n-1})\mathbb{E}_{T_{n-1}}(B(T_n)^{-1} \max(V_n(T_n, \lambda_n^{\text{lo}} R), V_n(T_n, \lambda_n^{\text{hi}} R))). \end{aligned}$$

Dividing through by R , using (19.18), and introducing the abbreviated notation $V_n(T_n, 1) = V_n(T_n)$, we obtain the valuation equation

$$\begin{aligned} V_{n-1}(T_{n-1}) &= P(T_{n-1}, T_n)X_{n-1} \\ &\quad + B(T_{n-1})\mathbb{E}_{T_{n-1}}(B(T_n)^{-1} \max(V_n(T_n)\lambda_n^{\text{lo}}, V_n(T_n)\lambda_n^{\text{hi}})). \end{aligned}$$

Clearly, if $V_n(T_n)$ is positive, then

$$\max(V_n(T_n)\lambda_n^{\text{lo}}, V_n(T_n)\lambda_n^{\text{hi}}) = V_n(T_n) \max(\lambda_n^{\text{lo}}, \lambda_n^{\text{hi}}) = V_n(T_n)\lambda_n^{\text{hi}},$$

and if $V_n(T_n)$ is negative, then

$$\max(V_n(T_n)\lambda_n^{\text{lo}}, V_n(T_n)\lambda_n^{\text{hi}}) = V_n(T_n) \min(\lambda_n^{\text{lo}}, \lambda_n^{\text{hi}}) = V_n(T_n)\lambda_n^{\text{lo}}.$$

In total, we therefore have the ultimate valuation recursion

$$\begin{aligned} V_{n-1}(T_{n-1}) &= P(T_{n-1}, T_n)X_{n-1} \\ &\quad + B(T_{n-1})\mathbb{E}_{T_{n-1}}(B(T_n)^{-1}V_n(T_n)(\lambda_n^{\text{lo}}1_{\{V_n(T_n)<0\}} + \lambda_n^{\text{hi}}1_{\{V_n(T_n)>0\}})). \end{aligned} \quad (19.20)$$

The value at time 0 is given by $R_0 V_0(T_0)$; only one PDE plane is required to calculate it, unlike for (19.17) which requires multiple R -planes. We note that the standard cancelable swap valuation recursion is recovered with $\lambda_n^{\text{lo}} = 0$, $\lambda_n^{\text{hi}} = 1$.

19.5.3 Marginal Exercise Value Decomposition

It is instructive to see what the marginal exercise value decomposition of Section 18.2.3 looks like for a flexi-swap with local bounds. We rewrite (19.20) in a slightly different way,

$$\begin{aligned} V_{n-1}(T_{n-1}) &= P(T_{n-1}, T_n) X_{n-1} \\ &\quad + B(T_{n-1}) \mathbb{E}_{T_{n-1}} \left(B(T_n)^{-1} [\lambda_n^{\text{lo}} V_n(T_n)^- + \lambda_n^{\text{hi}} V_n(T_n)^+] \right). \end{aligned}$$

Then, since

$$V_n(T_n)^+ = V_n(T_n) - V_n(T_n)^-$$

we have

$$\begin{aligned} V_{n-1}(T_{n-1}) - \lambda_n^{\text{hi}} V_n(T_{n-1}) &= P(T_{n-1}, T_n) X_{n-1} + B(T_{n-1}) \mathbb{E}_{T_{n-1}} \left(B(T_n)^{-1} (\lambda_n^{\text{hi}} - \lambda_n^{\text{lo}}) (-V_n(T_n)^-) \right) \\ &= P(T_{n-1}, T_n) X_{n-1} + B(T_{n-1}) \mathbb{E}_{T_{n-1}} \left(B(T_n)^{-1} (\lambda_n^{\text{hi}} - \lambda_n^{\text{lo}}) (-V_n(T_n))^+ \right) \end{aligned}$$

and, taking discounted expected values to time 0,

$$\begin{aligned} V_{n-1}(0) - \lambda_n^{\text{hi}} V_n(0) &= \mathbb{E} \left(B(T_n)^{-1} X_{n-1} \right) \\ &\quad + \mathbb{E} \left(B(T_n)^{-1} (\lambda_n^{\text{hi}} - \lambda_n^{\text{lo}}) (-V_n(T_n))^+ \right). \end{aligned}$$

This holds for $n = 1, \dots, N$. Weighting the n -th equality with

$$\alpha_n^{\text{hi}} \triangleq \prod_{i=1}^{n-1} \lambda_i^{\text{hi}}$$

(with $\alpha_1^{\text{hi}} = 1$), summing all terms, and observing that

$$\sum_{n=1}^N \alpha_n^{\text{hi}} (V_{n-1}(0) - \lambda_n^{\text{hi}} V_n(0)) = V_0(0) - \alpha_N^{\text{hi}} \lambda_N^{\text{hi}} V_N(0) = V_0(0),$$

we obtain

$$\begin{aligned} V_0(0) &= \sum_{n=1}^N \alpha_n^{\text{hi}} \mathbb{E} \left(B(T_n)^{-1} X_{n-1} \right) \\ &\quad + \sum_{n=1}^N \alpha_n^{\text{hi}} (\lambda_n^{\text{hi}} - \lambda_n^{\text{lo}}) \mathbb{E} \left(B(T_n)^{-1} (-V_n(T_n))^+ \right). \quad (19.21) \end{aligned}$$

More generally

$$\begin{aligned}
V_n(T_n) &= \frac{1}{\alpha_{n+1}^{\text{hi}}} \sum_{i=n+1}^N \alpha_i^{\text{hi}} \mathbb{E} (B(T_i)^{-1} X_{i-1}) \\
&\quad + \frac{1}{\alpha_{n+1}^{\text{hi}}} \sum_{i=n+1}^N \alpha_i^{\text{hi}} (\lambda_i^{\text{hi}} - \lambda_i^{\text{lo}}) \mathbb{E} (B(T_i)^{-1} (-V_i(T_i))^+).
\end{aligned} \tag{19.22}$$

19.5.4 Narrow Band Limit

The decomposition (19.21) turns out to be useful to study the flexi-swap when the notional range $|\lambda_n^{\text{hi}} - \lambda_n^{\text{lo}}|$ is small, which is often the case as clients look for cheaper means to hedge their balance-guarantee swaps (recall that narrow range implies less optionality and lower cost). Let $\epsilon = (\lambda_n^{\text{hi}} - \lambda_n^{\text{lo}})/\lambda_n^{\text{hi}}$ be small, and denote by U_n the value of all coupons fixing on or after T_n weighted by α_i^{hi} so that

$$U_n(t) = B(t) \sum_{i=n}^{N-1} \alpha_{i+1}^{\text{hi}} \mathbb{E}_t (B(T_{i+1})^{-1} X_i).$$

This is the value of the portion of the (amortizing) swap after T_n assuming that on each exercise date the option holder always chooses the multiplier λ_i^{hi} . Then, it follows from (19.22) that to first order in ϵ ,

$$\alpha_{n+1}^{\text{hi}} V_n(T_n) - U_n(T_n) = O(\epsilon) \tag{19.23}$$

and we obtain from (19.21) that

$$\begin{aligned}
V_0(0) &= U_0(0) + \epsilon \sum_{n=1}^N \alpha_{n+1}^{\text{hi}} \mathbb{E} (B(T_n)^{-1} (-V_n(T_n))^+) \\
&= U_0(0) + \epsilon \sum_{n=1}^N \mathbb{E} (B(T_n)^{-1} (-U_n(T_n))^+) \\
&\quad + \epsilon \sum_{n=1}^N \mathbb{E} (B(T_n)^{-1} [\alpha_{n+1}^{\text{hi}} (-V_n(T_n))^+ - (-U_n(T_n))^+]).
\end{aligned}$$

The last line is of the second order in ϵ per (19.23), so that

$$V_0(0) = U_0(0) + \epsilon \sum_{n=1}^N \mathbb{E} (B(T_n)^{-1} (-U_n(T_n))^+) + O(\epsilon^2). \tag{19.24}$$

We recognize the terms in the sum above as European swaptions on the amortizing swap $-U_n(T_n)$, so the value of a narrow-band flexi-swap with local constraints is approximately equal to the underlying amortizing swap plus a strip of European swaptions on the remaining parts of the reverse of the underlying amortizing swap.

19.6 Monte Carlo Valuation

With our discussion so far, we have demonstrated that low-dimensional models, if appropriately calibrated, can be used effectively for Bermudan swaption valuation with PDE or tree-based methods. Still, it is sometimes useful to be able to price Bermudan swaptions with Monte Carlo methods, e.g. to compare prices computed in a low-dimensional model to those of a larger globally-calibrated model (such as the LM model). The mechanics of Monte Carlo valuation follow our discussion in Section 18.3 closely, so here we merely point out the simplifications made possible by the simpler structure of Bermudan swaptions, as compared to the general class of callable Libor exotics (CLEs).

19.6.1 Regression Methods

Bermudan swaptions can be valued by Monte Carlo simulation in straightforward fashion, using the general regression-based methods of Section 18.3. There are, however, a number of shortcuts that are worth pointing out. First, observe that the exercise values of a Bermudan swaption can be calculated directly off the yield curve at the time of exercise, whereby there is no need to use regression methods to estimate exercise values. It follows that we can use the simple algorithm of Section 18.3.1 directly.

Another advantage that Bermudan swaptions enjoy over more complex CLEs is the relative ease with which good explanatory variables can be selected. It is clear that for regressing the hold value of a Bermudan swaption on a given exercise date, the value of the underlying swap is important, suggesting that the overall level of interest rates on each exercise date — as represented by either the swap rate or the value of the swap starting on the exercise date and maturing on the final date of the Bermudan swaption — should always be included in the set of explanatory variables. The slope of the yield curve on each of the exercise dates turns out to be relevant as well, as it is actually the forward-starting swap — that is, the swap that underlies the European swaption expiring on the next exercise date — that impacts the hold value, and the difference in value between a spot-starting and a forward-starting swaps clearly originates with the slope of the yield curve. We can either include the forward starting swap as an additional explanatory variable on each exercise date or, better yet, include the spot Libor rate for the next period on each exercise date. The latter suggestion achieves nearly the same result as the spot-starting swap could be decomposed into an FRA for the next period and a forward-starting swap. Note that empirical evidence shows that it is not advisable to use the forward-starting swap alone as the sole explanatory variable per exercise date as it appears both the level and the slope of the yield curve should be represented in the set, especially in the setting of a multi-factor model. We investigate and compare various concrete exercise strategies in Section 19.6.2 below.

Another observation that can be fruitfully explored in the LS regression algorithm is the fact that prices of European options on the underlying swaps, i.e. European swaptions, could be calculated (or approximated in a computationally efficient manner) in most models of interest. As mentioned in Section 18.3.9.3, usage of (proxies of) European swaptions allows us to better incorporate convexity in the hold values in the regressions, improving the final value estimate. Additionally, we can draw on all tricks and enhancements from Section 18.3.10, including, in particular, policy improvement based on the carry argument of Section 18.3.10.2 (we extend the carry results for Bermudan swaptions in Section 19.7.2).

19.6.2 Parametric Boundary Methods

One hallmark of the LS regression approach is its “semi-automatic” nature: once we have identified some potentially meaningful variables and assumed a particular form of the regression basis functions, we let the regression algorithm work its magic to sort out a reasonable exercise strategy. In contrast, to be effective, the boundary optimization technique introduced in Section 3.5.2 requires more careful thought about the functional form of the exercise boundary. As there is considerable intuition to be gained from the results of the boundary optimization technique, let us spend some time on the application of this method to Bermudan swaption pricing. The material in this section draws on results in Andersen [2000a], where the reader can look up many additional details and numerical results that we do not list in our brief treatment here.

19.6.2.1 Sample Exercise Strategies for Bermudan Swaptions

Perhaps the simplest exercise strategy for a Bermudan swaption (and for many other options with early exercise rights) is to “exercise when the option is sufficiently deep in the money”. Mathematically speaking, if $\iota(\cdot)$ is the exercise indicator function¹¹ then our first proposed strategy is

$$\text{Strategy I: } \iota(T_n) = 1_{\{V_{n,N-n}(T_n) \geq h_I(T_n)\}}, \quad (19.25)$$

where $V_{n,N-n}(T_n)$ is the underlying swap value (see Section 19.4.1) and $h_I(\cdot)$ is some unknown deterministic function. Assuming that the Bermudan swaption is of the payer type, the swap value $V_{n,N-n}(T_n)$ in (19.25) can be computed directly from the yield curve as

$$V_{n,N-n}(T_n) = A_{n,N-n}(T_n) (S_{n,N-n}(T_n) - k).$$

It is clear that the function $h_I(\cdot)$ must be strictly non-negative, a constraint that should be checked and enforced in the search for h_I . As discussed in

¹¹Recall from Section 18.3.8.2 that $\iota(T_n) = 1_{\{U_n(T_n) > H_n(T_n)\}}$ defines the rule for exercising at T_n , assuming we have not exercised previously.

Section 3.5.2, the search for the $N - 1$ values $h_I(T_{N-1}), h_I(T_{N-2}), \dots, h_I(T_1)$ can be conducted in backwards fashion from a set of Monte Carlo pre-trials, starting from the known condition $h_I(T_{N-1}) = 0$ and using the fact that a Bermudan swaption with first exercise date T_1 will have the same optimal exercise indicator function at time T_n as will a Bermudan swaption with first exercise date T_n (as long as both are written on swaps with identical coupons and terminal maturity, of course). For each value of n , establishing $h_I(T_n)$ involves a one-dimensional optimization only, to be done either by outright sorting or by a derivatives-free one-dimensional optimizer. We emphasize that all pre-trials should be cached for numerical efficiency; when using Strategy I, for each path it suffices to store at every date T_n , $n = 1, \dots, N - 1$, the intrinsic swap value as well as the numeraire (e.g. the spot numeraire $B(T_n)$ if working in the spot measure), for a total of $2(N - 1)$ double-precision numbers per path.

As the function $h_I(T_n)$ tends to be decreasing roughly linearly as a function of T_n , it often will suffice to assume that $h_I(\cdot)$ is piecewise linear on the interval $[T_1, T_{N-1}]$, with a low-dimensional number b of break-points $t_1 < t_2 < \dots < t_b$, satisfying $t_1 = T_1$ and $t_b = T_{N-1}$. The $b - 1$ values $h_I(t_1), \dots, h_I(t_{b-1})$ can be found by a series of one-dimensional optimizations as described earlier, with the values of $h_I(\cdot)$ at coupon dates T_1, T_2, \dots, T_{N-1} easily computed by linear interpolation. The piecewise linear representation of the exercise rule not only improves numerical efficiency by reducing the number of optimizations to be performed from $N - 1$ to $b - 1$, but also makes the overall algorithm more robust by assigning more explanatory value to each quantity that is optimized over. Indeed, the fewer parameters that have to be estimated by optimization, the less Monte Carlo pre-trials are necessary to get a smooth, noise free estimation of the exercise boundary. Andersen [2000a] demonstrates that even very low values of b (e.g. 2 or 3) will often suffice, a consequence of the well-known fact that prices of options with early exercise rights tend to be quite insensitive to the precise location of the exercise barrier.

We shall show some numerical results for Strategy I in (19.25) shortly, but let us first introduce some more advanced strategies. Recalling that Bermudan swaptions cannot be worth less than the most expensive core European swaption (see the bound (19.2)), it is reasonable to contemplate the application of a policy improvement step in (19.25) to enforce this constraint. Let us therefore define

$$V_{\text{swaption}, M(n)}^{\max}(T_n) = \max_{i=n+1, \dots, M(n)} V_{\text{swaption}, i}(T_n; k)$$

where $M(n)$ is some n -dependent upper bound for the number of European swaptions to include in the max-operation (and k is the strike, see (19.2)). Then, a second exercise strategy is

$$\text{Strategy II: } \iota(T_n) = \begin{cases} 1, & V_{n,N-n}(T_n) > \max(h_{II}(T_n), V_{\text{swaption}, M(n)}^{\max}(T_n)), \\ 0, & \text{otherwise.} \end{cases} \quad (19.26)$$

Setting $M(n) = N - 1$ for all n would ensure that our strategy never breaks the hard value bound (19.2), but could also make computation of the strategy computationally expensive, particularly in models where European swaption pricing requires non-trivial work. As typically only the first few European swaptions are candidates for the maximum in (19.26), to cut down on numerical work¹² it may make sense to write

$$M(n) = \min(N - 1, n + 1 + m), \quad (19.27)$$

where m is some relatively small integer, e.g. 1 or 2.

A strategy related to (19.26) is

$$\text{Strategy III: } \iota(T_n) = \begin{cases} 1, & V_{n,N-n}(T_n) > h_{III}(T_n) + V_{\text{swaption}, M(n)}^{\max}(T_n), \\ 0, & \text{otherwise.} \end{cases} \quad (19.28)$$

Strategy III replaces the absolute trigger condition of Strategy I with a relative one, where exercise takes place when the intrinsic value is sufficiently high relative to the most expensive core European swaption. To some extent a Bermudan swaption can be viewed as a multi-factor best-of option (that is, an option to choose the most expensive of several assets, see Section 19.2), and Strategy III allows one to impose the well-known condition that exercise never takes place when the underlying assets are too close to each other, irrespective of their magnitudes¹³. Notice that if we enforce that h_{III} be strictly non-negative, Strategy III would automatically enforce the policy improvement condition in (19.26). By considering multiple component swaptions, Strategies II and III effectively embed more information about the detailed state of the yield curve into the exercise decision than Strategy I. Strategies II and III can thus be expected to be most useful in a multi-factor model¹⁴. Note that both Strategies II and III can be modified the same way as Strategy I to allow for a piecewise linear representation of the trigger functions h_{II} and h_{III} on some low-dimensional grid $\{t_i\}_{i=1}^b$. Also note that the storage requirements for pre-simulations of Strategies II and III will involve $3(N - 1)$ number per pre-simulation, as we must store on each exercise date T_n i) the numeraire value; ii) the intrinsic swap value; and iii) the maximum core European swaption value.

¹²An alternative, and even less expensive, technique to apply policy improvement would rely on the carry argument developed in Section 19.7.2. The resulting bound requires no option price computations.

¹³For a discussion of exercise strategies for best-of options (also known as *MAX-options*), see Broadie and Detemple [1997].

¹⁴Indeed, we notice that Strategy I is, in fact, optimal for a 1-factor Markov short rate model. For a 1-factor LM model, however, Strategy I is not optimal (although, as we shall see later, it appears to perform very well).

Strategies I–III all involve only sequences of one-dimensional optimizations to uncover the scalar functions h_I , h_{II} , h_{III} ; all optimizations start with the boundary condition $h_I(T_{N-1}) = h_{II}(T_{N-1}) = h_{III}(T_{N-1}) = 0$. Higher-dimensional strategies are possible, too, although they rarely seem worth the extra effort. In Andersen [2000a], the strategies (19.25) and (19.28) are combined into

$$\text{Strategy IV: } \iota(T_n) = \begin{cases} 1, & V_{n,N-n}(T_n) > \max(h_{IV}^1(T_n), h_{IV}^2(T_n) \\ & + V_{\text{swaption}, M(n)}^{\max}(T_n)), \\ 0, & \text{otherwise,} \end{cases}$$

where now two functions, h_{IV}^1 and h_{IV}^2 , have to be determined by optimization. In Andersen [2000a] it is found that this strategy results in no statistically significant pick-up in Bermudan value compared to the simpler strategies above.

19.6.2.2 Some Numerical Tests

To test the exercise strategies outlined above, we shall use simple one- and two-factor log-normal LM models. Specifically, we consider a setting where 3 month Libor rates satisfy

$$dL_k(t)/L_k(t) = O(dt) + \lambda_k(t)^{\top} dW(t),$$

where $W(t)$ is a vector Brownian motion and where the drift term depends on the probability measure, see Chapter 14. We consider two settings of $\lambda_k(t)$:

$$\text{Scenario A : } \lambda_k(t) = 20\%,$$

$$\text{Scenario B : } \lambda_k(t) = \left(15\%, 15\% - \sqrt{0.009(T_k - t)}\right)^{\top}.$$

In Tables 19.5 and 19.6 are numerical results for various Bermudan swaptions, using several of the strategies outlined earlier. We used an initial forward curve that was flat at 10% (quarterly compounded). In computing the lower bounds in the tables, we used 5,000 pre-trials to establish the trigger functions h_I , h_{II} , and h_{III} on a time line with $b = 4$ break-points; in Strategies II and III, we used $M(n) = N - 1$ for all n . 50,000 independent pricing paths were subsequently drawn to compute the lower bounds for each strategy. The tables also include upper bound duality results, computed from Strategy I using the nested simulation algorithm in Section 18.3.8.2, with $K_U = 750$ outer paths and $K_{\text{nest}} = 300$ inner paths. The 95% confidence interval (CI) listed in the tables were computed as outlined in Section 18.3.8.4.

To comment on the tables, we first notice from Table 19.5 that the duality gap computed from Strategy I is never more than 1–2 basis points, leading

Type	Strike	Strategy I	Strategy II	Strategy III	$\hat{\Delta}$	95% CI
15M/3M	8%	184.6 (0.1)	184.6 (0.1)	184.6 (0.1)	0.02	184.5 - 184.8
15M/3M	10%	49.1 (0.1)	49.1 (0.1)	48.9 (0.1)	0.02	48.7 - 49.2
15M/3M	12%	8.9 (0.1)	8.9 (0.1)	8.7 (0.1)	0.004	8.5 - 8.9
3Y/1Y	8%	355.6 (0.4)	355.6 (0.4)	355.1 (0.4)	0.07	354.3 - 355.9
3Y/1Y	10%	157.8 (0.5)	157.8 (0.5)	156.8 (0.5)	0.2	156.0 - 158.0
3Y/1Y	12%	61.8 (0.4)	61.8 (0.4)	61.0 (0.3)	0.04	60.2 - 61.7
6Y/1Y	8%	807.2 (0.9)	807.2 (0.9)	808.0 (0.9)	0.23	805.9 - 809.8
6Y/1Y	10%	417.8 (0.9)	417.8 (0.9)	416.9 (0.9)	0.63	13.7 - 418.0
6Y/1Y	12%	212.7 (0.9)	212.7 (0.9)	212.6 (0.9)	0.33	11.4 - 215.2
11Y/1Y	8%	1381.6 (1.6)	1381.6 (1.6)	1380.2 (1.6)	1.33	1378.5 - 1386.3
11Y/1Y	10%	812.9 (1.4)	812.9 (1.4)	813.2 (1.4)	1.26	810.1 - 817.1
11Y/1Y	12%	495.8 (1.5)	495.8 (1.5)	496.7 (1.4)	0.71	495.3 - 502.1
6Y/3Y	8%	493.2 (0.8)	493.7 (0.8)	493.3 (0.8)	0.08	492.3 - 495.7
6Y/3Y	10%	293.6 (0.9)	294.6 (0.9)	293.0 (0.9)	0.65	292.4 - 296.7
6Y/3Y	12%	170.3 (0.8)	170.3 (0.8)	169.9 (0.8)	0.53	168.9 - 172.8

Table 19.5. Upper and lower bound results for the one-factor model in Scenario A. The initial forward curve is flat at 10%, quarterly compounded. All values are computed using Euler-style discretization and are reported in upfront basis points; numbers in parentheses are sample Monte Carlo errors. “Type” refers to the maturity/lockout period of the Bermudan swaption. “ $\hat{\Delta}$ ” is the upper-lower duality gap estimate and “95% CI” is the 95% confidence interval for the Bermudan swaption price. The computational setup is described in more detail in the text.

Type	Strike	Strategy I	Strategy II	Strategy III	$\hat{\Delta}$	95% CI
15M/3M	8%	184.0 (0.0)	184.0 (0.0)	184.0 (0)	0.05	183.9 - 184.1
15M/3M	10%	43.3 (0.1)	43.4 (0.1)	43.2 (0.1)	0.06	43.1 - 43.6
15M/3M	12%	5.6 (0.1)	5.6 (0.1)	5.6 (0.1)	0.01	5.5 - 5.7
3Y/1Y	8%	339.7 (0.2)	339.8 (0.2)	339.4 (0.2)	0.4	339.2 - 340.6
3Y/1Y	10%	125.8 (0.3)	125.9 (0.3)	125.7 (0.3)	0.7	125.1 - 127.2
3Y/1Y	12%	36.9 (0.2)	36.8 (0.2)	36.6 (0.2)	0.2	36.4 - 37.6
6Y/1Y	8%	750.2 (0.6)	749.6 (0.6)	751.6 (0.6)	3.7	749.0 - 755.2
6Y/1Y	10%	317.0 (0.7)	315.9 (0.7)	319.4 (0.7)	5.0	315.6 - 323.5
6Y/1Y	12%	127.7 (0.6)	128.0 (0.6)	129.2 (0.6)	2.6	126.5 - 131.6
11Y/1Y	8%	1247.3 (1.2)	1250.9 (1.2)	1253.7 (1.3)	18.1	1245.1 - 1269.0
11Y/1Y	10%	620.8 (1.1)	627.1 (1.1)	633.2 (1.3)	20.8	618.4 - 645.0
11Y/1Y	12%	327.1 (1.2)	331.8 (1.1)	337.0 (1.2)	14.8	324.7 - 345.0
6Y/3Y	8%	444.7 (0.6)	444.4 (0.6)	445.2 (0.6)	0.8	443.6 - 446.6
6Y/3Y	10%	226.9 (0.7)	227.2 (0.7)	227.5 (0.7)	1.2	225.5 - 229.5
6Y/3Y	12%	107.1 (0.6)	107.1 (0.6)	107.6 (0.6)	0.8	105.9 - 109.0

Table 19.6. Upper and lower bound results for two-factor model in Scenario B. All values are in upfront basis points; numbers in parentheses are sample Monte Carlo errors. Labels are identical to those of Table 19.5.

.

us to conclude that Strategy 1 very accurately captures the correct exercise decision for the model setup in Table 19.5. Supporting this conclusion is the fact that Strategies II and III lead to no statistically significant increase in the Bermudan swaption value. In the two-factor scenario in Table 19.6, the duality gaps are, not surprisingly, wider than for the one-factor case, although still relatively small for most of the contracts examined. Reasonably significant spreads, in the order of 15 to 20 basis points, can be observed for the 11 year contract with 1 year lockout. Intuitively, for the correlation effects introduced by the two-factor model to matter, the exercise period must be quite long; otherwise, even a two-factor model would imply near-perfect correlation of the different swaps the option holder can exercise into. The suboptimality of exercise based on Strategy I for the 11 no-call 1 Bermudan swaption is also reflected in the fact that the more complicated Strategies II and, especially, III here pick up significant additional value relative to Strategy I. In fact, Strategy III produces prices that lie close to the average of the upper and lower bound, suggesting that this strategy is likely quite close to optimal. Using Strategy III (rather than Strategy I) to form an upper bound confirms this: the duality gap for the 11 year contract with 1 year lockout is reduced to 7.3, 6.3, and 3.5 basis points for coupons of 8%, 10%, and 12%, respectively.

While one should not read too much into the limited set of test data presented above, our results do suggest that for models without stochastic volatility, Strategy I is sufficient for short-dated Bermudan swaptions and for models with high forward rate correlation. For longer-dated structures and for multi-factor models, Strategy III is a safer bet. In a LS regression setting, this reinforces our observations of Section 19.6.1 on the importance of including variables that represent both the level and the slope of the yield curve on exercise dates.

19.6.2.3 Additional Comments

A number of papers in the literature elaborate on the analysis in Andersen [2000a]. For instance, in an LM model setting Jensen and Svenstrup [2003] conclude that for Strategy III just setting $m = 1$ in (19.27) typically yields Bermudan swaption values that are indistinguishable from those computed using $M(n) = N - 1$. Jensen and Svenstrup [2003] also compare the parametric boundary optimization technique against an LS regression where the basis functions include the first two powers of the intrinsic swap value and the spot numeraire, as well as their cross product. For an LM model without stochastic volatility, the parametric boundary technique with Strategy III is found to slightly outperform this particular setup of the LS regression. A similar conclusion is reached in Pedersen [1999], where more details on the LS regression for Bermudan swaptions can also be found.

The analysis in Andersen [2000a] (and our discussion in the previous section) concerns itself only with models that contain no stochastic volatility

component. Jensen and Svenstrup [2003] examine an LM model with stochastic volatility and conclude that in this case a rather small, but economically significant, duality gap opens up for Strategy III, especially when the volatility of variance is large and/or the mean reversion speed of volatility is low. Not surprisingly, an LS regression where the variance level itself is included in the set of regressors manages to lower this duality gap. In general, for models with stochastic volatility, explicitly specifying the functional form of the exercise boundary seems difficult, and the best approach is typically to use a regression approach to uncover it, as we described in Section 18.3.9.3.

19.7 Other Topics

19.7.1 Robust Bermudan Swaption Hedging with European Swaptions

As we explain in more detail in Chapter 22, risk management of exotic derivatives such as Bermudan swaptions generally involves both delta hedging (offsetting sensitivity to the yield curve by dynamically trading in swaps) and vega hedging (offsetting sensitivity of Bermudan swaptions to changes in volatility). Vega hedging of Bermudan swaptions and other CLEs is typically done by trading in European swaptions, but as transaction costs for options can be relatively high, dealers would prefer hedges that do not require frequent rebalancing. A good example of such a hedge would be the static hedging of CMS-linked derivatives with European swaptions at multiple strikes, as specified by the replication method in Section 16.6.1. The resulting hedge not only needs no rebalancing, but is model-independent (up to the annuity mapping function selection).

For Bermudan swaptions we are not aware of any known model-independent static hedge position in European swaptions, but some insights can be gained from the marginal exercise value decomposition of Section 18.2.3. Even though each of the European options in the decomposition, $E(B(T_n)^{-1}(U_n(T_n) - H_n(T_n))^+)$, is not a standard European swaption, the replication method of Section 16.6.1 tells us that it can easily be represented as a static position of European swaptions over a continuum of strikes. This position, however, is *not* model-independent, as each payoff $(U_n(T_n) - H_n(T_n))^+$ is sensitive (through $H_n(T_n)$) to the model volatilities and, in particular, the *forward volatilities* produced by the model. Also, as volatilities inevitably change over time, the effective payoffs are liable to change, as is therefore the composition of the European swaption portfolio that the Bermudan swaption is decomposed into. This implies that the hedge portfolio will need rebalancing over time, which of course rules it out as a truly static hedging portfolio.

Another approach that could be pursued is the semi-static decomposition of barrier options into European swaptions developed by Andersen et al.

[2002]. While developed specifically for barrier options, the technique of this paper also applies to Bermudan swaptions, as these can be interpreted as barrier options with a knock-in barrier set to the optimal exercise boundary (in fact, we already used this representation in developing valuation algorithms). Unfortunately, this line of attack also fails to produce a static hedge, for the same reasons as for the marginal exercise value decomposition: the hedge portfolio depends strongly on the model-specific volatility structure, and is also likely to need rebalancing over time as volatility moves around randomly.

While no theoretically airtight static hedge for Bermudan swaptions is known (as far as we are aware), various pragmatic strategies — often collectively known as the *portfolio replication approach* — have been more successful. Let us briefly summarize the main idea here, while acknowledging the fact that it can be implemented in many different ways. As an aside, we note that the approach can be applied to many exotic derivatives, although its performance would ultimately depend on the specific risk characteristics of the specific derivative under consideration.

We start by identifying a universe of potential hedging instruments. For a Bermudan swaption of given maturity T_N , we typically would select all European swaptions with expiry + tenor less than or equal T_N , with strikes chosen to span a reasonably wide range. Having identified the hedging instruments, we formulate market data scenarios that we would want the hedge portfolio to cover. These would typically be scenarios of joint moves of the yield curve and the volatility surface. For Bermudan swaptions, it is probably sufficient to choose parallel moves in the yield curve, moves by a pre-specified amount within a given range, although one can add more complicated ones, e.g. the yield curve twists and “bends” suggested by principal components analysis (see Section 14.3.1). Similar types of volatility scenarios could be used, such as parallel and non-parallel shifts across all swaption expiries and maturities.

Suppose we have defined M scenarios and chosen K hedging instruments. Let ΔV_{Berm} denote the M -dimensional vector of value changes of the Bermudan swaption in all M scenarios, and let the vector of value changes of the k -th hedging instrument be ΔV_k , $k = 1, \dots, K$. Then, on the last step of the portfolio replication method, we look for a vector of weights $\chi = (\chi_1, \dots, \chi_K)^\top$ such that the portfolio of hedges defined by these weights immunizes the changes in values of the Bermudan swaption in all scenarios. This is usually formalized as a least-squares optimization problem,

$$\left\| \Delta V_{\text{Berm}} - \sum_{k=1}^K \chi_k \Delta V_k \right\|^2 \rightarrow \min.$$

Variations are possible, including weighting different scenarios differently or adding additional terms to the objective function to express user preferences,

such as minimizing the total notional of all swaptions, penalizing excessive use of deep out-of-the-money swaptions, and so forth.

We generally find the portfolio replication method to be an effective risk management tool for Bermudan swaptions. Anecdotal evidence provided by traders suggests that the method often outperforms the standard delta/vega hedging approach relying on “local” sensitivities. It also provides a relatively straightforward way of dealing with the well-known gamma-theta mismatch that plagued many dealers’ Bermudan swaption portfolios in the 1980s and 1990s. As it turns out, if one uses the volatility sensitivity information from certain simple models to hedge the vega (volatility sensitivity) of a Bermudan swaption by selling European swaptions, the resulting position will sometimes be short gamma (second order yield curve sensitivity) *and* short time decay (theta). This, however, runs counter to the “standard” Black-Scholes theory in which the gamma and theta balance each other: a long gamma position is always short time decay and makes money in volatile markets and loses money in calm markets; and a short-gamma, long-theta position does the opposite. The unenviable position of being short gamma and short theta will tend to lose money in *all* markets, volatile and calm, and historically resulted in a number of Bermudan swaption book disasters over the years. The portfolio replication method can help resolve the gamma-theta problem by effectively hedging the exposure to forward volatility or to inter-temporal correlation with European swaptions across multiple expiries and tenors¹⁵ — something a globally calibrated LM model would do, for example. On the other hand, vega hedging positions computed in short rate models calibrated according to the views held at the time (where mean reversion was often considered superfluous and either excluded from consideration or linked, directly or indirectly, to volatility as in, e.g., the BDT model from Section 11.1.1) would generally suggest incorrect European swaption hedges for the unobservable volatility positions. For example if one does not explicitly link mean reversion to market values of (off-diagonal) swaptions, the forward volatility exposure could either be not hedged at all, or might (wrongly) be linked to the diagonal European swaptions.

19.7.2 Carry and Exercise

In Section 18.3.10.2 we showed that if on a given exercise date the next (net) coupon of a cancelable note is positive, then it is not optimal to cancel the note on that date. The net coupon of a derivative security is often referred to as its *carry*, so we can state that carry-positive cancelable notes should never be exercised.

¹⁵Of course, when applying the portfolio replication method, we should explicitly link mean reversion to market inputs through the local projection method of Sections 19.2 and 18.4.

Cancelable notes and callable Libor exotics are, of course, intricately linked (see for example Section 18.3.3), so it should come as no surprise that carry-based restrictions on exercise decisions exist for Bermudan swaptions. In fact a result more general than (18.56) holds for cancelable notes and CLEs, but we present it in this chapter, since only for Bermudan swaptions is this more general result actually easy to apply.

We start by recalling that the n -th exercise value of a Bermudan payer swaption can be written as

$$U_n(t) = A_n(t)(S_n(t) - k).$$

A simple relation follows,

$$\begin{aligned} U_n(t) &= [P(t, T_n) - P(t, T_{n+1}) - k\tau_n P(t, T_{n+1})] + U_{n+1}(t) \\ &= \tau_n P(t, T_{n+1}) \left[\frac{P(t, T_n) - P(t, T_{n+1})}{\tau_n P(t, T_{n+1})} - k \right] + U_{n+1}(t) \\ &= A_{n,1}(t)[S_{n,1}(t) - k] + U_{n+1}(t), \end{aligned}$$

and, more generally,

$$U_n(t) = A_{n,m}(t)[S_{n,m}(t) - k] + U_{n+m}(t), \quad m \geq 1, \quad (19.29)$$

where we used the notation $A_{n,m}(t)$ and $S_{n,m}(t)$ for the annuity and the swap rate, respectively, for a swap starting at T_n and covering m periods. Clearly, for the hold values we have ($m \geq 1$),

$$H_n(t) \geq H_{n+m-1}(t) \geq U_{n+m}(t),$$

hence it follows from (19.29) that

$$U_n(t) \leq A_{n,m}(t)[S_{n,m}(t) - k] + H_n(t), \quad m \geq 1.$$

Taking a minimum over all m we obtain the following result.

Proposition 19.7.1. *For a given Bermudan payer swaption and a given exercise date T_n , we have*

$$U_n(T_n) - H_n(T_n) \leq \min_{m \geq 1} \{A_{n,m}(T_n)[S_{n,m}(T_n) - k]\}, \quad (19.30)$$

and so if any of the swaps that start at T_n and have maturities up to the final maturity of the Bermudan swaption have negative value at T_n , it is never optimal to exercise at time T_n .

Proof. If there exists m such that $S_{n,m}(T_n) - k < 0$ then by (19.30), the exercise value is strictly less than the hold value, and the exercise is not optimal. \square

As annuity factors are always positive, it follows from Proposition 19.7.1 above that a Bermudan payer swaption should never be exercised if any swap rate of any still-alive swap is less than the fixed coupon k . A similar result holds for Bermudan receiver swaptions — we trust the reader can derive it himself.

19.7.3 Fast Pricing via Exercise Premia Representation

There are situations when the speed of valuation of Bermudan swaptions is key yet the accuracy could be sacrificed. One example is robust hedging of Bermudan swaptions in Section 19.7.1 where high-precision pricing is not particularly important. Another is the calculation of a *credit value adjustment* (CVA), an adjustment to the value of a derivative that takes into account the possibility that the counterparty could default on its payments. While CVA calculations are outside the scope of this book (see Gregory [2009] for a good overview), in essence the evaluation of CVA requires prices of a portfolio of Bermudan swaptions at many future dates under many simulated market conditions. Again, speed of valuation here is very important. One can speed up valuation by using a simple model, such as the one-factor Gaussian model of Section 10, but even higher performance is often desired. In this section we consider a useful approximation based on the representation of a Bermudan swaption as a stream of coupons paid in the exercise region, an adaptation of the representation we developed in Section 1.10.3 for American options.

Recall the definition of $\iota(\cdot)$, the exercise indicator, from Section 19.6.2.1 and Section 18.3.8.2. Let us denote by $V(t)$ the value of the Bermudan swaption at time t as in (18.32).

Proposition 19.7.2. *The following holds for any $n = 1, \dots, N - 1$,*

$$\begin{aligned} \mathbb{E}_{T_n} \left(\frac{V(T_n)}{B(T_n)} - \frac{V(T_{n+1})}{B(T_{n+1})} \right) &= \mathbb{E}_{T_n} (\iota(T_n) B(T_{n+1})^{-1} X_n) \\ &+ \mathbb{E}_{T_n} (\iota(T_n) (1 - \iota(T_{n+1})) B(T_{n+1})^{-1} (U_{n+1}(T_{n+1}) - H_{n+1}(T_{n+1}))), \end{aligned} \quad (19.31)$$

where we have used the convention that $\iota(T_N) \equiv 1$ and $V(T_N) \equiv 0$. In particular,

$$\begin{aligned} V(0) &= \mathbb{E} \left(\sum_{n=1}^{N-1} \iota(T_n) B(T_{n+1})^{-1} X_n \right) \\ &+ \mathbb{E} \left(\sum_{n=1}^{N-1} \iota(T_n) (1 - \iota(T_{n+1})) B(T_{n+1})^{-1} (U_{n+1}(T_{n+1}) - H_{n+1}(T_{n+1})) \right). \end{aligned} \quad (19.32)$$

Proof. We have that

$$V(T_n) = \iota(T_n) U_n(T_n) + (1 - \iota(T_n)) H_n(T_n)$$

and

$$\begin{aligned} & \frac{V(T_n)}{B(T_n)} - \frac{V(T_{n+1})}{B(T_{n+1})} \\ &= \iota(T_n) \left(\frac{U_n(T_n)}{B(T_n)} - \frac{V(T_{n+1})}{B(T_{n+1})} \right) + (1 - \iota(T_n)) \left(\frac{H_n(T_n)}{B(T_n)} - \frac{V(T_{n+1})}{B(T_{n+1})} \right). \end{aligned}$$

Taking expected value conditioned on \mathcal{F}_{T_n} and using the fact that $\iota(T_n)$ is \mathcal{F}_{T_n} -measurable and

$$H_n(T_n) = B(T_n) \mathbb{E}_{T_n} \left(\frac{V(T_{n+1})}{B(T_{n+1})} \right)$$

by definition, we get

$$\mathbb{E}_{T_n} \left(\frac{V(T_n)}{B(T_n)} - \frac{V(T_{n+1})}{B(T_{n+1})} \right) = \iota(T_n) \mathbb{E}_{T_n} \left(\frac{U_n(T_n)}{B(T_n)} - \frac{V(T_{n+1})}{B(T_{n+1})} \right).$$

Moreover,

$$\begin{aligned} & \iota(T_n) \left(\frac{U_n(T_n)}{B(T_n)} - \frac{V(T_{n+1})}{B(T_{n+1})} \right) \\ &= \iota(T_n) \left(\frac{U_n(T_n)}{B(T_n)} \right. \\ &\quad \left. - \frac{1}{B(T_{n+1})} (\iota(T_{n+1}) U_{n+1}(T_{n+1}) + (1 - \iota(T_{n+1})) H_{n+1}(T_{n+1})) \right) \\ &= \iota(T_n) \left(\frac{U_n(T_n)}{B(T_n)} - \frac{U_{n+1}(T_{n+1})}{B(T_{n+1})} \right) \\ &\quad + \iota(T_n) (1 - \iota(T_{n+1})) \frac{1}{B(T_{n+1})} (U_{n+1}(T_{n+1}) - H_{n+1}(T_{n+1})). \end{aligned}$$

Then, since

$$\mathbb{E}_{T_n} \left(\frac{U_n(T_n)}{B(T_n)} - \frac{U_{n+1}(T_{n+1})}{B(T_{n+1})} \right) = \mathbb{E}_{T_n} \left(\frac{X_n}{B(T_{n+1})} \right)$$

(see (18.2)), the result (19.31) follows. Finally, the result (19.32) follows by summing up equalities (19.31) for $n = 1, \dots, N - 1$ and taking the (unconditional) expected value. \square

Let us consider the second term on the right-hand side of (19.31). It represents the contribution of those paths that are in the exercise region at time T_n ($\iota(T_n) = 1$) and are in the hold region at time T_{n+1} ($\iota(T_{n+1}) = 0$). One can argue that there are “not too many” of such paths, especially if T_n and T_{n+1} are relatively close. Moreover, quantities that are actually evaluated for those paths, the differences between exercise and hold values $U_{n+1}(T_{n+1}) - H_{n+1}(T_{n+1})$, will be small because the exercise value is close to the hold value on the border between exercise and hold regions (by definition of the exercise boundary). Indeed, in the continuous-exercise limit these

terms simply disappear, as should be clear from comparing (19.32) to (1.77). These considerations lead us to suggest an approximation to the value of a Bermudan swaption in which we simply disregard the second sum on the right-hand side of (19.32):

Corollary 19.7.3. *The value of a Bermudan swaption (or, indeed, any callable Libor exotic) is approximately equal to the sum of (net) coupons that are paid only in the exercise region, i.e.*

$$\begin{aligned} V(0) &\approx E \left(\sum_{n=1}^{N-1} \iota(T_n) B(T_{n+1})^{-1} X_n \right) \\ &= E \left(\sum_{n=1}^{N-1} 1_{\{U_n(T_n) \geq H_n(T_n)\}} B(T_{n+1})^{-1} X_n \right). \end{aligned} \quad (19.33)$$

The error of approximation is given by the second term in (19.32); the error will decrease as the frequency of exercise of the Bermudan swaption is lowered.

At this point the reader may recall that we have already derived a similar-looking representation for CLEs, namely the marginal exercise value decomposition of Proposition 18.2.1, where we can rewrite (18.8) as

$$V(0) = E \left(\sum_{n=1}^{N-1} 1_{\{U_n(T_n) \geq H_n(T_n)\}} B(T_n)^{-1} (U_n(T_n) - H_n(T_n)) \right). \quad (19.34)$$

Not surprisingly, (19.33) could also be derived from (19.34) if we observe that, when the Bermudan is “deep in the money” at time T_n then

$$\begin{aligned} H_n(T_n) &= B(T_n) E_{T_n} (B(T_{n+1})^{-1} \max(U_{n+1}(T_{n+1}), H_{n+1}(T_{n+1}))) \\ &\approx B(T_n) E_{T_n} (B(T_{n+1})^{-1} U_{n+1}(T_{n+1})) \\ &= U_{n+1}(T_n) = U_n(T_n) - P(T_n, T_{n+1}) X_n, \end{aligned}$$

and, therefore, we see that the marginal exercise value decomposition implies that

$$\begin{aligned} V(0) &\approx E \left(\sum_{n=1}^{N-1} \iota(T_n) B(T_n)^{-1} P(T_n, T_{n+1}) X_n \right) \\ &= E \left(\sum_{n=1}^{N-1} \iota(T_n) B(T_{n+1})^{-1} X_n \right) \end{aligned}$$

which is (19.33) in Corollary 19.7.3.

Everything we have discussed so far is valid for general CLEs. Let us now specialize our setup to Bermudan swaptions, with the net coupon X_n

given by $\tau_n(L_n(T_n) - k)$. The coupon is a function of the Libor rate $L_n(T_n)$; critically, in pretty much all one-factor models (and certainly in the one-factor Gaussian model) the exercise boundary at time T_n could also be parameterized by the same Libor rate and expressed in the form

$$\iota(T_n) = 1_{\{L_n(T_n) \geq h(T_n)\}}$$

for some deterministic function $h(\cdot)$. Then we can rewrite (19.33) as

$$V(0) \approx \sum_{n=1}^{N-1} P(0, T_{n+1}) \tau_n E^{T_{n+1}} ((L_n(T_n) - k) 1_{\{L_n(T_n) \geq h(T_n)\}}). \quad (19.35)$$

Each term on the right-hand side can be expressed as a combination of caplets and digital caplets on the Libor rate L_n . Notice that the expected value is taken under the T_{n+1} -forward measure, a measure under which the Libor rate is a martingale. In the one-factor Gaussian model the distribution of $L_n(T_n)$ under the T_{n+1} -forward measure is well approximated by the Gaussian distribution, and each term could be evaluated rapidly with just a few applications of the Bachelier formula (7.16). Similar approximations could be derived in many other models. For the Libor market model, in particular, the distributions of Libor rates is often known exactly.

The exercise boundary function $h(\cdot)$ in (19.35) is not known *a-priori* and needs to be found as part of valuation. This can be done efficiently in a backward induction algorithm that utilizes the representation (19.35) for a Bermudan swaption at future times. In particular, we can find the value $h(T_k)$ by solving

$$H_k(T_k)|_{L_k(T_k)=h(T_k)} = U_k(T_k)|_{L_k(T_k)=h(T_k)},$$

where $H_k(T_k)$ is calculated by an analog to (19.35) with $h(T_{k+1}), \dots, h(T_{N-1})$ already determined from previous steps. The recursion can be accelerated further by search for the exercise boundary only on the subset of exercise dates, with the missing points filled by interpolation (see Ju [1998]). The final algorithm turns out to be quick and robust, and is well-suited for situations where valuation speed is the primary consideration.

19.A Appendix: Forward Volatility and Correlation

When European swaption prices are kept fixed, increasing correlation between forward rates (e.g. by moving from a two- to a one-factor model) will tend to increase forward volatility. There are a number of ways to explain this effect. One approach relies on swaption “triangles” (see Section 20.3), another on the BGM/HJM model formalism. In this appendix we explore the latter.

First, a bit of notation. Let Libor rates $L_k(t) = L_k(t, T_k, T_{k+1})$ satisfy

$$dL_k(t) = O(dt) + \sigma_k(t) dW_k(t), \quad k = 1, \dots, N - 1,$$

where σ_k 's are scalar and the W_k 's are correlated Brownian motions. We assume that our calibration is global (see Section 14.5) and therefore the model calibrates properly to caplets (or, equivalently, swaptions on short-tenor swaps). It follows that the quantities

$$\int_0^{T_k} \sigma_k(u)^2 du \tag{19.36}$$

must be invariants, independent of the correlation between the W_k 's.

Consider now the market for short-expiry swaptions — which must also be calibrated in our setup (which is global) — and let us study two different settings of the average correlation between the W_k 's, “high” and “low”, indicated by appropriate subscripts in what follows. To match the short-expiry, short-maturity swaptions (i.e., caplets), we fundamentally need

$$\sigma_1^{\text{hi}}(0) = \sigma_1^{\text{lo}}(0).$$

On the other hand, to match long-tenor (but still short-expiry) swaptions, we need

$$\sigma_k^{\text{hi}}(0) < \sigma_k^{\text{lo}}(0), \quad k = 2, \dots, N - 1, \tag{19.37}$$

a relationship that follows from the fact that swaption volatilities are effectively volatilities of sums of Libor rates. Specifically, as the volatility of a sum increases in correlation, we need to lower the volatility of the “components of the sum” (that is, the Libor rates) to preserve swaption volatilities.

Let us pick some $k > 1$. If we look at (19.36) and (19.37), it is obvious that to satisfy both conditions simultaneously $\sigma_k^{\text{hi}}(t)$ must¹⁶ ultimately “overtake” $\sigma_k^{\text{lo}}(t)$ as t is increased from 0 to T_k . As this holds for all k , it is clear that forward volatilities of both caps and swaptions will, as promised above, be higher in the high-correlation model than in the low-correlation model.

19.B Appendix: A Primer on Moment Matching

19.B.1 Basics

Let there be given d log-normal random variables X_1, \dots, X_d , with known distribution parameters m_i and s_i :

$$\ln(X_i) \sim \mathcal{N}(m_i, s_i^2), \quad i = 1, \dots, d.$$

¹⁶See Figures 1–3 in Andersen and Andreasen [2001] for visual confirmation of this as well as of (19.37).

We assume that the $d \times d$ correlation matrix ρ of logarithms is known,

$$\rho_{i,j} \triangleq \text{Corr}(\ln(X_i), \ln(X_j)), \quad i, j = 1, \dots, d.$$

From standard results for log-normal variables, the first two moments of the X_i can be computed as

$$\mathbb{E}(X_i) = \exp\left(m_i + \frac{s_i^2}{2}\right), \quad (19.38)$$

$$\text{Var}(X_i) = \mathbb{E}(X_i)^2 (\exp(s_i^2) - 1). \quad (19.39)$$

Also,

$$\begin{aligned} \mathbb{E}(X_i X_j) &= \mathbb{E}(\exp(\ln(X_i) + \ln(X_j))) \\ &= \mathbb{E}\left(\exp\left(m_i + s_i Z + m_j + s_j \left(\rho_{i,j} Z + \sqrt{1 - \rho_{i,j}^2} Y\right)\right)\right) \end{aligned}$$

where Z and Y are independent standard Gaussian variables, and where we have used the Cholesky decomposition. Therefore, using the result (19.38),

$$\mathbb{E}(X_i X_j) = \mathbb{E}(X_i) \mathbb{E}(X_j) \exp(\rho_{i,j} s_i s_j). \quad (19.40)$$

We note in passing that therefore (see e.g. (17.66))

$$\begin{aligned} \text{Cov}(X_i, X_j) &= \mathbb{E}(X_i X_j) - \mathbb{E}(X_i) \mathbb{E}(X_j) \\ &= \mathbb{E}(X_i) \mathbb{E}(X_j) (\exp(\rho_{i,j} s_i s_j) - 1). \end{aligned} \quad (19.41)$$

Suppose now that we are interested in approximating the moments of the weighted sum

$$\widehat{X} = \sum_{i=1}^d w_i X_i, \quad (19.42)$$

where the w_i 's are given positive constants. Clearly

$$\mathbb{E}(\widehat{X}) = \sum_{i=1}^d w_i \mathbb{E}(X_i) \quad (19.43)$$

with $\mathbb{E}(X_i)$ given in (19.38). Also,

$$\begin{aligned} \mathbb{E}(\widehat{X}^2) &= \sum_{i=1}^d \sum_{j=1}^d w_i w_j \mathbb{E}(X_i X_j) \\ &= \sum_{i=1}^d \sum_{j=1}^d w_i w_j \mathbb{E}(X_i) \mathbb{E}(X_j) \exp(\rho_{i,j} s_i s_j) \\ &= \sum_{i=1}^d w_i^2 \mathbb{E}(X_i)^2 e^{s_i^2} + 2 \sum_{i=1}^d \sum_{j=i+1}^d w_i w_j \mathbb{E}(X_i) \mathbb{E}(X_j) e^{\rho_{i,j} s_i s_j}. \end{aligned} \quad (19.44)$$

In many applications, we are interested in representing \widehat{X} as being approximately log-normal, i.e. we would like to write

$$\ln(\widehat{X}) \sim \mathcal{N}\left(m_{\widehat{X}}, s_{\widehat{X}}^2\right).$$

Using a moment-matching principle, we would determine $m_{\widehat{X}}$ and $s_{\widehat{X}}$ from the equations

$$\begin{aligned} \exp\left(m_{\widehat{X}} + \frac{s_{\widehat{X}}^2}{2}\right) &= E(\widehat{X}), \\ \exp\left(m_{\widehat{X}} + \frac{s_{\widehat{X}}^2}{2}\right)^2 \exp\left(s_{\widehat{X}}^2\right) &= E(\widehat{X}^2), \end{aligned}$$

which can be solved to yield

$$s_{\widehat{X}} = \sqrt{\ln(E(\widehat{X}^2)) - \ln(E(\widehat{X})^2)}, \quad m_{\widehat{X}} = \ln(E(\widehat{X})) - \frac{s_{\widehat{X}}^2}{2}. \quad (19.45)$$

In these formulas $E(\widehat{X})$ and $E(\widehat{X}^2)$ should be computed from formulas (19.43) and (19.44), respectively.

Note that if we are willing to relax the requirement that \widehat{X} be approximately log-normal, we can obtain more accurate approximations. A popular choice here is to assume that \widehat{X} is approximately *displaced* log-normal. This introduces one more degree of freedom in the matching distribution (the displacement parameter) which, together with the mean and variance, could be used to match three, rather than two, moments. We leave the details of this for the reader to work out, and in the examples below we stick with simple log-normal moment matching.

19.B.2 Example 1: Asian Option in BSM Model

Let $I(t)$ be some asset following the simple process

$$dI(t)/I(t) = -b(t)dt + \sigma(t)dW(t), \quad (19.46)$$

where $W(t)$ is a scalar Brownian motion in the risk-neutral measure. For certain weights w_i , we form the weighted average

$$M(T_d) = \sum_{i=1}^d w_i I(T_i),$$

on some schedule $0 < T_1 < T_2 < \dots < T_d$. An Asian option pays

$$V_{\text{Asian}}(T_{\text{pay}}) = (M(T_d) - K)^+, \quad T_{\text{pay}} \geq T_d,$$

where typically the weights are $w_i = 1/d$ for all i . Standing at time 0, we wish to use moment-matching to model the T_d -observed average $M(T_d)$ as a log-normal variable. From (19.46), it is clear that $I(T_i)$ is a log-normal random variable, since

$$I(T_i) = I(0)l(T_i) \exp\left(-\frac{1}{2}v(T_i)^2T_i + \int_0^{T_i} \sigma(u) dW(u)\right),$$

where we have defined

$$l(T_i) \triangleq \exp\left(-\int_0^{T_i} b(u) du\right), \quad v(T_i)^2 \triangleq T_i^{-1} \int_0^{T_i} \sigma(u)^2 du.$$

If we define $X_i = I(T_i)$, it follows that, in the notation of Section 19.B.1,

$$\begin{aligned} m_i &= \ln(l(T_i)) - \frac{1}{2}v(T_i)^2T_i, \\ s_i &= v(T_i)\sqrt{T_i}, \end{aligned}$$

and that $M(T_d) = \widehat{X}$, where \widehat{X} is defined in (19.42). To use the results (19.45) it only remains to find the correlation matrix ρ . But clearly

$$\begin{aligned} \rho_{i,j} &= \text{Corr}\left(\int_0^{T_i} \sigma(u) dW(u), \int_0^{T_j} \sigma(u) dW(u)\right) \\ &= \frac{\int_0^{\min(T_i, T_j)} \sigma(u)^2 du}{\sqrt{\int_0^{T_i} \sigma(u)^2 du} \sqrt{\int_0^{T_j} \sigma(u)^2 du}} \\ &= \frac{\min(v(T_i)^2 T_i, v(T_j)^2 T_j)}{v(T_i)v(T_j)\sqrt{T_i T_j}} \\ &= \frac{\min(v(T_i)\sqrt{T_i}, v(T_j)\sqrt{T_j})}{\max(v(T_i)\sqrt{T_i}, v(T_j)\sqrt{T_j})}. \end{aligned}$$

Applying (19.43), (19.44), and finally (19.45) then allows us to write, approximately,

$$\ln(M(T_d)) \sim \mathcal{N}(m_{\widehat{X}}, s_{\widehat{X}}^2)$$

for computed constants $m_{\widehat{X}}$ and $s_{\widehat{X}}$. Assuming deterministic interest rates, standard Black-Scholes arguments (see Section 1.9) allow us to finally approximate the time 0 option price as

$$V_{\text{Asian}}(0) \approx P(0, T_{\text{pay}}) \left(e^{m_{\widehat{X}} + \frac{1}{2}s_{\widehat{X}}^2} \Phi(d_+) - K \Phi(d_-) \right), \quad (19.47)$$

$$d_{\pm} = \frac{\ln\left(e^{m_{\widehat{X}} + \frac{1}{2}s_{\widehat{X}}^2}/K\right) \pm \frac{1}{2}s_{\widehat{X}}^2}{s_{\widehat{X}}} = \frac{m_{\widehat{X}} + \frac{1}{2}s_{\widehat{X}}^2 - \ln(K) \pm \frac{1}{2}s_{\widehat{X}}^2}{s_{\widehat{X}}},$$

where Φ is the Gaussian CDF and $P(0, T_{pay})$ is a risk-free discount factor to time T_{pay} . We note that we may, of course, rewrite this expression in the perhaps slightly more convenient form

$$V_{\text{Asian}}(0) \approx P(0, T_{pay}) (\mathbb{E}(M(T_d)) \Phi(d_+) - K \Phi(d_-)), \quad (19.48)$$

$$d_{\pm} = \frac{\ln(\mathbb{E}(M(T_d))/K) \pm \frac{1}{2}s_{\hat{X}}^2}{s_{\hat{X}}},$$

where

$$\mathbb{E}(M(T_d)) = \sum_{i=1}^d w_i \mathbb{E}(I(T_i)) = \sum_{i=1}^d w_i I(0) l(T_i).$$

Note that this form does not require us to compute $m_{\hat{X}}$, as only $s_{\hat{X}}$ is needed.

19.B.3 Example 2: Basket Option in BSM Model

Consider d risk-neutral processes

$$dI_i(t)/I_i(t) = -b_i(t) dt + \sigma_i(t) dW_i(t), \quad i = 1, \dots, d,$$

where we assume that $\langle dW_i(t), dW_j(t) \rangle = \rho_{i,j} dt$. Also consider the payout of a basket option

$$V_{\text{basket}}(T_{pay}) = (\widehat{I}(T) - K)^+, \quad T_{pay} \geq T,$$

where

$$\widehat{I}(T) = \sum_{i=1}^d w_i I_i(T),$$

with the understanding that all basket weights w_i are positive. In the framework of Section 19.B.1, we now set $X_i = I_i(T)$, such that, at time 0, X_i is log-normal with parameters

$$m_i = \ln(l_i(T)) - \frac{1}{2}v_i(T)^2 T, \quad s_i = v_i T \sqrt{T},$$

where we have defined, for $T > 0$,

$$l_i(T) \triangleq \exp\left(-\int_0^T b_i(u) du\right), \quad v_i(T)^2 \triangleq T^{-1} \int_0^T \sigma_i(u)^2 du.$$

In the notation of Section 19.B.1, clearly $\widehat{I}(T) = \widehat{X}$ and we may proceed as in Example 1 above to find $m_{\widehat{X}}$ and $s_{\widehat{X}}$, at which point the formula (19.47) (or (19.48)) will price the basket option at time 0.

TARNs, Volatility Swaps, and Other Derivatives

Having completed our discussion of callable Libor exotics, in this chapter we turn our attention to a few remaining types of exotic interest rate derivatives that are popular in the market. Our analysis gives us the opportunity to provide additional examples of the local projection method introduced in Chapter 18 which, along with the out-of-model adjustment methods in Chapter 21, are the cornerstone techniques for the situations where computational efficiency constraints prohibit the usage of large, globally calibrated models.

20.1 TARNs

20.1.1 Definitions and Examples

As explained in Section 5.15.2, a TARN (Targeted Redemption Note) pays structured coupons in exchange for Libor coupons until the cumulative amount of structured coupon payments exceeds a pre-agreed target, at which point the derivative terminates. While many coupon types could be used in a TARN, we focus our discussion on inverse floating coupons indexed to the Libor rate. Recall (Section 5.13.1) that an inverse floating coupon with strike s , gearing g , a zero floor and no cap is defined as

$$C_n = (s - g \times L_n(T_n))^+, \quad (20.1)$$

with the underlying rate observed (fixed) at time T_n and the coupon paid at T_{n+1} . We shall use the specific structured coupon (20.1) as an example throughout this section; in defining it, we have used the usual notation for spanning Libor rates

$$L_n(t) = L(t, T_n, T_{n+1}) = \frac{P(t, T_n) - P(t, T_{n+1})}{\tau_n P(t, T_{n+1})}, \quad n = 0, \dots, N-1,$$

and have also introduced a tenor structure

$$0 = T_0 < T_1 < \dots < T_N, \quad \tau_n = T_{n+1} - T_n.$$

In the TARN, the structured coupon fixed at time T_n is only paid if the sum of coupons fixing before (but not including) time T_n is below a given total return R . Thus, from the investor viewpoint, the value of the TARN at time 0 under is given by

$$V_{\text{tarn}}(0) = \mathbb{E} \left(\sum_{n=1}^{N-1} B(T_{n+1})^{-1} \tau_n (C_n - L_n(T_n)) \mathbf{1}_{\{Q_n < R\}} \right), \quad (20.2)$$

$$Q_n = \sum_{i=1}^{n-1} \tau_i C_i, \quad Q_1 = 0,$$

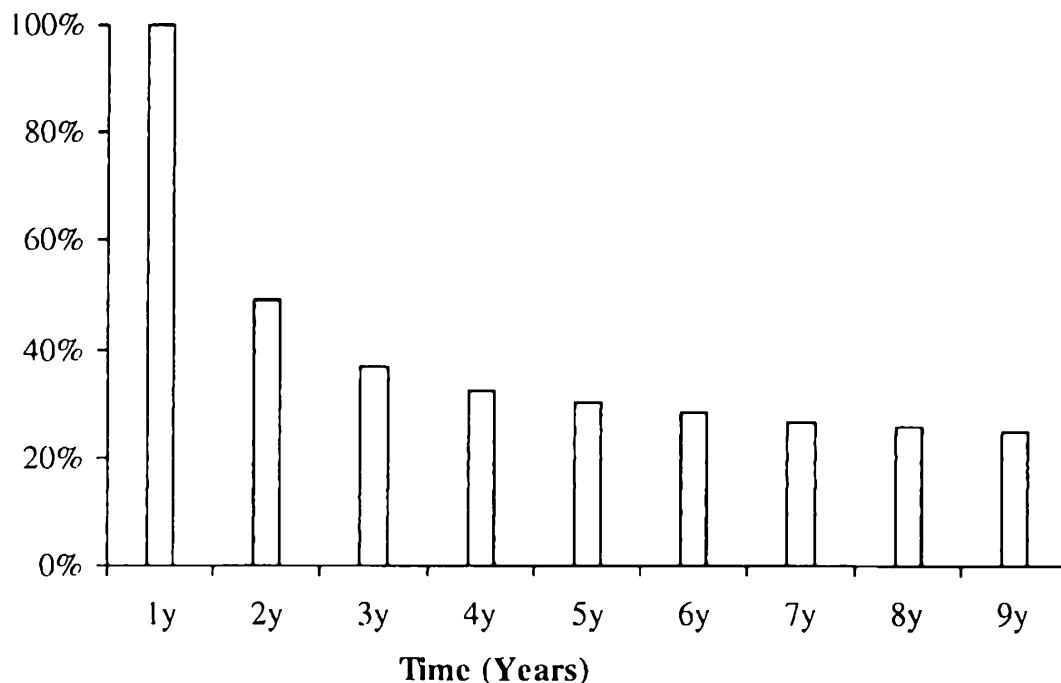
where we, arbitrarily, have used the spot measure numeraire $B(t)$ (and \mathbb{E} therefore denotes expectation in measure Q^B). We recall that a TARN typically pays fixed coupons to an investor before the knock-out feature starts; these coupons can be valued separately and are not included in the TARN definition above.

To make the discussion a bit more concrete, let us warm up by considering a typical example. Let the total maturity T_N be 10 years, let the target return R be 3%, and let the strike s and gearing g in (20.1) be 11.5% and 2, respectively. Also suppose the TARN pays annual coupons ($\tau_n = 1$ year). Using a yield curve with continuously compounded yields that grow from 3.5% in 1 day to 6.50% in 10 years and a displaced log-normal LM model with skew parameter 0.6 and calibrated to flat 35% swaption ATM Black volatilities, the value of the TARN with these parameters implies an attractive fixed coupon of 11% in the first year. If the TARN knocks out after the second year (at T_2), the investor would have received 14% return over two years (11% fixed coupon up front plus 3% targeted return), and is repaid the principal upon termination. This scenario comes true provided C_1 is above 3%, which according to (20.1) is equivalent to $L_1(T_1)$ fixing below 4.25%. More generally, the TARN will terminate early if interest rates are low. On the flip side, if, say, the rates go above 5.75% and stay there for the entire 10 year life of the TARN, all coupons C_n pay zero, and the investor receives nothing for 10 years. Yet, he has to pay Libor (by essentially, forfeiting interest on the principal) for 10 years, so the high-rate scenario is obviously not advantageous to the investor.

For reference, Figure 20.1 plots the probability (in spot measure) of the TARN being alive at future points in time, using the same market data and the model as above. According to the figure, the TARN stays alive for 10 years (bad for the investor) with about 25% probability, and knocks out after the first two years (good for the investor) with about 65% probability. Loosely speaking, the TARN investor therefore makes good money with (risk-neutral) probability of 65%, and loses a significant amount

with probability of 25%. This demonstrates how a high leverage inherent in TARNs allows them to pay attractive (i.e., high) coupons in scenarios that favor the investor. The leverage in any particular TARN depends on many factors, but is primarily a function of the target return R , with TARNs having smaller target return R providing higher leverage, *ceteris paribus*.

Fig. 20.1. Probability of TARN Being Alive at Future Years



Notes: Model-implied spot measure probability of a TARN being alive after a given number of years. The TARN contract details and the model are described in the text.

20.1.2 Valuation and Risk with Globally Calibrated Models

Using a flexible model (e.g. a Libor market model or a multi-factor quasi-Gaussian model) calibrated to the full swaption volatility grid is always a relatively safe choice for TARN pricing, since a globally calibrated multi-factor model can typically be counted on to capture the majority of all possible market risk factors. As shall be explained later, faithful reproduction of volatility smiles of various Libor rates is important for TARN valuation, so among all possible LM or qG models, we recommend the versions with stochastic volatility (Sections 14.2.5 and 13.3.2), as these models have enough flexibility to provide good fits to volatility smiles for a collection of forward Libor rates.

Pricing TARNs in LM and qG models is conceptually straightforward: as TARNs are purely path-dependent derivatives with no optimal exercise

features, standard Monte Carlo simulation techniques apply. However, since the knock-out feature of the TARN will introduce “digital” discontinuities into the payoff, Monte Carlo errors of the contract value and, especially, its risk sensitivities can be quite large, see Section 3.3.1. The number of paths required to get reasonably accurate estimates of risk sensitivities of a derivative with a discontinuous payoff could be high, which may sometimes render the application of a Monte Carlo based market model impractical. We review methods that help us obtain risk sensitivities in Monte Carlo for TARNs later in the book, see Sections 23.2.4, 23.4.4, and 25.2. Ultimately, however, the full power of LM and qG models may not be required for TARNs; despite appearances, TARNs turn out to be relatively simple instruments amendable to treatment by less complex — and more performant — models. We pursue this topic next.

20.1.3 Local Projection Method

In Chapters 18 and 19 we introduced the local projection method and applied it to Bermudan swaptions and other callable Libor exotics. We recall that the method is based on finding a relatively simple, local model that is calibrated to a global model (such as an LM model) in such a way as to approximate the value of the global model for a particular derivative. In particular, the local model should be calibrated to the parts of the global model volatility structure that are relevant to the derivative being valued. Let us apply this approach to TARNs. To start, it is informative to rewrite the TARN value as follows,

$$V_{\text{tarn}}(0) = E \left(\sum_{n=1}^{N-1} B(T_{n+1})^{-1} \tau_n \left((s - g \times L_n(T_n))^+ - L_n(T_n) \right) \times 1_{\{\sum_{i=1}^{n-1} \tau_i (s - g \times L_i(T_i))^+ < R\}} \right) \quad (20.3)$$

(with the usual convention that $\sum_{i=1}^0 = 0$). Scrutinizing the payoff, we notice that it depends on the values

$$\tilde{L} = (L_1(T_1), L_2(T_2), \dots, L_{N-1}(T_{N-1}))$$

of Libor rates on their fixing dates *only* (for the discrete money market numeraire $B(t)$ this follows from (4.24)). With the values of Libor rates at intermediate times irrelevant, only the distribution properties of the $(N-1)$ -dimensional vector \tilde{L} must be captured in whatever model we decide to use. Clearly this is a major simplification from a typical valuation problem. Notice, for instance, that a Bermudan swaption would depend on values of Libor rates at various dates on *and before* their fixing dates. A similar principle also holds approximately true for more complicated TARNs linked

to swap rates (rather than Libor rates): only the distribution properties of the $(N - 1)$ -dimensional vector of swap rates observed on their fixing dates needs to be captured. In stating this principle, we have relied on the fact that the dependence on Libor rates through discounting with the spot numeraire is rather mild and has only limited impact on the value of TARNs.

Focusing on the covariance characteristics only (we will deal with volatility smiles later), and assuming log-normal distributions for market rates for the time being, we see that if two models assign the same values to the term variances of Libor rates $\text{Var}(\ln L_n(T_n))$, $n = 1, \dots, N - 1$, and inter-temporal correlations of Libor rates $\text{Corr}(\ln L_n(T_n), \ln L_m(T_m))$, $n, m = 1, \dots, N - 1$, then the values of a TARN in the two models would be the same. With this in mind, we can apply the local projection method as follows. First, we calibrate, say, a Libor market model to the full swaption volatility grid (and, of course, one's views on the proper dynamics of the volatility structure). Second, we use the calibrated LM model to calculate the relevant term volatilities and inter-temporal correlations needed for the TARN. Third, we pick a simpler model and calibrate it to the volatilities and correlations extracted from the LM model. Finally, we use the calibrated local model for valuing the TARN. Of course, when computing risk sensitivities, we would update the volatilities and correlations produced by the global LM model for each shock of market data.

In the procedure above, the local model needed for the third and final steps needs enough flexibility in its volatility structure specification to calibrate to the set of TARN volatility information we identified earlier. Fortunately, the set is not very extensive and, as we have seen before, can be effectively captured even by models as simple as a one-factor Gaussian model, see Sections 13.1.7 and 13.1.8.3. While adequate for capturing the volatility structure, the smile capabilities of the Gaussian model are, however, quite limited, and we shall consider more advanced alternatives below.

20.1.4 Volatility Smile Effects

To investigate the effects of the volatility smile on TARNs, let us consider the TARN value on date T_1 as a function of the Libor rate $L_1(T_1)$:

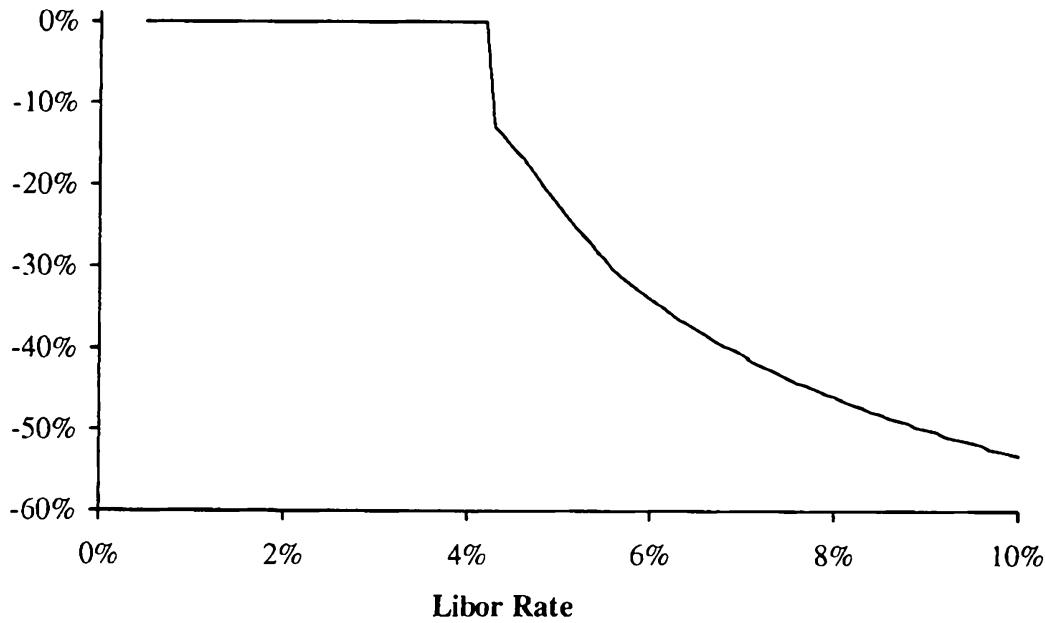
$$V_{\text{tarn}}(T_1, x) = E \left(B(T_1) \sum_{n=1}^{N-1} B(T_{n+1})^{-1} \tau_n (C_n - L_n(T_n)) \mathbf{1}_{\{Q_n < R\}} \middle| L_1(T_1) = x \right).$$

We plot $V_{\text{tarn}}(T_1, x)$ as a function of x in Figure 20.2 for the same TARN example and market/model data used in Section 20.1.1. Since $V_{\text{tarn}}(0)$ is given by the integral of $V_{\text{tarn}}(T_1, x)$ over the distribution of $L_1(T_1)$, the features of the payoff $V_{\text{tarn}}(T_1, x)$ highlight the characteristics of the distribution of $L_1(T_1)$ that are important for valuation. Clearly, $V_{\text{tarn}}(T_1, x)$ has an outright discontinuity at a barrier $L_1(T_1) = b_1$ implicitly given by

$$\tau_1 (s - gb_1)^+ = R, \quad (20.4)$$

and a call-option type singularity (a kink) at s/g . Moreover, values of future coupons are non-linear functions of $L_1(T_1)$, so the payoff $V_{\text{tarn}}(T_1, x)$ is non-linear in x . From the replication argument of Proposition 8.4.13, we recall that a model generally needs to faithfully incorporate the whole distribution of $L_1(T_1)$ as implied from caplet prices across a range of strikes, and not just some summary information such as an implied volatility at a certain strike.

Fig. 20.2. Value of TARN on First Knockout Date



Notes: Value of a TARN on the first knockout date as a function of the spot Libor rate on that date. TARN and model details are given in the text.

Some might argue for focusing all attention on the knockout barrier b_1 and disregarding the rest of the volatility smile, believing that to value a TARN properly, it suffices to choose a model that values a digital caplet with strike b_1 consistently with the market. While this argument has some merit for the first date T_1 , it is simply not valid for any subsequent knock-out dates. For instance, it is obvious that the value of $L_2(T_2)$ at which the derivative knocks out would depend on the realized fixing of $L_1(T_1)$, a value that is unknown at time $t = 0$. Since the location of the knock-out barrier at time T_2 is unknown at time $t = 0$, we cannot find a single strike that would faithfully represent relevant features of the volatility smile at T_2 ; thus a model that only matches the level, or slope, of the implied volatility of $L_2(T_2)$ at a single strike will be inadequate.

From this discussion, and from what we have learned about the local projection method in previous chapters, it is clear that a successful candidate for the local model should have the ability to calibrate to volatility smiles of all Libor rates, in addition to having a low number of state variables and enough flexibility to calibrate to inter-temporal correlations of Libor rates. One reasonable candidate is the one-factor quasi-Gaussian (qG) model with stochastic volatility (13.64). To calibrate this model, we would first fix its mean reversion function to match the inter-temporal correlations of Libor rates $\text{Corr}(\ln L_n(T_n), \ln L_m(T_m))$, $n, m = 1, \dots, N - 1$, as explained in Section 13.1.8.3. Subsequently, we would use the methods of Section 13.2 to calibrate to the volatility smiles of all Libor rates that appear in the payoff formula. With the formulas developed there, the time-dependent local volatility function $\sigma_r(t, x, y)$, and the time-dependent volatility of variance function $\eta(t)$ (see (13.64) for notation) could be chosen to match the implied SV parameters of relevant caplets.

The qG-SV model is a good choice of a local model for TARNs as it has just enough — but not more — flexibility to calibrate to all relevant covariance and smile information. Other suitable model candidates include the one-factor Markov-functional model from Appendix 11.A of Chapter 11, as it could be calibrated in a similar way. Finally, we could also use a two-factor version of the quadratic Gaussian model of Section 12.3. As all these three local models have sufficient flexibility to reproduce the TARN-specific correlation and smile properties of a globally calibrated multi-factor model, we would expect them to produce similar values for TARNs. Still, the models generate volatility smiles using different mechanisms, which may change correlations in subtle ways, as we saw an example of in Appendix 17.A of Chapter 17. While this effect on TARNs is quite minor, it can be significant for other classes of derivatives, see Section 20.2.4 below.

20.1.5 PDE for TARNs

Suppose that we succeeded in applying the local projection method to calibrate a low-dimensional Markov model targeted to TARN valuation. Actual pricing of the TARN structure could then obviously be accomplished by standard Monte Carlo techniques. As low-dimensional models typically allow for particularly efficient path discretization, the resulting scheme would be substantially faster than, say, simulation of a full Libor market model. Even more attractive, if the Markov model has a dimension less than 3 or 4, the local projection method allows for the usage of PDE-based TARN valuation schemes. In a finite-difference setting, the path-dependent nature of TARNs can be dealt with using the now-familiar method of augmenting the state variable space, see for instance Sections 2.7.5 and 18.4.5. In doing so, we shall implicitly assume that all C_n 's are non-negative, as is almost universally the case for structured notes.

Let $V(t, I)$ be the value of the TARN at time t assuming that the total accumulated coupon at time t is $I = I(t)$, where we have defined

$$I(t) = \sum_{i=1}^{N-1} \tau_i C_i(T_i) 1_{\{T_i \leq t\}}.$$

We start by initializing the value of the TARN at T_N to 0,

$$V(T_N^-, I) = 0. \quad (20.5)$$

Then, for each $n = N - 1, \dots, 0$, we perform the following steps.

1. Roll back the value of the TARN

$$V(T_n, I) = E_{T_n} \left(\frac{B(T_n)}{B(T_{n+1})} V(T_{n+1}^-, I) \right)$$

by solving an appropriate model PDE for each value of I .

2. Apply the continuity condition

$$V(T_n^-, I) = V(T_n, I + \tau_n C_n)$$

across I -planes, corresponding to the update of the total return Q_n at time T_n .

3. Add the time T_n coupon times the survival indicator,

$$V(T_n^-, I) = V(T_n^-, I) + P(T_n, T_{n+1}) \tau_n (C_n - L_n(T_n)) 1_{\{I < R\}}. \quad (20.6)$$

4. Starting from the new terminal condition (20.6), repeat Steps 1 through 3 with $n \rightarrow (n - 1)$.

The final value is given by $V_{\text{tarn}}(0) = V(0, 0)$.

The discretization scheme based on this algorithm would require specifying bounds, and potentially the discretization grid, for the extra state variable I . The lower bound for I is clearly 0. From (20.6) it follows that the coupon update equation for $I > R$ is trivial so one would think that the upper bound for I should be R . Yet if we look at the continuity condition in Step 2 we see that, on the right-hand side, we may actually need the values of $V(T_n, I)$ for $R < I < R + \tau_n C_n$. Hence the upper bound should be somewhat higher than R and is best determined as

$$R + \max_{n, \omega} \{\tau_n C_n(\omega)\}, \quad (20.7)$$

where by $C_n(\omega)$ we denoted the value of n -th coupon in the TARN over a realization of an interest rate path ω . For an inverse floater TARN, as well as for many other TARN types, the coupons are globally bounded and the expression in (20.7) makes sense. For TARNs with unbounded coupons

this strategy will obviously not work and the global maximum will need to be replaced with a maximum over a set of sufficiently high probability. Needless to say rather crude calculations are sufficient here. For example, if the coupon C_n is a deterministic function of a Libor rate $L_n(T_n)$, then the, say, 99% confidence interval for $L_n(T_n)$ could be established from its forward value and market-observed volatility; this confidence interval on the rate could then be translated into a confidence interval on the coupon.

While trade-specific analysis for the bounds of I is not conceptually difficult, it is always preferable to have a generic scheme that works for a large class of instruments. For example, one could imagine a simple scheme where, prior to solving a PDE, a Monte Carlo simulation with a low number of paths is run and an empirical distribution of I is estimated. Using this distribution, not only the bounds of given probabilistic coverage on I could be established, but we could also use it to set up a discretization grid. For example, we can discretize more finely in the region where realized values of I are dense, and use coarse discretization elsewhere to save calculation time.

Section 5.15.2 lists a few potential tweaks to the standard TARN specification; they could be included in the PDE scheme without much difficulty. For example, to include the “capped at trigger” feature we would replace (20.6) with

$$\begin{aligned} V(T_n-, I) &= V(T_n-, I) \\ &+ P(T_n, T_{n+1}) \tau_n (\min(C_n, (R - I)) - L_n(T_n)) 1_{\{I < R\}}. \end{aligned}$$

And to account for the “make whole” provision we would replace the initialization (20.5) with

$$V(T_N-, I) = (R - I)^+.$$

20.2 Volatility Swaps

We now shift our attention to volatility swaps introduced in Section 5.16. Valuation of many flavors of volatility swaps is straightforward in Monte Carlo, so a globally-calibrated LM model is a reasonable choice¹. As always, however, performance considerations suggest that we seek methods that are faster.

¹There is some evidence that many participants in the volatility swap market tend to use fairly naive, low-dimensional models for valuation. As a result, if correlations for the LM model are extracted from spread options, say, the LM model may produce forward volatilities that are lower than the market consensus (see Appendix 19.A for the rationale). Any arbitrages induced by such “segmentation” between the markets for spread options and volatility swaps are hard to exploit in practice, so the market differences can be quite persistent.

20.2.1 Local Projection Method

We recall that the structured coupon paid at T_{n+1} of a typical volatility swap has the form

$$C_n = |S_{n+1}(T_{n+1}) - S_n(T_n)|, \quad (20.8)$$

where $S_n(t)$, $n = 1, \dots, N$, are the reference rates of the swap. The payoff of a volatility swap shares certain characteristics with a TARN payoff, and this makes it amendable to the same treatment as what we applied to TARNs. In particular, it is clear from the valuation equation (5.26) that the value of a volatility swap depends on the values of the rates S_n on their fixing dates only. As such, the specific local projection method developed in Section 20.1.3 may be applied to volatility swaps as well. We do not repeat the analysis here, but just emphasize that we can use one-factor models to value volatility swaps as long as we calibrate them to the marginal rate distributions and the correlation structure of $(S_1(T_1), \dots, S_N(T_N))$. The former typically would come from the market and the latter from a globally calibrated model, e.g., the LM model.

As a properly calibrated one-factor model is appropriate for valuation, one may wonder whether we can use PDE, rather than Monte Carlo, methods. Indeed, this is the case, as each “swaplet” (20.8) in the structured leg of the volatility swap can be valued in a finite difference grid by introducing of an extra state variable to track the “strike” $S_n(T_n)$ in the swaplet payoff. In this particular case, the extra state variable method amounts to calculating, via a PDE, the value of the coupon at time T_n ,

$$\begin{aligned} V_{\text{swaplet}}(x, T_n) &\triangleq \mathbb{E}_{T_n}^{T_{n+1}} (C_n | S_n(T_n) = x) \\ &= \mathbb{E}_{T_n}^{T_{n+1}} (|S_{n+1}(T_{n+1}) - x| | S_n(T_n) = x), \end{aligned} \quad (20.9)$$

for a selection of values x (we use T_{n+1} -forward measure in this example). To obtain the time 0 value of the coupon, we can then calculate, again in a PDE, the expected value

$$\mathbb{E} (\beta(T_n)^{-1} P(T_n, T_{n+1}) V_{\text{swaplet}}(S_n(T_n), T_n)), \quad (20.10)$$

where $\beta(t)$ is the money market numeraire.

For some versions of the volatility swap payoff, we can go further and derive approximate closed-form expressions. We will discuss this in more detail later but, briefly, the basic approach here is to calculate the value of each swaplet payout with a two-dimensional integration or a suitable approximation for a spread option. In doing so, we rely on a model to pre-calculate the term volatilities and the (inter-temporal) correlation of $S_n(T_n)$ and $S_{n+1}(T_{n+1})$ in (20.8).

20.2.2 Shout Options

As pointed out in Section 5.16.2, volatility swaps often give the receiver of the structured coupons an option to *shout*, i.e. to choose the observation time of the rate $S_{n+1}(\cdot)$ in (20.8). The coupon in (20.8) is then replaced with

$$C_n = |S_{n+1}(\eta_n) - S_n(T_n)|,$$

where the stopping time $\eta_n \in [T_n, T_{n+1}]$ is chosen by the party receiving the coupon². Importantly, the payoff is still paid at time T_{n+1} , even if η_n is strictly less than T_{n+1} . As viewed from time T_n , the option then looks like an American option on an at-the-money straddle, with exercise value $P(\eta_n, T_{n+1})|S_{n+1}(\eta_n) - S_n(T_n)|$; the presence of the discount factor $P(\eta, T_{n+1})$ in the exercise value reflects the fact that payment of the coupon will always take place at time T_{n+1} , irrespective of the time of exercise. Notice that $S_{n+1}(t)$ has almost no drift³ so Jensen's inequality implies that

$$\begin{aligned} P(\eta_n, T_{n+1}) \mathbb{E}_{\eta_n}^{T_{n+1}} (|S_{n+1}(T_{n+1}) - S_n(T_n)|) \\ \geq P(\eta_n, T_{n+1}) |\mathbb{E}_{\eta_n}^{T_{n+1}} (S_{n+1}(T_{n+1})) - S_n(T_n)| \\ \approx P(\eta_n, T_{n+1}) |S_{n+1}(\eta_n) - S_n(T_n)|, \end{aligned}$$

i.e. the exercise value is (approximately) dominated by the hold value and the value of the early exercise is negligible. As a consequence, the shout option can safely be ignored for valuation purposes, and the coupons could be valued as if (20.8) were the actual payoff.

The situation is somewhat less clear cut in a reasonably popular case of a *capped coupon* with a shout,

$$C'_n = \min (|S_{n+1}(\eta_n) - S_n(T_n)|, c) \quad (20.11)$$

for some $c > 0$. Clearly, if for some $t \in [T_n, T_{n+1}]$ the rate $S_{n+1}(t)$ is outside the interval $[S_n(T_n) - c, S_n(T_n) + c]$, then the holder should exercise at that point as he will never get a higher value for the coupon (but if he waits he may end up with a lower value at expiration). So early exercise is optimal in some cases. This may seem like a major complication as any application of Monte Carlo would apparently require the estimation of the optimal exercise rule, potentially requiring the full suite of regression-based methods of Chapter 18. Fortunately the situation is much simpler; we formalize this result as a proposition.

²Sometimes the coupon is linked to a swap rate that starts at shout time η_n rather than at the end of the period T_{n+1} . This modification to our discussion is easy to incorporate, and we do not consider this case separately.

³In the T_{n+1} -forward measure, $S_{n+1}(t)$ will typically have a small convexity-induced drift, as $S_{n+1}(t)$ here represents some CMS rate. However, the period $[T_n, T_{n+1}]$ rarely exceeds one year, and over one year the drift of a CMS rate will normally be quite close to zero.

Proposition 20.2.1. *The value of the American option on a capped straddle with the payoff (20.11) is equal to the value of a straddle with a barrier, so that $E_{T_n}^{T_{n+1}}(C'_n) = E_{T_n}^{T_{n+1}}(C''_n)$, where*

$$\begin{aligned} C''_n &= c \times 1_{\left\{\max_{t \in [T_n, T_{n+1}]} (|S_{n+1}(t) - S_n(T_n)|) \geq c\right\}} \\ &\quad + |S_{n+1}(T_{n+1}) - S_n(T_n)| \times 1_{\left\{\max_{t \in [T_n, T_{n+1}]} (|S_{n+1}(t) - S_n(T_n)|) < c\right\}}. \end{aligned}$$

Remark 20.2.2. The proposition tells us that the optimal exercise strategy is known analytically: for the period $[T_n, T_{n+1}]$ one should simply exercise the shout option on the first time t when $S_{n+1}(t)$ hits either of the barriers $S_n(T_n) \pm c$.

Proof. We content ourselves with a sketch of the proof; a more formal argument is developed in Broadie and Detemple [1995]. Let us denote by η_n^c the first hitting time of the double barrier $S_n(T_n) \pm c$,

$$\eta_n^c = \inf \{t \in [T_n, T_{n+1}] : |S_{n+1}(t) - S_n(T_n)| = c\} \wedge T_{n+1}.$$

Then, clearly,

$$E_{T_n}^{T_{n+1}}(C''_n) = E_{T_n}^{T_{n+1}}(\min(|S_{n+1}(\eta_n^c) - S_n(T_n)|, c))$$

(we use T_{n+1} -forward measure for valuation here). On one hand,

$$\begin{aligned} E_{T_n}^{T_{n+1}}(\min(|S_{n+1}(\eta_n^c) - S_n(T_n)|, c)) \\ \leq E_{T_n}^{T_{n+1}}(\min(|S_{n+1}(\eta_n) - S_n(T_n)|, c)) \end{aligned} \quad (20.12)$$

because, by definition, η_n is the *optimal* stopping time for the American capped straddle. On the other hand, for each $t \in [T_n, T_{n+1}]$,

$$E_t^{T_{n+1}}(\min(|S_{n+1}(\eta_n^c) - S_n(T_n)|, c)) \geq \min(|S_{n+1}(t) - S_n(T_n)|, c)$$

as Figure 20.3 demonstrates. Therefore,

$$\begin{aligned} &E_{T_n}^{T_{n+1}}(\min(|S_{n+1}(\eta_n) - S_n(T_n)|, c)) \\ &= E_{T_n}^{T_{n+1}} \left(\int_{T_n}^{T_{n+1}} \min(|S_{n+1}(t) - S_n(T_n)|, c) \delta(\eta_n - t) dt \right) \\ &\leq E_{T_n}^{T_{n+1}} \left(\int_{T_n}^{T_{n+1}} E_t^{T_{n+1}}(\min(|S_{n+1}(\eta_n^c) - S_n(T_n)|, c)) \delta(\eta_n - t) dt \right) \\ &= E_{T_n}^{T_{n+1}} \left(\min(|S_{n+1}(\eta_n^c) - S_n(T_n)|, c) \int_{T_n}^{T_{n+1}} \delta(\eta_n - t) dt \right) \\ &= E_{T_n}^{T_{n+1}}(\min(|S_{n+1}(\eta_n^c) - S_n(T_n)|, c)), \end{aligned} \quad (20.13)$$

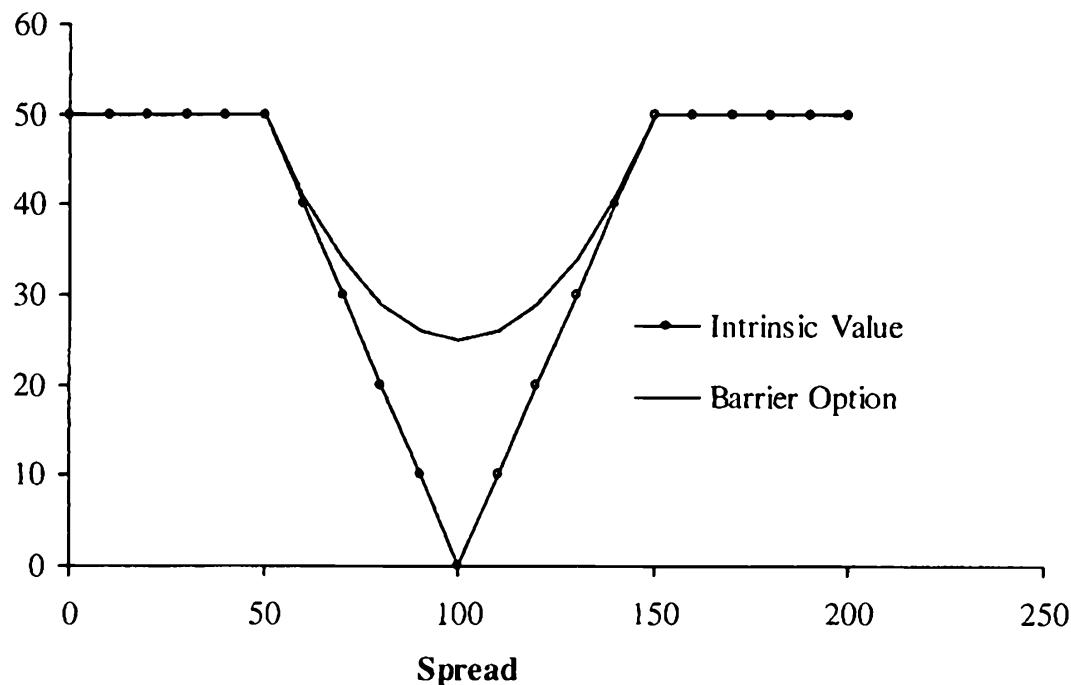
where the second-to-last equality follows by the law of iterated conditional expectations and \mathcal{F}_t -measurability of $\delta(\eta_n - t)$. Comparing (20.12) and (20.13) we see that

$$\mathbb{E}_{T_n}^{T_{n+1}}(C'_n) = \mathbb{E}_{T_n}^{T_{n+1}}(C''_n)$$

and the optimal exercise strategy is actually given by η_n^c , as stated earlier.

□

Fig. 20.3. Value of a Barrier Option on a Capped Straddle



Notes: The payoff and the present value of a barrier option on a capped straddle, vs. the underlying. The barrier option value dominates the payoff in all states of the world.

Having replaced an American option with a barrier option, we can now value capped volatility swaps in standard Monte Carlo, as no optimal exercise features need to be incorporated. Still, we do have some complications — namely two continuous barriers — that now need to be handled in Monte Carlo. Here the techniques developed in Section 3.2.9 come in handy. Deserving a special mention is the method of Broadie et al. [1997] which replaces a continuously-observed barrier with a discretely-observed one that is shifted by a certain amount, see Theorem 3.2.2. Other methods from Section 3.2.9 could also be used, although all of them assume that the underlying follows a simple process such as a Brownian motion; this is typically not a problem as the dynamics of $S_{n+1}(t)$ for $t \in [T_n, T_{n+1}]$ could be closely approximated as such, irrespective of the underlying model used, since $T_{n+1} - T_n$ tends to be relatively small (one year or less).

PDE methods can also be used for barrier options, with the same trick of using the strike $S_n(T_n)$ as an extra state variable as we discussed in the European case (20.9)–(20.10), only now valuing a (double) barrier option for each value of the strike $S_n(T_n)$. This is probably the best we can do as far as valuation speed is concerned, as it should be obvious that closed-form approximations to values of capped straddle coupons with a shout option would be rather hard to develop.

20.2.3 Min-Max Volatility Swaps

Defined in Section 5.16.3, min-max volatility swaps replace the straddle coupon (20.8) with a coupon that measures the maximum move of a given rate over a given period,

$$C_n = M_n - m_n, \quad (20.14)$$

where

$$M_n = \max_{s \in [T_n, T_{n+1}]} S_n(s), \quad m_n = \min_{s \in [T_n, T_{n+1}]} S_n(s).$$

At first glance the min-max coupon appears significantly more “exotic” than the plain straddle (20.8), to the point that one may wonder if the min-max coupon has significantly different risk characteristics. For example, a superficial analysis could suggest that the min-max coupon has significantly higher “forward skew” exposure, i.e. the exposure to the slope of the volatility smile as observed at the time when the minimum (or maximum) of the process is reached. It is all the more interesting that, in fact, we can show that the two coupons, (20.14) and (20.8), are quite alike.

Let us start by assuming that for $t \in [T_n, T_{n+1}]$, $S_n(t)$ follows a driftless Brownian motion with some constant volatility σ , i.e.

$$dS_n(t) = \sigma dW(t),$$

where $W(t)$ is a Brownian motion in the T_{n+1} -forward measure $Q^{T_{n+1}}$. As we have already pointed out, this is not a bad approximation as the period $[T_n, T_{n+1}]$ is often rather short, on the order of one year. By the reflection principle from Section 2.6 of Karatzas and Shreve [1997], we have for any $b \geq S_n(T_n)$,

$$Q_{T_n}^{T_{n+1}} (M_n \geq b) = 2Q_{T_n}^{T_{n+1}} (S_n(T_{n+1}) \geq b). \quad (20.15)$$

As

$$\int_{b_{\min}}^{b_{\max}} 1_{\{x \geq b\}} db = \min(\max(x - b_{\min}, 0), b_{\max} - b_{\min})$$

we obtain, integrating (20.15) over $b \in [S_n(T_n), \infty)$ that

$$E_{T_n}^{T_{n+1}} (M_n - S_n(T_n)) = 2E_{T_n}^{T_{n+1}} ((S_n(T_{n+1}) - S_n(T_n))^+).$$

Similarly,

$$\mathbb{E}_{T_n}^{T_{n+1}} (S_n(T_n) - m_n) = 2\mathbb{E}_{T_n}^{T_{n+1}} \left((S_n(T_n) - S_{n+1}(T_n))^+ \right)$$

and, adding the last two equations together, we obtain

$$\mathbb{E}_{T_n}^{T_{n+1}} (M_n - m_n) = 2\mathbb{E}_{T_n}^{T_{n+1}} (|S_n(T_{n+1}) - S_n(T_n)|). \quad (20.16)$$

Therefore, the value of the min-max coupon is (approximately) equal to twice the value of the straddle coupon, and the min-max volatility bond could be valued, and risk managed, in the same way as the standard volatility swap in Section 20.2.

The starting point of our proof, equation (20.15), has an interesting financial interpretation. On the left-hand side we have the value of a “one-touch” option, an option that pays 1 if the underlying process ever touches the barrier b (before T_{n+1}). The equality suggests that this continuously monitored barrier option can somehow be hedged with two European-style digital call options (the right-hand side). To show that this is indeed the case, consider buying two digital calls struck at the barrier. If the underlying process never hits the barrier, both the one-touch and the two digitals expire worthless. On the other hand, if the process touches the barrier, then we sell one of the digital calls and buy one digital put, i.e. an option that pays 1 at T_{n+1} if and only if $S_n(T_{n+1}) \leq b$. The value of the digital call and the digital put are the same due to the (assumed) symmetry of the Brownian motion process for S_n ; hence we can trade at zero cost. After the trade, our replicating portfolio consists of one digital call and one digital put, which will produce a payoff of 1 irrespective of the final value of the process $S_n(T_{n+1})$. Note that this is exactly equal to the payoff of the one-touch in this case. Therefore the one-touch and the replicating portfolio have the same payoffs in all states of the world, as claimed. The replicating portfolio (or the inverse of it, a hedging portfolio) is called *semi-static* to reflect the fact that the replicating strategy may involve some (costless) trading activity during the life of the trade.

In deriving (20.16) we represented the min-max payoff as a (continuous) integral of one-touch payoffs. It should then come as no surprise that we can set up a semi-static replicating portfolio for a min-max coupon that starts with two European straddles (see footnote 22 in Section 5.16.3). We leave it as an exercise to the reader to write down the explicit trading strategy for the replication; as a hint we mention that it involves holding at each time $t \in [T_n, T_{n+1}]$ a portfolio with the payoff

$$|S_n(T_{n+1}) - M_n(t)| + |m_n(t) - S_n(T_{n+1})|,$$

where $M_n(t)$, $m_n(t)$ are the *running* maximum and minimum,

$$M_n(t) = \max_{s \in [T_n, t]} S_n(s), \quad m_n(t) = \min_{s \in [T_n, t]} S_n(s).$$

The replication of a min-max coupon with two standard straddles is not model-independent, as it relies on approximating the process for the rate $S_n(t)$ with a (driftless) Gaussian process. As a result of this assumption, ATM puts and calls will have identical prices, a relationship often known as *arithmetic put-call symmetry*. The hedging arguments above can be extended to all processes for which arithmetic put-call symmetry holds at a barrier hitting time, i.e. processes for which the distribution of $S_n(T_{n+1})$ observed at any stopping time in $[T_n, T_{n+1}]$ is symmetric. In some settings, it is most useful to assume *geometric put-call symmetry*, which essentially means that the Black volatility smile is symmetric in log-moneyness; a simple example of a process satisfying this assumption is the geometric Brownian motion process without drift, or the (drift-free) Heston model with zero asset/volatility correlation. It is not difficult to prove that a semi-static hedge also exists for this case, although the hedge is somewhat more complicated than two straddles (European options at a full continuum of strikes are required). For further discussions on the topic, see Carr and Lee [2009a] which surveys (and generalizes) the considerable amount of work in the literature on applications of put-call symmetry.

While one can experiment with various assumptions to find more accurate valuation formulas, ultimately the main utility of the result such as (20.16) lies in demonstrating that the risk characteristics of a min-max volatility swap are largely the same as those of a standard volatility swap. Such qualitative understanding is useful when making model selections, even when direct replication arguments no longer work. As a typical example we can mention the *capped* min-max volatility swap, a swap with coupons that have a payoff

$$C_n = \min(M_n - m_n, c),$$

for some $c > 0$.

20.2.4 Impact of Volatility Dynamics on Volatility Swaps

With the analysis above as background, let us now ponder the question of what is, ultimately, the appropriate model for volatility swaps. While the discussion of Section 20.2.1 has made it clear that (properly calibrated) single-factor models can be safely used, we still need to decide on other features of potential models, such as a faithful reproduction of volatility smile and, perhaps, its dynamics. As always, we look for a model that captures the main risk factors for a given type of derivatives, yet avoids introducing complicated features that may not be relevant.

To show that the model choice is not entirely trivial we point to Figure 20.4. Here, we have plotted the value of volatility swaplets for *fixed-tenor* and *fixed-expiry* volatility swaps (see Section 5.16.1) in three different models. Both swaps are of 10 year maturity and have annual coupons. For the fixed-tenor swap the underlying rate is the 10 year CMS rate; for the fixed-expiry

swap it is a 10 year swap rate fixing in 10 years time. For model calibration we use Euro market data from the summer of 2008. The three models are: i) the SV version of the quasi-Gaussian (qG) model as in Section 13.2; ii) the two-factor quadratic Gaussian (QG) model from Section 12.3; and iii) a local volatility version of the qG model, with local volatility a quadratic function of the short rate (not something we normally recommend, see Section 13.1.5) and time. All three models have been calibrated to volatility smiles of the 20 year coterminal swaption strip. To give a sense of market data used in calibration, the vanilla SV model (see Section 16.1.3) used to mark swaptions was set up to have the volatility between 14% and 16%, the skew between -10% and 10% and the volatility of variance between 100% and 200%, with mean reversion of volatility at 20%. All three models have the same mean reversion parameter of 2% in a (loose) attempt to make the inter-temporal correlations of relevant swap rates invariant across models.

We see that the differences in the values of individual coupons for the three models are quite significant. As we have calibrated the models to the same spot market data (volatility smiles), we must conclude that it is the different dynamics of the volatility structure (and, perhaps, volatility smiles) in the three models that are responsible for the valuation differences. In fact, we will argue that the difference in multi-dimensional distributions of the swap rates in the three models lead to differences in the meaning of the mean reversion parameter which imply different forward volatilities in different models. To understand this better, let us consider the issue of pricing an individual coupon in more details.

Let $S(t)$ be a forward swap rate that corresponds to a swap that starts at time T (with some, unspecified, maturity). Also, define $A(t)$ to be the corresponding annuity. We consider a contract that pays

$$|S(t) - S(u)| \quad (20.17)$$

at time t , where $0 < u < t \leq T$, a contract commonly called a *forward CMS straddle*. This contract corresponds to a coupon of a fixed-expiry volatility swap, the first graph in Figure 20.4. Let us find an approximate expression for the value of (20.17) in order to study its dependence on the various market quantities. The value can be expressed in the annuity measure induced by S and equals

$$V(0) = A(0)E^A(|S(t) - S(u)| / A(t)).$$

Using the linear TSR model of Section 16.6.4 we obtain

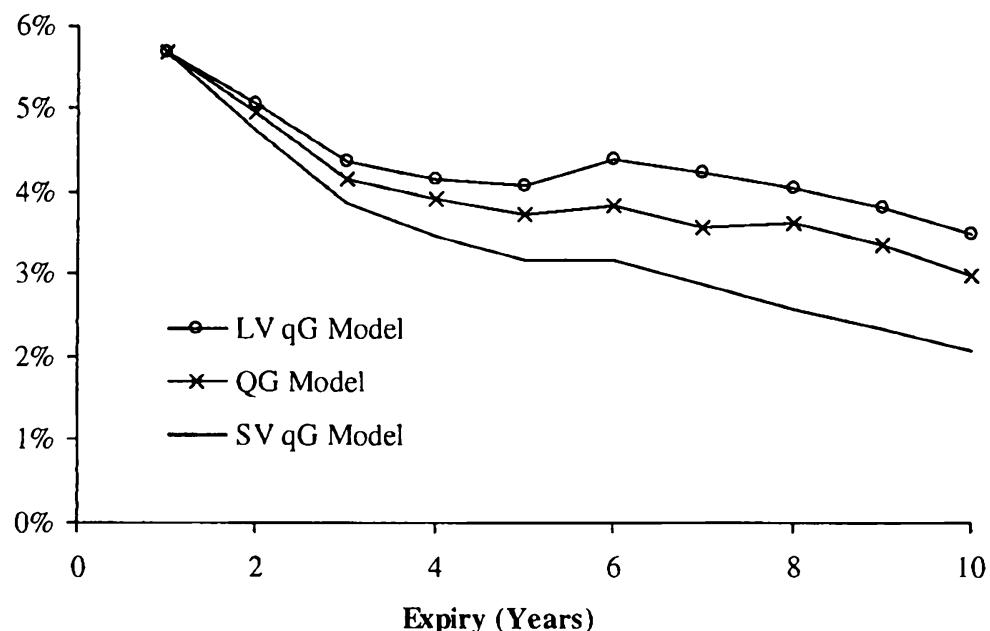
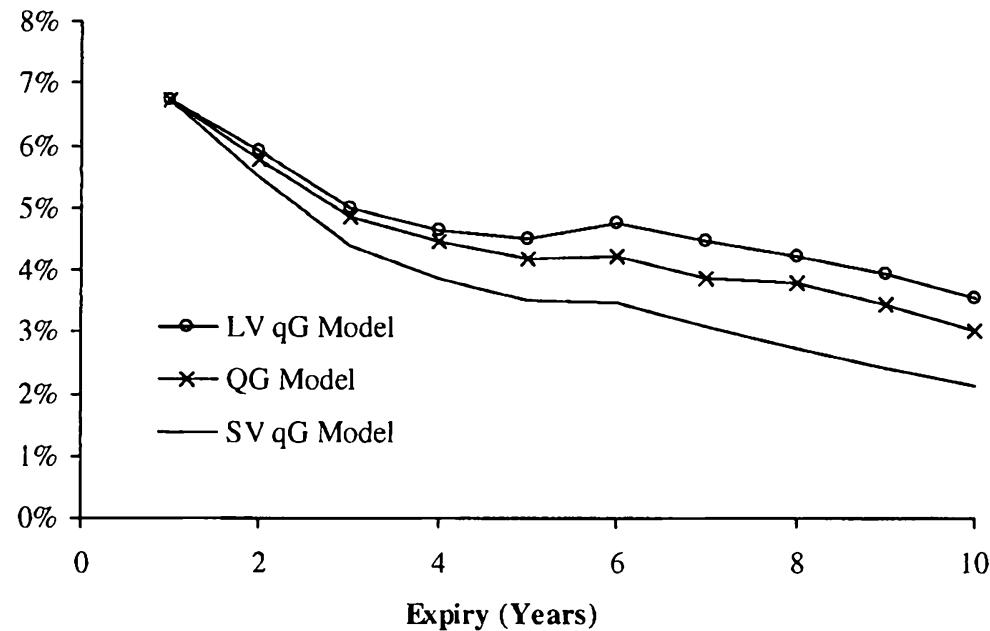
$$V(0) = A(0)E^A\left(|S(t) - S(u)| \times \left(\frac{1}{A(0)} + \alpha_1(S(t) - S(0))\right)\right)$$

for some $\alpha_1 > 0$, so

$$V(0) = V_1 + V_2 + V_3, \quad (20.18)$$

where we have defined

Fig. 20.4. Values of Volatility Swaplets For Fixed-Tenor and Fixed-Expiry Volatility Swaps



Notes: Values of volatility swaplets for fixed-tenor (first panel) and fixed-expiry (second panel) 10 year maturity annual volatility swaps in three different models, as described in the text.

$$\begin{aligned}
 V_1 &= (1 - \alpha_1 S(0) A(0)) E^A (|S(t) - S(u)|), \\
 V_2 &= \alpha_1 A(0) E^A (|S(t) - S(u)| (S(t) - S(u))), \\
 V_3 &= \alpha_1 A(0) E^A (|S(t) - S(u)| S(u)).
 \end{aligned}$$

We shall see that V_1 is linked to the future volatility of the swap rate S , V_2 is linked to the future skew, and V_3 is largely determined by the convexity adjustment as defined by Section 16.6.

Let us take a closer look at the first term V_1 . Conditioning on the sigma-algebra at time u we obtain

$$\begin{aligned} \mathbb{E}^A(|S(t) - S(u)|) &= A(0)^{-1} \mathbb{E}(\beta(u)^{-1} A(u) \mathbb{E}_u^A(|S(t) - S(u)|)) \quad (20.19) \\ &= A(0)^{-1} \mathbb{E}(\beta(u)^{-1} V_{\text{straddle}}(u, S(u), t)), \end{aligned}$$

where

$$V_{\text{straddle}}(u, K, t) = A(u) \mathbb{E}_u^A(|S(t) - K|)$$

can be recognized as the time u value of a swaption straddle, i.e. a sum of K -strike payer and receiver European swaptions for delivery at time t of a swap that starts at $T \geq t$, as entered at time u . The quantity $V_{\text{straddle}}(u, S(u), t)$ is then the value at time u of the at-the-money (ATM) straddle. We note that the contract, entered at time 0, that pays

$$A(t) |S(t) - S(u)|$$

at time t is known as a *forward swaption straddle*. The forward swaption straddle differs from the forward CMS straddle in (20.17) in that it pays the difference of swap values rather than swap rates.

Let us denote by $\sigma_N(u, S(u); t, K)$ the value of the implied basis-point (or Normal) volatility of the swap rate S , as observed at time u for swaptions with expiry t and strike K (and again, on a swap that starts at some $T \geq t$). From the Bachelier pricing formula (7.16), we see that

$$V_{\text{straddle}}(u, S(u), t) = A(u) \sqrt{\frac{2\tau}{\pi}} \sigma_N(u, S(u); t, S(u)), \quad (20.20)$$

so the ATM straddle value is equal to the (scaled) value of the implied basis-point volatility. From (20.19) it follows that

$$V_1 = a_1 \mathbb{E}^A(\sigma_N(u, S(u); t, S(u)))$$

for some constant a_1 . As mentioned earlier, the value V_1 is given by the expected value of a future (at-the-money) basis-point volatility of a swap rate S over a period $[u, t]$.

As we have discussed before, the information about such future volatilities is, by and large, not contained in the market data available at time 0, but is mostly driven by the dynamics of a model used for valuation. While calibrated to the same marginal distributions at time 0, the three models used in Figure 20.4 have different dynamics and therefore different forward volatilities. Probably the easiest way to understand it is to recall (see Dupire [1997]) that the price of the contract that pays instantaneous forward variance is model-independent (as long as all European options are matched),

which implies, by Jensen's inequality, that the price of a contract that pays forward *volatility* will depend on the distribution (mainly variance) of the volatility itself.

We shall discuss forward swaption straddles in more detail in Section 20.3, but first we turn our attention back to the remaining two terms in the decomposition (20.18). We can rewrite the term V_2 as

$$V_2 = a_2 E^A(f(S(t) - S(u))),$$

where a_2 is some constant and $f(x) = x|x|$. The function $f(x)$ is concave for $x < 0$ and convex for $x > 0$, suggesting that the value V_2 is largely driven by *forward skew*, i.e. the slope of the volatility smile $\sigma_N(u, S(u); t, K)$ at time u . This is most easily seen from the replication result of Proposition 8.4.13. as we have

$$\begin{aligned} E^A(f(S(t) - S(u))) &= \\ &\int_0^\infty (E^A((S(t) - (S(u) + y))^+) - E^A(((S(u) - y) - S(t))^+)) dy, \end{aligned}$$

i.e. V_2 is a sum of forward call spreads for different strike offsets y from the time u at-the-money rate $S(u)$. The value of each call spread is largely (to the first order in volatility) determined by the difference in the appropriate implied volatilities,

$$\begin{aligned} E^A((S(t) - (S(u) + y))^+) - E^A(((S(u) - y) - S(t))^+) \\ \approx c \cdot (\sigma_N(u, S(u); t, S(u) + y) - \sigma_N(u, S(u); t, S(u) - y)), \end{aligned}$$

which is clearly related to the slope of the volatility smile as observed at time u . As is the case for forward volatility, forward skew is strongly model-dependent, which helps to explain the differences in forward CMS straddle values in Figure 20.4.

The third term, V_3 in (20.18) provides a relatively small contribution to the value of the forward CMS straddle. To analyze it in more detail, it is convenient to condition on $S(u)$ and use (20.20) to obtain

$$V_3 = a_3 E^A(S(u) \sigma_N(u, S(u); t, S(u)))$$

for some constant a_3 . For many models the ATM volatility $\sigma_N(u, S; t, S)$ is a linear function of S ,

$$\sigma_N(u, S; t, S) \approx a_4 + a_5 S,$$

so then

$$V_3 = a_3 a_4 S(0) + a_3 a_5 E^A(S(u)^2)$$

and we see that V_3 is mostly influenced by the convexity adjustment of $S(u)$ (see Section 16.6.4). The value of the convexity adjustment is linked to the

spot volatility information (volatility smile for options on $S(u)$ as observed at time 0), and models that are calibrated to the (spot) volatility smile on $S(u)$ should give identical values to this term. This is, for example, the case in Figure 20.4 for the fixed-tenor volatility swaps; in the fixed-expiry case the models would typically be calibrated to (spot) volatility smiles on $S(T)$ (T here is the start date for the swap) and may imply different smiles for time u . Anyway, as we already mentioned, numerical experiments show that the V_3 term is not particularly significant.

As the value of the forward CMS straddle is largely defined by the forward volatility term V_1 , it may seem puzzling that the values on display in Figure 20.4 are strongly sensitive to the model choice. After all, we used the same mean reversion for the three models (2%), which seemingly implies the same inter-temporal correlations between the swap rates involved, which, as we argued before, should essentially lock in the forward volatilities to the same levels in the three models. However, material differences here arise from the fact that different smile mechanisms in the three models in Figure 20.4 change the meaning of the mean reversion parameter. We have seen a similar effect before in Section 17.A that showed that the effective correlation in a displaced log-normal model depends on the skew. The same happens here: even though the mean reversion parameter is the same, the actual effective inter-temporal correlations (or, equivalently, forward volatilities) are different because of different multi-dimensional distributions of the swap rates involved (even though marginal distributions are nearly identical). To estimate these smile effects on correlation we would face the difficult task of going beyond the Gaussian-type approximations used in, e.g., Section 13.1.8.3.

Let us summarize. The values of forward CMS straddles — i.e. coupons of volatility swaps — depend on the level and shape of volatility smiles at future times. As these are largely defined by the volatility dynamics of the model used, different models can produce significantly different values, even if calibrated to identical (spot) volatility information. In particular the impact of volatility smiles on inter-temporal correlations/forward volatilities is of significant importance, and it is advisable to use models with different mechanisms of smile generation to monitor and control it.

20.3 Forward Swaption Straddles

Besides being closely linked to volatility swaps (see Section 20.2.4), forward swaption straddles are themselves traded as stand-alone products, with most of the demand coming from hedge funds interested in expressing views on future implied volatility (see (20.20)). In this stand-alone traded format forward swaption straddles tend to be relatively short-dated, with typical expiries around 3–5 years. As such, these securities are often treated as vanilla, rather than exotic, derivatives, and it is common to use simple

vanilla-type models for their valuation. We proceed to describe a typical approach.

First, to be able to use our standard swap rate notations (4.8), (4.10), we assume that a tenor structure

$$0 < T_0 < T_1 < \dots < T_N, \quad \tau_n = T_{n+1} - T_n,$$

is given. To introduce matching notation for forward straddles, let

$$V_{n,m}(v; T_k)$$

be the time v value of the forward straddle payoff

$$A_{n,m}(T_n) |S_{n,m}(T_n) - S_{n,m}(T_k)| \quad (20.21)$$

paid at T_n , $k < n$, $n + m < N$. We fix a particular forward straddle; to tie our discussion to Section 20.2.4, we assume that

$$u = T_s, \quad t = T = T_e$$

for some indices⁴ s, e , $0 \leq s < e < N$, and that the final payment date of the swap is T_N . So this contract delivers at time T_s the value of an at-the-money straddle with expiry T_e on a swap that starts on T_e and matures on T_N . In other words, the contract has the payoff

$$V_{\text{straddle}}(T_s, S_{e,N-e}(T_s), T_e)$$

paid at T_s or, equivalently, the payoff (20.21) with $k = s$, $n = e$, $m = N - e$. By (20.20), the value at time 0 is equal to

$$\begin{aligned} V_{e,N-e}(0; T_s) &= A_{e,N-e}(0) \sqrt{\frac{2(T_e - T_s)}{\pi}} \\ &\times E^A(\sigma_N(T_s, S_{e,N-e}(T_s); T_e, S_{e,N-e}(T_s))) \end{aligned} \quad (20.22)$$

(the measure Q^A here is actually $Q^{A_{e,N-e}}$ but we simplify the notation for brevity). As we already mentioned, in most models that we use the at-the-money basis-point volatility $\sigma_N(T_s, S; T_e, S)$ can be approximated to excellent precision as a linear function. This allows us to write

$$\begin{aligned} \sigma_N(T_s, S_{e,N-e}(T_s); T_e, S_{e,N-e}(T_s)) &\approx \sigma_N(T_s, S_{e,N-e}(0); T_e, S_{e,N-e}(0)) \\ &+ \frac{\partial}{\partial S} \sigma_N(T_s, S; T_e, S) \Big|_{S=S_{e,N-e}(0)} (S_{e,N-e}(T_s) - S_{e,N-e}(0)). \end{aligned}$$

As $S_{e,N-e}$ is a Q^A -martingale we then simply have

⁴Here s also stands for “strike setting” and e for “expiry”.

$$V_{e,N-e}(0; T_s) \approx A_{e,N-e}(0) \sqrt{\frac{2(T_e - T_s)}{\pi}} \lambda_{e,N-e}(T_s, T_e), \quad (20.23)$$

where we have abbreviated

$$\lambda_{e,N-e}(T_s, T_e) \triangleq \sigma_N(T_s, S_{e,N-e}(0); T_e, S_{e,N-e}(0))$$

and added subscripts to highlight the rate this volatility corresponds to. As discussed earlier, $\lambda_{e,N-e}(T_s, T_e)$ is the forward volatility of the swap rate $S_{e,N-e}$ over the period $[T_s, T_e]$. While this quantity cannot be observed directly in the market, it can be linked to quantities that can. Indeed, if we approximate the swap rate $S_{e,N-e}$ as following a Gaussian process in measure Q^A , then splitting the total variance of $S_{e,N-e}(T_e)$ over the periods $[0, T_s]$ and $[T_s, T_e]$ we get

$$\lambda_{e,N-e}(T_s, T_e)^2 (T_e - T_s) = \lambda_{e,N-e}(0, T_e)^2 T_e - \lambda_{e,N-e}(0, T_s)^2 T_s. \quad (20.24)$$

Clearly, $\lambda_{e,N-e}(0, T_e)$ is observable, since the value of a standard spot (time 0) starting at-the-money straddle with expiry T_e is given by

$$V_{e,N-e}(0; 0) \approx A_{e,N-e}(0) \sqrt{\frac{2T_e}{\pi}} \lambda_{e,N-e}(0, T_e).$$

The other term in (20.24), $\lambda_{e,N-e}(0, T_s)$, is the volatility of the swap starting at T_e over the period $[0, T_s]$. Equivalently, it is the volatility implied from the value of an *option on a forward starting swap*; at option expiry T_s the holder has the right to enter into a swap starting at some future time $T_e > T_s$. Such options are not traded (or, rather, not liquid), but note that a forward-starting swap can be represented as a combination of spot-starting swaps. We do the calculation for swap rates:

$$\begin{aligned} S_{e,N-e}(T_s) &= \frac{P(T_s, T_e) - P(T_s, T_N)}{A_{e,N-e}(T_s)} \\ &= \frac{P(T_s, T_s) - P(T_s, T_N)}{A_{e,N-e}(T_s)} - \frac{P(T_s, T_s) - P(T_s, T_e)}{A_{e,N-e}(T_s)} \\ &= w_1(T_s) S_{s,N-s}(T_s) - w_2(T_s) S_{s,e-s}(T_s), \end{aligned} \quad (20.25)$$

where

$$w_1(T_s) = \frac{A_{s,N-s}(T_s)}{A_{e,N-e}(T_s)}, \quad w_2(T_s) = \frac{A_{s,e-s}(T_s)}{A_{e,N-e}(T_s)}.$$

To proceed, we approximate the ratios of PVBP by their values at time 0,

$$w_1(T_s) \approx w_1 = \frac{A_{s,N-s}(0)}{A_{e,N-e}(0)}, \quad w_2(T_s) \approx w_2 = \frac{A_{s,e-s}(0)}{A_{e,N-e}(0)}, \quad (20.26)$$

and assume that $S_{s,N-s}(T_s)$, $S_{s,e-s}(T_s)$ are approximately Gaussian with correlation ρ . Then from (20.25) we obtain

$$\begin{aligned}\lambda_{e,N-e}(0, T_s)^2 &\approx w_1^2 \lambda_{s,N-s}(0, T_s)^2 \\ &\quad - 2w_1 w_2 \lambda_{s,N-s}(0, T_s) \lambda_{s,e-s}(0, T_s) \rho + w_2^2 \lambda_{s,e-s}(0, T_s)^2,\end{aligned}\quad (20.27)$$

where the volatilities $\lambda_{s,N-s}(0, T_s)$, $\lambda_{s,e-s}(0, T_s)$ are now observable as they correspond to ATM (spot) starting swaptions with expiry T_s on $(N-s)$ -period and $(e-s)$ -period swaps, respectively. Invoking (20.24) and putting it all together, we have

$$\begin{aligned}V_{e,N-e}(0; T_s) &\approx A_{e,N-e}(0) \sqrt{\frac{2}{\pi}} \left(\lambda_{e,N-e}(0, T_e)^2 T_e - (w_1^2 \lambda_{s,N-s}(0, T_s)^2 \right. \\ &\quad \left. - 2w_1 w_2 \lambda_{s,N-s}(0, T_s) \lambda_{s,e-s}(0, T_s) \rho + w_2^2 \lambda_{s,e-s}(0, T_s)^2) T_s \right)^{1/2},\end{aligned}\quad (20.28)$$

where the only unobserved parameter is the correlation ρ . This parameter can, for instance, be estimated from a properly-calibrated LM model or left as an “exotic” parameter for traders to tweak.

As the two relevant swap rates $S_{s,N-s}(T_s)$, $S_{s,e-s}(T_s)$ fix on the same date, the correlation ρ is usually quite high; in fact, it is not uncommon to simply assume that $\rho = 1$. Additionally, sometimes one approximates $A_{e,N-e} \approx T_N - T_e$, $A_{e,s-e} \approx T_s - T_e$, and so forth. This results in the following well-known approximation for forward volatility:

$$\begin{aligned}\lambda_{e,N-e}(T_s, T_e) &= \left(\lambda_{e,N-e}(0, T_e)^2 \frac{T_e}{T_e - T_s} \right. \\ &\quad \left. - \left(\frac{T_N - T_s}{T_N - T_e} \lambda_{s,N-s}(0, T_s) - \frac{T_e - T_s}{T_N - T_e} \lambda_{s,e-s}(0, T_s) \right)^2 \frac{T_s}{T_e - T_s} \right)^{1/2}.\end{aligned}\quad (20.29)$$

While (20.29) may occasionally be useful for back-of-the-envelope computations, (20.27) is still preferable.

The simple expression (20.28) for the value of a forward swaption straddle makes its vega exposure quite transparent. A position in the forward straddle is equivalent to a long position in the spot-starting straddle on the same rate $S_{e,N-e}(T_e)$, minus a spread option on two swap rates $S_{s,N-s}(T_s)$ and $S_{s,e-s}(T_s)$. It is important to realize, however, that the vega hedge suggested by this decomposition is not static: as rates move, the vega of a forward swaption straddle does not change, while the vegas of standard swaptions implicitly used in the decomposition (20.28) do change, and may disappear altogether if the swaptions become sufficiently far in or out of the money. The vega hedge consequently must be rebalanced quite frequently over the life of the forward swaption straddle, often at fairly significant expense. As for other risk sensitivities, the forward swaption straddle has no delta⁵ and

⁵In the sense that there is no sensitivity to the yield curve, provided that all basis-point (Gaussian) volatilities are kept fixed. If we assume that the basis-point

(almost) no gamma until the time of the strike fix (T_s); so it is, indeed, an instrument with pure volatility exposure.

The formula (20.28) was obtained using Gaussian approximation. In Section 20.2.4, on the other hand, we highlighted the importance of accounting for volatility smile in pricing forward swaption straddles. This, however, does not invalidate (20.28), as we recall that the main issue with the smile is its impact on the meaning of the mean reversion parameter — and, ultimately, the effective correlation in the model. In (20.28), we control de-correlation directly, through an exogenous correlation parameter ρ , which “bundles” smile effects and correlation effects in one parameter.

The task of choosing a reasonable value for ρ is fairly straightforward since, as we recall, forward swaption straddles are usually rather short dated. Still, if we want to study the smile effects separately from correlation, we can extend the model to explicitly include the smile. Many routes could be taken here; let us outline a possible approach. First, we note that the forward swaption straddle value is given by the expected value of the payoff

$$|S_{e,N-e}(T_e) - S_{e,N-e}(T_s)| \quad (20.30)$$

in the annuity measure in which $S_{e,N-e}$ is a martingale. The distribution of $S_{e,N-e}(T_e)$ is known directly in this measure from the swaption values across strikes (recall Chapter 16). The distribution of $S_{e,N-e}(T_s)$ is, however, not known. However, by (20.25) and (20.26) it can be represented as a weighted difference of two swap rates $S_{s,N-s}(T_s)$, $S_{s,e-s}(T_s)$ whose full distributions are, again, observable. So we can rewrite the payoff (20.30) as

$$|S_{e,N-e}(T_e) - w_1 S_{s,N-s}(T_s) + w_2 S_{s,e-s}(T_s)|$$

and then use any of the copula methods from Chapter 16, methods that allow direct inclusion of full marginal distributions of the rates involved. We let the reader fill in the remaining details.

Before concluding, let us mention a few variations of the basic forward swaption straddle product. We already mentioned options on forward starting swaps, which are closely related to forward swaption straddles and can be priced along the same lines as above. Another related contract pays the value of the implied basis-point volatility (for an at-the-money straddle with expiry T_e) at T_s , i.e. a contract with the payoff

$$\sigma_N(T_s, S_{e,N-e}(T_s); T_e, S_{e,N-e}(T_s))$$

at T_s . Its value can be linked to that of a forward swaption straddle as we have

volatility surface moves with the yield curve (see discussion in Section 16.1.1 on backbones), then a “shadow delta” could, of course, come into play.

$$\begin{aligned} & E(\beta(T_s)^{-1}\sigma_N(T_s, S_{e,N-e}(T_s); T_e, S_{e,N-e}(T_s))) \\ & \quad = E^A \left(\frac{\sigma_N(T_s, S_{e,N-e}(T_s); T_e, S_{e,N-e}(T_s))}{A_{e,N-e}(T_s)} \right). \end{aligned}$$

Apart from some deterministic scaling, the difference from (20.22) is in the convexity term $1/A_{e,N-e}(T_s)$. We can link it to the value of the swap rate $S_{e,N-e}(T_s)$ using methods from Chapter 16. While we are not able to derive a simple formula such as (20.28), copula methods obviously still apply.

Out-of-Model Adjustments

When valuing exotic derivatives, like CLEs and TARNs, that are written on an underlying structured swap, it is natural to desire that the structured swap is priced in line with the market¹. Sometimes such consistency is easy to achieve, as when the term structure model used for exotic derivatives pricing happens to coincide (either exactly or to good approximation) with the vanilla model(s) used to define the “market”. For instance, the stochastic volatility versions of the quasi-Gaussian model (see Chapter 13) and the Libor market model (see Chapters 14 and 15) are consistent with the vanilla SV model of Chapter 8 when it comes to European swaption pricing. Such consistency is, however, not always feasible. For instance, when using volatility smile parameterizations such as SABR or SVI for vanilla swaption marking (see Sections 8.6 and 16.1.5), we will typically not be perfectly consistent with any of the standard term structure models above. Similarly, it is hard to imagine a term structure model that would be exactly consistent with some of the copulas we introduced for multi-rate vanilla derivatives pricing in Chapter 17. Even if a sophisticated calibration routine is employed, such lack of consistency generally implies that the vanilla and term structure models will disagree in an economically meaningful way on the prices of structured swaps underlying exotic derivatives, a situation that is typically seen as undesirable.

In this chapter we review various methods to *force* the value of a structured swap in a term structure model to match the market (or vanilla model) value, through outright manipulation of some quantity that affects the derivative price. In selecting the quantity (or quantities) to alter, any of the key “ingredients” in a derivative model price are potential candidates: the model, the market data, and the trade. We consider all three possibilities in what follows. There is a strong ad-hoc flavor to all the methods we present,

¹For exotic swaps that require a model for their valuation, the “market” is often understood in a broad sense to represent the values of underlying exotic swaps in a vanilla model of choice.

and theoretical justification is typically rather weak. Nevertheless, if applied judiciously, risk management and pricing accuracy can sometimes benefit from the methods of this chapter. To be clear, the methods we present are not designed to cover for gross mis-calibrations or mis-valuations of underlying swaps, and should not be used as such. We only (cautiously) endorse them as ways of correcting for “small” mismatches in valuation. While it is difficult to make general statements on how small is small, one should use a combination of common sense, experience and rigorous testing in making the judgment.

21.1 Adjusting the Model

We start out by considering adjustments, where model-derived information is used to adjust the value of an exotic derivative to account for mispricing of the underlying. Here, and throughout the chapter, we denote the coupons of the structured swap used as an underlying for a given derivative by C_n , $n = 1, \dots, N - 1$, and the exotic derivative itself — a callable Libor exotic or a TARN, for example — by H_0 . Note that we therefore depart slightly from our standard notation, whereas we would normally use C_n and H_0 to determine *values* of coupons and exotic derivatives; we do so to distinguish different values calculated by different methods. In particular, the market value of the n -th coupon is denoted by $V_{\text{mkt}}(C_n)$, $n = 1, \dots, N - 1$. The value of the same coupon in the term structure model is denoted by $V_{\text{mdl}}(C_n)$, $n = 1, \dots, N - 1$.

21.1.1 Calibration to Coupons

As we have seen on numerous occasions in this book, when pricing exotic derivatives, term structure models are typically calibrated to a multitude of European swaptions. Of course, this in itself does not necessarily guarantee that the prices of coupons of the underlying swap in the term structure model would match their market, or vanilla model, prices. A good example here is a fixed-rate callable range accrual (see Section 5.13.4), where the n -th coupon is given by

$$C_n = k \times \frac{1}{T_{n+1} - T_n} \sum_{t \in [T_n, T_{n+1}]} 1_{\{L(t) \in [l, u]\}},$$

with $L(t)$ being some Libor rate. A fixed-rate range accrual coupon is decomposable into a collection of digital options on the Libor rate, and as such can be valued in a vanilla model with slight timing-delay convexity adjustments, see Chapter 16 and in particular Section 16.5. It is likely, however, that $V_{\text{mdl}}(C_n) \neq V_{\text{mkt}}(C_n)$, due to, for example, differences in the

treatment of convexity effects or in the volatility smiles implied by the term structure and the vanilla models.

To guarantee that $V_{\text{mdl}}(C_n) = V_{\text{mkt}}(C_n)$ for all $n = 1, \dots, N - 1$, the model can explicitly be calibrated to the market prices of the underlying coupons. This *extended calibration* method is fairly benign as far as model adjustments go, and could well be considered an extension of the local projection method (see Section 18.4). The ability to calibrate to the prices of underlying coupons relies, of course, on the availability of efficient methods for calculating their values $V_{\text{mdl}}(C_n)$, $n = 1, \dots, N - 1$, in a term structure model used. For most “interesting” coupons, however, closed-form expression are generally unavailable in sophisticated term structure models (e.g. LM-type models), and one has to resort to numerical calculations for calibration, often requiring Monte Carlo simulations. Calibration by Monte Carlo simulation is not something we would typically recommend, but if this approach is chosen nevertheless, some fairly obvious precautions should be taken: all coupons should be computed in the same simulation loop, the simulation seed must be taken to be the same in all calibration iterations as well as for the main valuation, and so on. A body of literature on these so-called *stochastic optimization methods* exists (see e.g. Broadie et al. [2009], Andradóttir [1995], Andradóttir [1996]) and should be consulted before one attempts to use Monte Carlo simulation inside a calibration loop.

For term structure models amendable to PDE methods, calculating coupon values by numerical (PDE) methods for calibration is certainly a plausible strategy. For numerical efficiency, we recommend usage of the forward induction method (see Section 11.3.2.3), rather than the standard backward induction method, when calculating coupon values in a calibration algorithm.

Apart from numerical issues, the extended calibration method has certain other caveats. Calibration to non-standard targets requires special care, as one has to be mindful of using the right parameters in the calibration. For example, the value of the range accrual coupon, being a sum of digital options, is not a monotonic function of volatility, and trying to calibrate volatility of the model to the market prices of range accrual coupons may yield unrealistic volatility levels or fail outright. In this particular case, it is clear that the *skew* of the volatility smile is a primary driver of the range coupon value and, hence, it is the model skews (in a skew-enabled model such as a quasi-Gaussian local volatility model), and not the volatilities, that should be calibrated to the market prices of range accrual coupons.

Even if calibration to coupons does not fail, it may result in a set of model parameters that are inappropriate for valuing other optionality embedded in a given exotic derivative, such as callability. In the callable range accrual example above, had we mistakenly tried to calibrate the volatility of a term structure model to range accrual coupons, the model volatilities could end up being very high or very low, significantly over- or under-estimating the value of callability feature. Similarly, if we were to try to value a CMS spread

TARN (see Section 5.15) in a one-factor model by, say, calibrating mean reversion to the underlying CMS spread option values², the resulting mean reversion would very likely be inappropriate for valuing the trigger feature of the TARN.

21.1.2 Adjusters

While the extended calibration method can be attractive, brute-force calibration to coupon values may not always be feasible for numerical or other reasons. Fortunately, the problem can be simplified if we recall our main tenet that prudent application of out-of-model adjustments should be limited to correcting for *small* mismatches, in which case we should be able to linearize the problem and solve it with less effort. We call this idea the *adjusters method* after Hagan [2002] who popularized it.

Let ξ be some model parameter — a volatility function, a correlation parameter, a vector of mean reversions or even a yield curve — and ξ_0 be its calibrated value. In general ξ can be represented as a column vector, and to simplify our exposition we assume it is $(N - 1)$ -dimensional, with the n -th coordinate affecting the value of the n -th coupon. We make the dependence of model prices of various securities on ξ explicit and write $V_{\text{mdl}}(C_n; \xi)$.

Let ξ^* be the solution of

$$V_{\text{mdl}}(C_n; \xi^*) = V_{\text{mkt}}(C_n) \text{ for all } n = 1, \dots, N - 1.$$

For all $n = 1, \dots, N - 1$, ξ^* satisfies, to first order,

$$V_{\text{mdl}}(C_n; \xi_0) + \frac{\partial V_{\text{mdl}}}{\partial \xi}(C_n; \xi_0)(\xi^* - \xi_0) \approx V_{\text{mkt}}(C_n), \quad (21.1)$$

where $\partial V_{\text{mdl}}/\partial \xi$ is a row vector of $\partial V_{\text{mdl}}/\partial \xi_n$, $n = 1, \dots, N - 1$. Hence,

$$\xi^* \approx \xi_0 + \left[\frac{\partial \bar{V}_{\text{mdl}}}{\partial \xi}(C; \xi_0) \right]^{-1} (\bar{V}_{\text{mkt}}(C) - \bar{V}_{\text{mdl}}(C; \xi_0)), \quad (21.2)$$

where we use bars to denote column vectors,

$$\bar{V}_{\text{mdl}}(C; \xi) \triangleq (V_{\text{mdl}}(C_1; \xi), \dots, V_{\text{mdl}}(C_{N-1}; \xi))^T,$$

and so forth. In particular,

$$\frac{\partial \bar{V}_{\text{mdl}}}{\partial \xi}(C; \xi) \quad (21.3)$$

is an $(N - 1) \times (N - 1)$ matrix whose n -th row is $\partial V_{\text{mdl}}(C_n; \xi)/\partial \xi$.

Before proceeding, let us note that while we for exposition purposes assumed the dimension of the model parameter ξ to be the same as the

²Needless to say, this is not something that we generally recommend.

number of coupons, this need not be so in actual applications. In particular, if the dimensions do not match, we can think of the equation (21.1) as a linear regression problem and find ξ^* by the appropriate least-squares methods. This procedure has been used many times already in this book, see, for example, Section 6.4.3.

Having identified ξ^* , the adjusted model price of H_0 is given by

$$V_{\text{adj}}(H_0) = V_{\text{mdl}}(H_0; \xi^*),$$

and, expanding to first order and substituting (21.2),

$$\begin{aligned} V_{\text{adj}}(H_0) &\approx V_{\text{mdl}}(H_0; \xi_0) + \frac{\partial V_{\text{mdl}}}{\partial \xi}(H_0)(\xi^* - \xi_0) \\ &= V_{\text{mdl}}(H_0; \xi_0) \\ &\quad + \frac{\partial V_{\text{mdl}}}{\partial \xi}(H_0) \left[\frac{\partial \bar{V}_{\text{mdl}}}{\partial \xi}(C; \xi_0) \right]^{-1} (\bar{V}_{\text{mkt}}(C) - \bar{V}_{\text{mdl}}(C; \xi_0)). \end{aligned} \tag{21.4}$$

With these formulas, the adjusters method follows these steps.

1. Given the calibrated model parameter value ξ_0 , the unadjusted values of the exotic H_0 , and of all the underlying coupons C_n , $n = 1, \dots, N - 1$, are computed.
2. The sensitivities of the value of H_0 and values of the various C_n to the model parameter ξ (at ξ_0) are computed.
3. The matrix of parameter sensitivities in (21.3) is inverted.
4. The adjusted exotic value $V_{\text{adj}}(H_0)$ is calculated via the linear approximation (21.4).

Note that the calibration loop of Section 21.1.1 is now replaced by the calculation of the sensitivity matrix $\partial \bar{V}_{\text{mdl}}(C)/\partial \xi$. Given the actual structure of the problem, this matrix may be known to be of specific form, e.g. diagonal or lower-triangular, further simplifying its evaluation. Moreover, the matrix can often be cached and reused when calculating risk sensitivities, further improving the overall efficiency of the scheme. Other needed quantities such as $\partial V_{\text{mdl}}(H_0)/\partial \xi$, are typically calculated anyway for risk management purposes and should not, as a rule, add to the overall computational burden.

The adjusters method is not restricted to using the underlying coupons as adjusters, but can be applied more broadly. For example, the value of a Bermudan swaption in a model with no volatility smile capabilities could be “adjusted” for the smile by using European swaptions as adjusters. In a sense, we can see the adjustment as a type of a control variate method (see Section 3.4.3 or Chapter 25 below) with the values of the adjusters (coupons or other vanilla instruments) used as controls.

The adjusters method potentially applies to a variety of model parameters, and a key question concerns which parameter should be used as ξ — model

volatilities, skews, etc. The answer follows the same logic as in Section 21.1.1: we should apply the method to the parameter(s) that have the most impact on the values of coupons while affecting other features of the model as little as possible. For example, the level of a yield curve often affects coupon values directly and so we can use the yield curve as the adjuster. This case has a special name, the *delta-adjustment method*, and it is similar to some of the approaches we discuss below, such as the spread adjustment method (Section 21.2) and the strike adjustment (Section 21.3.3) methods. Volatilities, too, are often a good choice for adjustment as values of most “interesting” coupons depend on volatilities of relevant rates. Of course the situation could be more complicated as the example of a fixed range accrual coupon in Section 21.1.1 demonstrated; here the skew of the volatility smile, and not its overall level, was the most relevant adjuster. Overall, nothing replaces careful analysis of each type of exotic derivative before the adjusters method is applied.

Before we wrap up our discussion of adjusters we note that the term

$$\frac{\partial V_{\text{mdl}}}{\partial \xi} (H_0) \left[\frac{\partial \bar{V}_{\text{mdl}}}{\partial \xi} (C; \xi_0) \right]^{-1}$$

in (21.4) could be interpreted as the sensitivity of the value of an exotic to the values of the underlying coupons, an interesting measure of sensitivity in its own right.

21.1.3 Path Re-Weighting

In the case of Monte Carlo based models, an approach from Avellaneda et al. [2001] makes it possible to exactly match calibration targets to their desired values, while also correcting for numerical inaccuracies of the valuation method. Let us discuss the idea in some detail. As a start, we denote simulated Monte Carlo paths by ω_i , $i = 1, \dots, K$. As always, the value estimate of any payoff — be it a zero-coupon bond, a vanilla option, or the coupon C_n — is given by the average of the payoff values associated with each path ω_i , $i = 1, \dots, K$. Focusing exclusively on the problem of matching the model values of coupons C_n , $n = 1, \dots, N - 1$, to the market, we denote by C_n^i the value of the n -th coupon, $n = 1, \dots, N - 1$, along path ω_i , $i = 1, \dots, K$. Then, the basic Monte Carlo value estimate of the n -th coupon in the model is given by

$$V_{\text{mdl}} (C_n) = \frac{1}{K} \sum_{i=1}^K C_n^i. \quad (21.5)$$

The idea of the *path re-weighting method* is to assign non-equal probabilities to the different paths in order to match target values. Let the probability assigned to the path ω_i be p_i , satisfying the standard requirements

$$0 \leq p_i \leq 1 \quad \forall i = 1, \dots, K, \quad (21.6)$$

$$\sum_{i=1}^K p_i = 1. \quad (21.7)$$

Then the value of the n -th coupon is given by

$$\sum_{i=1}^K C_n^i p_i, \quad (21.8)$$

and is a linear function of the vector $p = (p_1, \dots, p_K)^\top$. Hence, one would expect that it should be fairly straightforward to find a vector p that matches model prices of all coupons to the market,

$$\sum_{i=1}^K C_n^i p_i = V_{\text{mkt}}(C_n), \quad n = 1, \dots, N - 1. \quad (21.9)$$

The resulting “probabilities” can subsequently be reused in the pricing of the exotic derivative.

The problem (21.6), (21.7), (21.9) is under-specified since the number of paths used — which also equals the dimension of the vector p — is typically (much) larger than the number of coupons. Hence, a suitable regularization target is needed if we want to have a unique solution. It is not unreasonable, for example, to try to keep the vector p as close as possible to the equi-weighted probabilities of (21.5). Working with probability distributions, a convenient measure of closeness is the so-called *Kullback-Leibler relative entropy* between the probability vector p and the equi-weighted prior. With this choice of norm, we can formalize the search for p as the following minimization problem:

$$I(p) \triangleq \sum_{i=1}^K p_i \ln(p_i) \rightarrow \min, \quad (21.10)$$

subject to the linear inequality constraints (21.6), as well as the linear equality constraints (21.7) and (21.9). Proponents of the principle of relative entropy optimization often justify the choice of norm in (21.10) from a perspective of information theory (e.g., as a way to ensure that we do not add information that we do not possess to the problem), but a more standard least-squares norm would likely do just as well³.

The range of model errors that the path re-weighting method can correct for is limited, since the $\sum_{i=1}^K C_n^i p_i$ will obviously always be between the minimum and the maximum path value of C_n , among the K paths. If the target $V_{\text{mkt}}(C_n)$ is outside this range, such a gross mismatch cannot

³We briefly consider least-squares norms later in the section.

be corrected. Should this situation ever arise, the difference between the model and market values would likely be of such magnitude that the path re-weighting scheme would fundamentally be inappropriate anyway, as we discussed in the beginning of this chapter.

Let us develop the solution to the entropy minimization problem above in a bit more detail. For this purpose, let $\lambda = (\lambda_1, \dots, \lambda_{N-1})$ be a vector of Lagrange multipliers for the constraints (21.9), and μ a Lagrange multiplier for the total probability constraint (21.7). Then the solution to the following unconstrained problem (the “dual” formulation of the constrained problem, see Cover and Thomas [2006]),

$$\min_{\lambda} \max_{p, \mu} J(p, \lambda, \mu), \quad (21.11)$$

where (note the negative sign in front of $I(p)$)

$$\begin{aligned} J(p, \lambda, \mu) &\triangleq -I(p) \\ &+ \sum_{n=1}^{N-1} \lambda_n \left(\sum_{i=1}^K C_n^i p_i - V_{\text{mkt}}(C_n) \right) + \mu \left(\sum_{i=1}^K p_i - 1 \right), \end{aligned} \quad (21.12)$$

if it happens to satisfy (21.6), would also solve (21.10) subject to (21.6), (21.7) and (21.9).

Proposition 21.1.1. *For a given vector λ , the solution of the inner maximization problem in (21.11) is given by*

$$\mu^* = 1 - \ln(Z(\lambda))$$

and

$$p_i^* = \frac{1}{Z(\lambda)} \exp \left(\sum_{n=1}^{N-1} \lambda_n C_n^i \right), \quad i = 1, \dots, K, \quad (21.13)$$

where the partition function $Z(\lambda)$ is given by

$$Z(\lambda) = \sum_{i=1}^K \exp \left(\sum_{n=1}^{N-1} \lambda_n C_n^i \right).$$

Proof. The necessary conditions for the inner maximum in (21.11) are given by

$$\frac{\partial J(p^*, \lambda, \mu^*)}{\partial \mu} = 0, \quad \frac{\partial J(p^*, \lambda, \mu^*)}{\partial p_i} = 0, \quad i = 1, \dots, K,$$

so that we have

$$-\ln(p_i^*) - 1 + \sum_{n=1}^{N-1} \lambda_n C_n^i + \mu^* = 0, \quad i = 1, \dots, K,$$

and $\sum_{i=1}^K p_i^* = 1$. The proposition follows. \square

We note that p_i^* 's defined by (21.13) always satisfy (21.6). The distribution of the form (21.13) is known as the *Boltzman-Gibbs distribution* for the partition function $Z(\lambda)$.

Now, substituting (21.13) into the definition of the objective function (21.12) we obtain

$$G(\lambda) \triangleq J(p^*, \lambda, \mu^*) = \ln(Z(\lambda)) - \sum_{n=1}^{N-1} \lambda_n V_{\text{mkt}}(C_n). \quad (21.14)$$

Now all we need to do is to minimize (21.14), i.e. solve the $(N-1)$ -dimensional optimization problem

$$\lambda^* = \operatorname{argmin}_{\lambda} (G(\lambda)). \quad (21.15)$$

Compared to the original formulation (21.10), the dimensionality of the problem has now been significantly reduced, as normally $N \ll K$. Moreover, (21.15) is unconstrained, and thus easier to solve by standard optimization techniques. In addition, it is a “nice” optimization problem as the function $G(\lambda)$ is globally convex with a single minimum, as stated in the following proposition.

Proposition 21.1.2. *The function $G(\lambda)$ is globally convex. In particular, the following holds for all $n, m = 1, \dots, N-1$,*

$$\frac{\partial G(\lambda)}{\partial \lambda_n} = E^\lambda(C_n) - V_{\text{mkt}}(C_n), \quad \frac{\partial^2 G(\lambda)}{\partial \lambda_n \partial \lambda_m} = \text{Cov}^\lambda(C_n, C_m), \quad (21.16)$$

where the measure Q^λ is defined on Monte Carlo paths by Boltzman-Gibbs weights p^* that correspond to λ as per (21.13), i.e. for any random variable X ,

$$E^\lambda(X) = \sum_{i=1}^K p_i^* X^i,$$

where X^i is the realization of the random variable X on the i -th path, $i = 1, \dots, K$.

Proof. Let us fix n . By straightforward differentiation,

$$\begin{aligned} \frac{\partial G(\lambda)}{\partial \lambda_n} &= \frac{1}{Z(\lambda)} \frac{\partial Z(\lambda)}{\partial \lambda_n} - V_{\text{mkt}}(C_n) \\ &= \frac{1}{Z(\lambda)} \sum_{i=1}^K \exp \left(\sum_{j=1}^{N-1} \lambda_j C_j^i \right) C_n^i - V_{\text{mkt}}(C_n) \\ &= \sum_{i=1}^K p_i^* C_n^i - V_{\text{mkt}}(C_n). \end{aligned}$$

Furthermore,

$$\frac{\partial^2 G(\lambda)}{\partial \lambda_n \partial \lambda_m} = \frac{1}{Z(\lambda)} \frac{\partial^2 Z(\lambda)}{\partial \lambda_n \partial \lambda_m} - \frac{1}{Z(\lambda)^2} \frac{\partial Z(\lambda)}{\partial \lambda_n} \frac{\partial Z(\lambda)}{\partial \lambda_m},$$

and we obtain that

$$\frac{\partial^2 G(\lambda)}{\partial \lambda_n \partial \lambda_m} = \text{Cov}^\lambda(C_n, C_m)$$

by straightforward calculations. The fact that $G(\lambda)$ is globally convex now follows from the representation of the second derivative of G as a covariance matrix in (21.16), and the fact that a covariance matrix is always nonnegative-definite. \square

We point out an interesting consequence of (21.16) is that the solution to the optimization problem (21.15) is given by such λ^* that

$$V_{\text{mkt}}(C_n) = E^{\lambda^*}(C_n), \quad n = 1, \dots, N-1.$$

As the objective function $G(\lambda)$ is globally convex and its first- and second-order derivatives are straightforward to calculate, most non-linear optimization algorithms as discussed in, for example, Section 14.5.7 would work well. For extra performance, specialized methods tuned for convex objective functions, such as the Nesterov-Nemirovskii algorithm (see Nesterov et al. [1994]), could be applied.

The constrained minimization formulation (21.10) as presented in Avelaneda et al. [2001] is not the only possible way to formalize the problem of path re-weighting. For example, we can replace the exact repricing criteria (21.9) by a suitably-defined least-squares target. In particular, denoting by v_n the penalty for violating (21.9) for a given n , $n = 1, \dots, N-1$, the problem can be re-formulated as

$$\sum_{i=1}^K p_i \ln(p_i) + \sum_{n=1}^{N-1} v_n \left(\sum_{i=1}^K C_n^i p_i - V_{\text{mkt}}(C_n) \right)^2 \rightarrow \min, \quad (21.17)$$

subject to (21.6), (21.7). Not surprisingly, the problem can also be solved by the partition function method along the lines of Proposition 21.1.1, a statement we leave to the reader to verify. Finally, an even simpler quadratic problem could be obtained by replacing relative entropy as an objective function by its second-order Taylor expansion around the equi-weighted probabilities, see e.g. Glasserman [2004], Section 4.5:

$$\sum_{i=1}^K (p_i - 1/K)^2 + \sum_{n=1}^{N-1} v_n \left(\sum_{i=1}^K C_n^i p_i - V_{\text{mkt}}(C_n) \right)^2 \rightarrow \min, \quad (21.18)$$

subject to (21.6), (21.7). Again, we leave it up to the reader to fill in relevant details.

It is worthwhile pointing out an interesting connection between entropy minimization methods and the problem of calculating risk sensitivities. It turns out that under the (rather unrealistic, admittedly) assumption that market data shocks do not affect generated Monte Carlo paths but only change the right-hand-side values in (21.9), the sensitivities of the exotic derivative to the prices of coupons/market shocks can be deduced via duality arguments from the solutions of the relevant optimization problems (21.10), (21.17) or (21.18). Details can be found in Avellaneda et al. [2001].

Most adjustment methods have undesirable side effects, and path re-weighting is no exception. With non-uniform weights assigned to paths, the prices of zero-coupon bonds in the model may no longer match their market values, with the model then allowing arbitrage. This in principle could be patched up by adding all relevant zero-coupon bonds to the set of constraints (21.9) to match, but, of course, at higher computational cost. Calibration to vanilla options could also unravel — remediation will, once again, involve enlargement of the set of constraints. Again, we remind the reader that over-using methods such as path re-weighting could be dangerous, as it is difficult to control all the consequences if large deviations from the equi-weighted paths are required. Should such situations arise, the model is most likely seriously mis-specified and any valuation results should be treated as suspect.

While introduced here as an adjustment technique, let us finally note that path re-weighting could be interpreted as a variance reduction technique, provided that the option prices we are matching our finite-sample estimates to are known to coincide with the true (infinite-sample) model values. Clearly, the resulting method would have strong similarities to the more familiar technique of control variates (see Chapter 25). Glasserman and Yu [2005] investigate this link further, and prove that the two techniques are essentially identical, for large enough sample sizes. For strict variance reduction purposes, the more straightforward method of control variates is therefore typically preferable.

21.1.4 Proxy Model Method

Suppose we have identified that a given exotic security is sensitive to an “exotic risk” factor. This factor may not be important for the valuation of vanilla securities, and implementing it into a term structure model may result in the model being so complex that analytical approximations used for calibration to the vanilla market fail to be accurate enough. On the other hand, suppose we also have a simpler term structure model that calibrates well to the vanilla market but does not have the required exotic risk factor. The following procedure, which we call the *proxy model method*, is sometimes used to combine the two models to measure the sensitivity to the exotic risk factor. First, we calculate the difference in value of the derivative in the complex model for different values of the exotic risk factor, $\xi = \xi_1$ and

$$\xi = \xi_0,$$

$$\Delta V_{\text{complex}} = V_{\text{complex}}(H_0; \xi_1) - V_{\text{complex}}(H_0; \xi_0).$$

Here typically ξ_0 corresponds to the base case, i.e. the value of the exotic risk factor that is more or less consistent with the simplified worldview of the simpler term structure model, and ξ_1 is our view of the actual market-observed value of the risk factor. Next, we calculate the base value of the derivative in the simple model calibrated to the market,

$$V_{\text{simple}}(H_0).$$

We would like to add $\Delta V_{\text{complex}}$ to $V_{\text{simple}}(H_0)$ to account for the risk factor impact; however, $\Delta V_{\text{complex}}$ is biased due to the problems of calibrating the complex model. To correct for the bias, we calibrate our simple model *to the vanilla prices generated by the complex model* in the two scenarios. Henceforth, we define $V_{\text{simple},0}(H_0)$ and $V_{\text{simple},1}(H_0)$ to be the prices of the derivative in question in the simple model calibrated to the vanilla prices as generated by the complex model with $\xi = \xi_0$ and $\xi = \xi_1$.

As the simple model is insensitive to the exotic risk factor, we would expect

$$\Delta V_{\text{simple}} = V_{\text{simple},1}(H_0) - V_{\text{simple},0}(H_0)$$

to solely represent the impact of mis-calibration of the complex model to vanillas on the value of the exotic derivative. Thus, it is not unreasonable to define the adjusted price by

$$V_{\text{adj}}(H_0) = V_{\text{simple}}(H_0) + \Delta V_{\text{complex}} - \Delta V_{\text{simple}}.$$

To make the discussion above a bit more concrete, consider the problem of assessing the impact of stochastic volatility de-correlation, which would be the exotic risk factor under consideration, on a callable CMS spread derivative. Suppose we have a suitable model, say an LM model from Section 15.7, which has multiple sources of volatility randomness. We would value the derivative with the correlation of volatilities set to 100% (ξ_0 case), and then set it to some other value that is less than 100% (ξ_1 case) which we obtain by, say, historical estimation. To correct for calibration errors induced by imperfect vanilla approximations, we would calibrate a simpler LM model with a single stochastic volatility factor to the vanilla prices produced by the complex model; the single volatility factor model serves as a model that is (one hopes) sufficiently “similar” to the complex model yet allows for accurate vanilla option approximations. In this case, in addition to the usual European swaptions, we should probably also include CMS spread options in the vanilla market, to make sure we control the extra de-correlation of rates that comes from de-correlating their volatilities.

The method outlined above is rarely accurate enough for trading and risk management purposes, but is useful for qualitative understanding of the impact of certain risk factors, as well as, say, reserve calculations.

21.1.5 Asset-Based Adjustments

Consider as an example an LM model applied to a CMS-style exotic derivative. The CMS convexity adjustment (see Section 16.6) as implied by the LM model may not be equal to the “market” CMS adjustment (as calculated by, say, a replication method from Section 16.6.1). This situation might arise in part because the volatility smiles generated by the LM model differ slightly from the ones implied by market prices. One way of compensating for the difference involves changing the trade definition, a method we discuss later in Section 21.3. Here we, instead, consider a different method in which we adjust the simulated dynamics of the relevant swap rate(s) in the LM model (see Van Steenkiste [2009]). The advantage of this *asset-based adjustment method* is that we are able to not only adjust the overall levels of the relevant swap rates, but also their volatilities and skews. This, in turn, could further aid us in closing the valuation gap for the underlying swap of a CMS-based exotic derivative.

For concreteness, consider a version of the LM model (14.4) with deterministic separable volatility (14.2.4). Suppose a given swap rate⁴ $S(t)$ is of special interest to us, because, say, some coupon of the underlying swap is a function of it,

$$C = C(S(T)). \quad (21.19)$$

The standard Monte Carlo scheme for the model involves simulating all Libor rates $\{L_n(T)\}$ per Section 14.6, calculating discount factors from simulated Libor rates and combining them to calculate the simulated value of the swap rate $S(T)$. The simulated value of the swap rate is then used to calculate the simulated value of a coupon in (21.19).

Instead of calculating the swap rate from simulated Libor rates, we can of course also just simulate $S(t)$ alongside the Libor rates directly. The exact dynamics of the swap rate under the measure used for simulating Libor rates (e.g., the spot measure) can be derived by Ito’s lemma, or from Proposition 14.4.2 by the appropriate measure change. In particular we have

$$dS(t) = \mu_S(t, \mathbf{L}(t)) dt + \varphi(S(t)) \sum_{n=1}^{N-1} w_n(t) \lambda_n(t)^\top dW^B(t), \quad (21.20)$$

where $\mathbf{L}(t)$ is the vector of all Libor rates at time t , $\mu_S(t, \mathbf{L}(t))$ is the appropriate drift, $w_n(t)$ ’s are given by (14.31), and $W^B(t)$ is a Brownian motion in the spot measure. We could discretize the SDE (21.20) in the same way we would discretize SDEs for the Libor rates, and simulate $S(t)$ together with the Libor rates; then we can use this simulated value in the payoff (21.19). Up to the discretization bias and simulation error, the value

⁴While we only consider one coupon and one swap rate, it is trivial to extend our discussion to the standard case of multiple coupons depending on different rates.

of the derivative computed in this scheme would be the same as in the standard scheme where $S(t)$ is calculated directly from Libor rates.

While the utility of the simulation scheme above by itself is questionable, it gives us a starting point for adjusting the dynamics of $S(t)$ as we see fit. Decoupling the dynamics of $S(t)$ in (21.20) from the dynamics of Libor rates allows us to treat the swap rate as a stand-alone market variable, or “asset” (hence the name of the method), whose model dynamics we can control independently. For example, suppose we would like to shift the mean of the simulated variable $S(T)$ to compensate for the differential of CMS convexity adjustments between the LM model and the market. Then we just use (21.20) with an initial condition shifted by some $c \neq 0$,

$$dS_{\text{adj}}(t) = O(dt) + \varphi(S_{\text{adj}}(t)) \sum_{n=1}^{N-1} w_n(t) \lambda_n(t)^T dW^B(t), \quad S_{\text{adj}}(0) = S(0) + c.$$

When valuing the S -dependent coupon (and the payoff of the entire exotic derivative) we then would use $S_{\text{adj}}(t)$ instead of $S(t)$ in payoff calculations.

The volatility of the swap rate could be adjusted in a similar way; for example we can specify that

$$dS_{\text{adj}}(t) = O(dt) + c\varphi(S_{\text{adj}}(t)) \sum_{n=1}^{N-1} w_n(t) \lambda_n(t)^T dW^B(t), \quad S_{\text{adj}}(0) = S(0),$$

for some volatility adjustment $c > 0$. Or, indeed, we can change the model skew of $S(t)$ by replacing (21.20) with

$$dS_{\text{adj}}(t) = O(dt) + \varphi(aS_{\text{adj}}(t) + b) \sum_{n=1}^{N-1} w_n(t) \lambda_n(t)^T dW^B(t), \quad S_{\text{adj}}(0) = S(0),$$

for some a, b . All three types of dynamics adjustments could, of course, be combined to provide a finer level of control over the distribution of $S(t)$; indeed, possibilities are limitless. The method is not restricted to adjustments of dynamics of swap rates only; we can apply the same trick to the *spread* of two swap rates to ensure that, say, the volatility of a CMS spread in this “adjusted” LM model matches that in the vanilla (multi-rate) model used.

While the asset-based adjustment method is rather flexible, it should be obvious that it comes at a serious cost of introducing arbitrage and making the model internally inconsistent. The swap rate simulated from one of the adjusted SDEs above will no longer equal the “true” swap rate as synthesized from (simulated) Libor rates. As a consequence, a European swaption would have different values in such a model depending on how the payoff is written, see equations (5.10) or (5.11) in Chapter 5. Taking this example to the extreme, one can imagine a trader equipped with such a model selling and buying identical swaptions booked in different formats and generating riskless “profits” on each trade.

21.1.6 Mapping Function Adjustments

Adjusting the dynamics of the swap rate $S(t)$ is not the only way to achieve desired changes to its distribution. As a possible alternative, we can modify its *terminal* distribution directly. In particular, instead of using the swap rate $S(T)$ when calculating the coupon value in (21.19), we would use $S_{\text{adj}}(T)$ defined by

$$S_{\text{adj}}(T) = \Lambda(S(T)).$$

Here $S(T)$ is the model-simulated value of the swap rate, and the *mapping function* $\Lambda(s)$ is chosen in such a way that $S_{\text{adj}}(T)$ has a desired distribution (e.g. one consistent with the swaption market, or perhaps with a particular vanilla model). This approach is sometimes called the *mapping function adjustment*.

When using the mapping function adjustment method, we would not adjust just one swap rate, but all rates used in calculating underlying coupons. If the product requires observations of the swap rate on different dates, we would, of course, use different mapping functions for different observation dates.

Determining the mapping function $\Lambda(s)$ for each required observation of each swap rate is conceptually simple. If $\Psi_{\text{mkt}}(s)$ is the market-implied cumulative distribution function of the swap rate $S(T)$ in the appropriate annuity measure (see Section 16.6.9), and $\Psi_{\text{mdl}}(s)$ is the same for the term structure model that we are adjusting, then we simply set

$$\Lambda(s) = \Psi_{\text{mkt}}^{-1}(\Psi_{\text{mdl}}(s)).$$

For most models we consider, efficient swaption pricing formulas exist, and the model CDF $\Psi_{\text{mdl}}(s)$ needed here is readily available. Needless to say, the mapping function(s) should be pre-computed and cached before valuing a given trade.

A careful reader will no doubt notice that this type of adjustment has a Markov-functional model flavor (see Appendix 11.A to Chapter 11). It is, however, *not* a full-blown Markov-functional model as, of course, we have made no provisions to retain the arbitrage-free characteristics of the adjusted model. Clearly, the usual caveats to usage of such non-arbitrage-free models apply, and the warnings at the end of Section 21.1.5 should be carefully considered before the mapping approach is used.

21.2 Adjusting the Market

Having finished with adjustments based on models, let us turn to using market data for that purpose. While several sources of market data could be used for adjustment, the most common target is the yield curve, where we can capitalize on the fact that the yield curve is frequently the only

parameter that is shared between the term structure model and the vanilla models. The resulting adjustment method is in many ways similar to that in Section 21.1.2, but there are a few twists that warrant a separate discussion.

Continuing with the notations of Section 21.1.2, we first specialize ξ to be the yield curve as used during valuation, and ξ_0 to be the yield curve as fit to the market prices of swaps, etc. (see Chapter 6). We formulate the adjustment problem as finding ξ^* such that

$$V_{\text{mdl}}(C_n; \xi^*) = V_{\text{mkt}}(C_n) \text{ for all } n = 1, \dots, N - 1. \quad (21.21)$$

With the view that the impact of the yield curve on the value of a coupon is roughly the same in the two models (vanilla and exotic), we define δ to be the (time-dependent) spread, to be applied to the yield curve, such that

$$V_{\text{vanilla}}(C_n; \xi_0 + \delta) = V_{\text{mdl}}(C_n, \xi_0) \text{ for all } n = 1, \dots, N - 1, \quad (21.22)$$

where $V_{\text{vanilla}}(C_n; \xi)$ is the value of coupon C_n in the vanilla model when using yield curve ξ . Then, by approximate linearity and (21.22),

$$V_{\text{mdl}}(C_n, \xi_0 - \delta) \approx V_{\text{vanilla}}(C_n; \xi_0 + \delta - \delta) = V_{\text{mkt}}(C_n),$$

and the approximate solution to (21.21) is given by

$$\xi^* \approx \xi_0 - \delta.$$

Note that solving (21.22) is normally much quicker than solving (21.21) directly.

The valuation of the exotic derivative proceeds with the adjusted yield curve $\xi^* = \xi_0 - \delta$,

$$V_{\text{adj}}(H_0) = V_{\text{mdl}}(H_0, \xi_0 - \delta).$$

We call this the *spread adjustment method*. Notably, the adjusted yield curve should only be applied to the structured leg — the Libor leg, if present, should use the original, unadjusted yield curve. Simultaneous modeling of two yield curves — the adjusted and the original — could follow the deterministic spread approach from Section 15.5.

While the idea behind the method is simple, the need to use multiple yield curves in valuation makes the method somewhat unwieldy, and it is not particularly popular. For derivatives that involve spread options, for obvious reasons we should adjust the slope, rather than the overall level, of the yield curve.

21.3 Adjusting the Trade

Adjusting the trade is probably the most common type of out-of-model adjustments. In this approach, some features of the coupons are changed

to line up the values of the adjusted coupons in the term structure model with the values of the original coupons in the vanilla model, i.e. their market values. The adjusted value of the exotic derivative is then calculated by applying the term structure model to a redefined contract with adjusted coupons.

Before discussing a few common approaches, it is worth pointing out the obvious, but sometimes overlooked, point that trade adjustments (and indeed any other type of adjustments) should be performed for each valuation of the trade — and, in particular, for each re-valuation during risk calculations. With trade adjustments in particular, it is tempting to calculate the adjustments once at trade initiation, and then book an adjusted trade in the booking system. However, even if booked trades are re-adjusted periodically, calculated risk measures would be consistently wrong, as they would *not* include the impact of market parameter shocks on the coupon adjustments.

21.3.1 Fee Adjustments

The additive fee adjustment owns its ease of applicability to the additive property of the pricing operator. Suppose we have in mind a payoff A_n to use in adjusting the coupon C_n , with A_n being paid at the payment date of the coupon, often T_{n+1} . Then we can always find a scaling α_n such that

$$\alpha_n V_{\text{mdl}}(A_n) = V_{\text{mkt}}(C_n) - V_{\text{mdl}}(C_n).$$

We define the adjusted coupon C_n^* as C_n plus the adjustment $\alpha_n A_n$. Then, clearly,

$$\begin{aligned} V_{\text{mdl}}(C_n^*) &= V_{\text{mdl}}(C_n + \alpha_n A_n) = V_{\text{mdl}}(C_n) + \alpha_n V_{\text{mdl}}(A_n) \\ &= V_{\text{mdl}}(C_n) + V_{\text{mkt}}(C_n) - V_{\text{mdl}}(C_n) = V_{\text{mkt}}(C_n). \end{aligned}$$

The adjusted value of H_0 is then given by the value, in the model, of H_0^* , where H_0^* is constructed from the adjusted coupons,

$$V_{\text{adj}}(H_0) = V_{\text{mdl}}(H_0^*).$$

This procedure is called the *fee adjustment method* because $\alpha_n A_n$ could be thought of as an extra “fee” that applies to the coupon C_n .

Fee adjustments require calculating $V_{\text{mdl}}(C_n)$ and $V_{\text{mdl}}(A_n)$, but only once as no iterative search is required. As a consequence, the method is computationally quite efficient even for Monte Carlo based models. Of course, for PDE-based models, the forward PDE valuation of coupons should still be favored over the backward PDE.

The simplest form of fee adjustment is the constant fee adjustment, i.e. using $A_n = 1$, paid at the payment date of the coupon C_n . Then the equation simplifies to be an equation on a scalar f_n such that

$$f_n V_{\text{mdl}}(1) = V_{\text{mkt}}(C_n) - V_{\text{mdl}}(C_n). \quad (21.23)$$

The adjusted coupon is given by $C_n^* = C_n + f_n$. This specialization is slightly faster than the general case as only $V_{\text{mdl}}(C_n)$ needs to be computed.

Another reasonable choice — at least from a computational complexity standpoint — of the adjustment payoff is the coupon itself, $A_n = C_n$. With this choice one would look for α_n such that

$$\alpha_n V_{\text{mdl}}(C_n) = V_{\text{mkt}}(C_n) - V_{\text{mdl}}(C_n), \quad (21.24)$$

which requires no more effort than finding a constant f_n in (21.23). The adjusted coupon is then given by

$$C_n^* = (1 + \alpha_n) C_n,$$

i.e. the adjustment has the same shape as the original coupon. Sometimes this is called a *multiplicative adjustment*.

The additive and multiplicative fee adjustment methods could be blended. Choosing ω_n , $0 \leq \omega_n \leq 1$, one can define the adjusted coupon by

$$C_n^* = C_n + (\omega_n f_n + (1 - \omega_n) \alpha_n C_n), \quad (21.25)$$

where f_n , α_n are given by (21.23), (21.24).

21.3.2 Fee Adjustment Impact on Exotic Derivatives

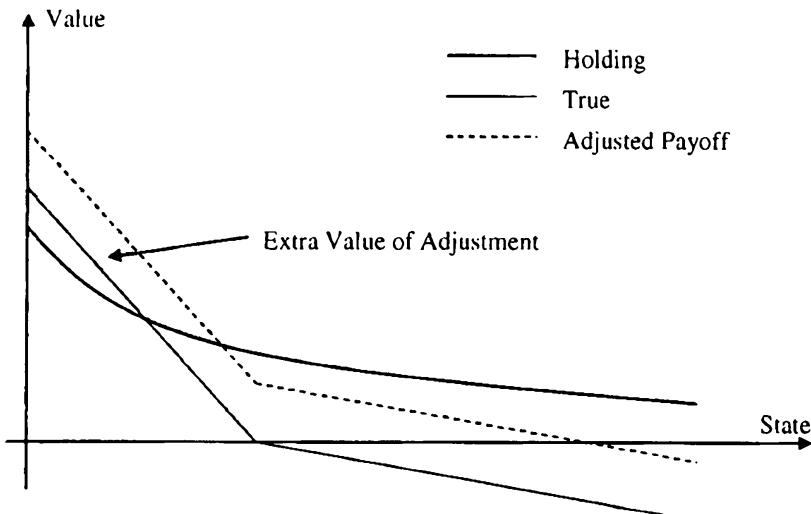
For different fee adjustment methods, the value of the structured swap underlying a given exotic derivative is invariant, by definition. This is not the case for the exotic derivative itself, as different adjustment methods would assign it different values. Generally, such differences originate with an asymmetric impact of a fee on the price of the exotic derivative. Considering a case of a callable derivative, only the changes to the underlying coupons *in the exercise region* will contribute to the price. Conversely, changes to the underlying coupons in the hold region are irrelevant. However, coupon adjustments are calculated to match the integral of the payoff *over the whole of the state space* to the market value. While the integrals of the adjusted exercise value over the whole state space are therefore independent of the type of payoff adjustment type, the same cannot be said for the integrals over the exercise region. To demonstrate, consider a callable inverse floater, i.e. a Bermudan style option to enter a swap to receive an inverse floating coupon $\max(s - gL_n(T_n), 0)$ and pay Libor $L_n(T_n)$. The underlying swap is a sum of net coupons

$$\max(s - (g + 1)L_n(T_n), -L_n(T_n)).$$

The exercise value is represented by the solid line in Figures 21.1 and 21.2. The dotted line represents the adjusted exercise value: an additive

adjustment in Figure 21.1 and a multiplicative one in Figure 21.2. While the integral of the difference of the dotted and solid lines is the same in both figures, their integrals over the exercise region, as represented by the grey area, are different. In this case, a multiplicative adjustment will assign a higher value to the callable inverse floater.

Fig. 21.1. Additive Adjustment for CLE



Notes: Effect of the additive fee adjustment on the exercise value of a callable inverse floater. “Holding” denotes the hold value of the callable inverse floater, as a function of the state of the model; “True” denotes the actual payoff of the coupon; and “Adjusted Payoff” represents the payoff adjusted according to the method described.

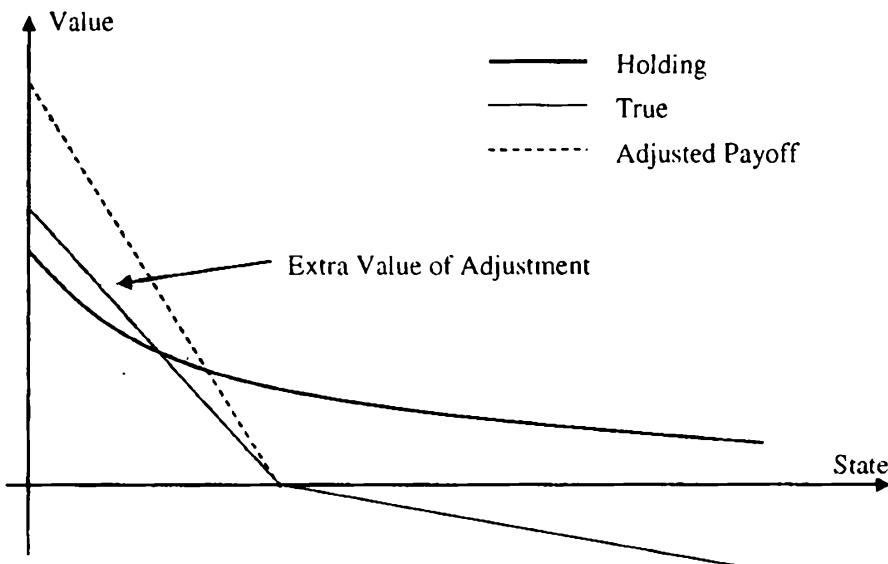
21.3.3 Strike Adjustment

Many coupon types have a natural “strike” parameter, as a quick recall of the definitions of capped/floored floaters, inverse floaters, etc. in Section 5.13 should confirm. Moreover, the value of a coupon is often a monotonic function of the strike. Denoting the strike by k , the n -th coupon as a function of strike by $C_n(k)$, and the actual value of the strike for the n -th coupon by k_n , we can therefore usually find k_n^* such that

$$V_{\text{mdl}}(C_n(k_n^*)) = V_{\text{mkt}}(C_n(k_n)) \quad (21.26)$$

for any $n = 1, \dots, N - 1$. As indicated by the notations, both k_n and k_n^* are coupon-specific, and depend on $n = 1, \dots, N - 1$. Then, denoting by H_0^* the exotic with the coupon strikes set to k_n^* , $n = 1, \dots, N - 1$, the adjusted value of the exotic is given by the model price of H_0^* ,

Fig. 21.2. Multiplicative Adjustment for CLE



Notes: Effect of the multiplicative fee adjustment on the exercise value of a callable inverse floater. See caption to Figure 21.1 for notations.

$$V_{\text{adj}}(H_0) = V_{\text{mdl}}(H_0^*).$$

This procedure is called the *strike adjustment method*.

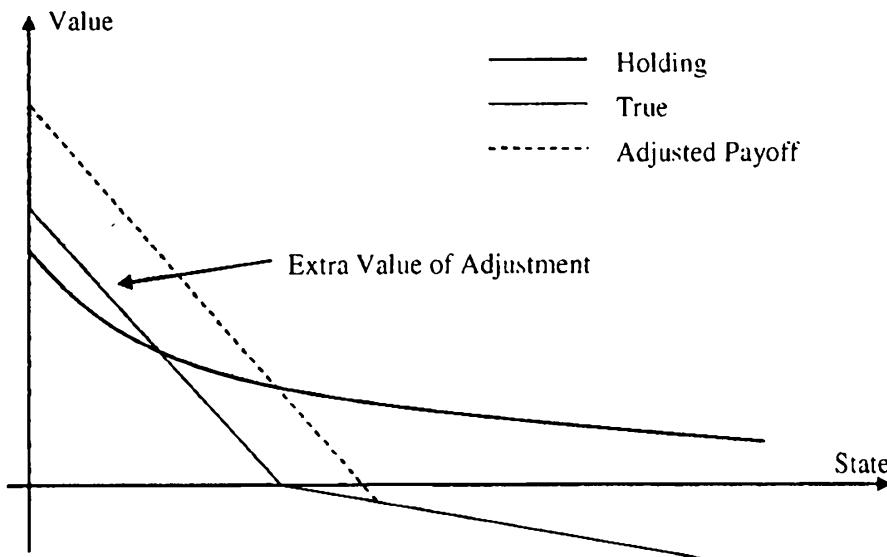
Other parameters can play the role of the strike in the method. For example, for range accrual coupons, an upper or lower range can be used, as the value of the coupon is monotone in those parameters.

Strike adjustments are more numerically intensive than fee adjustments, since solving (21.26) typically requires multiple calculations of $V_{\text{mdl}}(C_n(s))$ for different values of s . Despite higher computational costs and no discernible theoretical advantage over the fee adjustment method, the strike adjustment method remains popular, perhaps because traders are used to adjusting strikes/barriers for other purposes, such as improvement of risk management of barrier options and adjusting for sampling frequency effects (see e.g. Section 2.5.3, Theorem 3.2.2 and Broadie et al. [1997]).

As with the fee adjustment method, the effect of the strike adjustment on the value of an exotic derivative could be understood by looking at the impact of the adjustments in the relevant part of the state space. Continuing the example from the previous section, Figure 21.3 shows the impact of the strike adjustment on the price of a callable inverse floater.

As a final comment to this chapter, we note that there are undoubtedly many additional ingenious ways of adjusting models, market data and trades that could have been included in this chapter. At the end of the day, however, nothing replaces a good calibration of a well-specified term structure model to the vanilla market. Out-of-model adjustments are useful when applied sparingly, but can easily be abused. For example, it has been rumored that a

Fig. 21.3. Strike Adjustment for CLE



Notes: Effect of the strike adjustment on the exercise value of a callable inverse floater. See caption to Figure 21.1 for notations.

French bank used to risk manage its portfolio of callable CMS spread trades in a one-factor Gaussian model with trade adjustments. Needless to say that this is not something we would recommend. Even in more reasonable applications, the choice of the “right” adjustment could be a delicate exercise and continues to be more art than science.

Introduction to Risk Management

We have reached the point in the book where we are ready to discuss the problem of managing the market risk¹ exposure of interest rate derivatives portfolios. For our purposes here, the topic of primary interest is the quantification and computation of the risk exposure, a task that turns out to be quite challenging and shall require several chapters to cover. First, however, we devote a brief introductory chapter to a high-level overview of the risk management exercise, as practiced by a typical fixed income derivatives trading desk. As part of our analysis, we identify the most common “greeks” (risk sensitivities), and also provide some background on the role these play in hedging and risk management. As we shall see, actual hedging practices tend to deviate considerably from the theoretical ideals of pure delta hedging of Brownian increments (see Section 1.7). We discuss these issues here, and also provide some material on how the risk management and middle office teams in a bank may use market risk exposure information to compute summary statistics for overall risk exposure (the so-called *value-at-risk*), and to perform day-to-day analysis and breakout of portfolio profits and losses. The chapter serves to provide justification for the emphasis on greeks computation in the remainder of the book, and also elaborates on a number of discussions that have cropped up earlier, including Bermudan swaption risk management (Section 19.7.1) and computation of par-point yield curve risk reports (Section 6.4).

¹In addition, derivatives portfolios are exposed to *credit risk*, i.e. the risk that the counterparty to a derivatives transaction will default on its obligations. Management of credit risk is outside the scope of this book.

22.1 Risk Management and Sensitivity Computations

22.1.1 Basic Information Flow

To understand how a trading desk uses a model in practice, it is useful to introduce a bit of notation. Let $\Theta_{\text{mkt}}(t)$ be an N_{mkt} -dimensional vector representing the observable market data at time t . For a fixed income desk, the components of $\Theta_{\text{mkt}}(t)$ are typically swap and futures rates (for yield curve construction) and cap and swaption prices or implied volatilities at multiple strikes, tenors, and maturities (for volatility calibration). Second, let $\Theta_{\text{prm}}(t)$ denote the set of N_{prm} additional parameters that are not directly observed, but are estimated from historical data or are treated as “exotic” constants to be specified directly by the trader. Examples of such parameters include short rate mean reversion parameters, correlation parameterizations, stochastic volatility mean reversion speeds, local volatility parameters (e.g., the CEV power), and so forth. As we have seen in many chapters of this book, the question of how to split $\Theta_{\text{mkt}}(t)$ and $\Theta_{\text{prm}}(t)$ is often not clear-cut, as one can always attempt to add additional market variables to $\Theta_{\text{mkt}}(t)$ to allow us to deduce some of the elements of $\Theta_{\text{prm}}(t)$ by direct calibration, in which case these parameters can obviously be removed from $\Theta_{\text{prm}}(t)$. As we discussed in Section 14.5.9, one might for instance attempt to eliminate correlation information from $\Theta_{\text{prm}}(t)$ by introducing spread option price information into $\Theta_{\text{mkt}}(t)$; or one might try to calibrate short rate mean reversion from multiple swaption strips, rather than specify this parameter directly (see for instance Sections 13.1.8.2 and 13.1.8.3). As certain parameters are inherently difficult to extract in a stable and robust manner from market data, in practice it is rarely the case that $\Theta_{\text{prm}}(t)$ is completely empty.

Given $\Theta_{\text{mkt}}(t)$ and $\Theta_{\text{prm}}(t)$, the first step in pricing a derivative security typically involves a calibration procedure, where the vectors $\Theta_{\text{mkt}}(t)$ and $\Theta_{\text{prm}}(t)$ are turned into a vector² of model-appropriate parameters $\Theta_{\text{mdl}}(t)$ that contain the discount curve as well as the parameters that control its volatility structure and future dynamics. The calibration may itself require specification of certain control parameters, such as the smoothing weights used in a typical LM model calibration (Section 14.5); for simplicity, we consider these parameters part of $\Theta_{\text{prm}}(t)$. The model calibration itself will typically involve at least two steps: the construction of the discount bond curve, followed by calibration of a model for the dynamics of this curve. As the first step can typically be separated completely from the latter, it is informative to break the calibration in two parts, as in Figure 22.1 below. Notice that we here have introduced a pre-processing step where

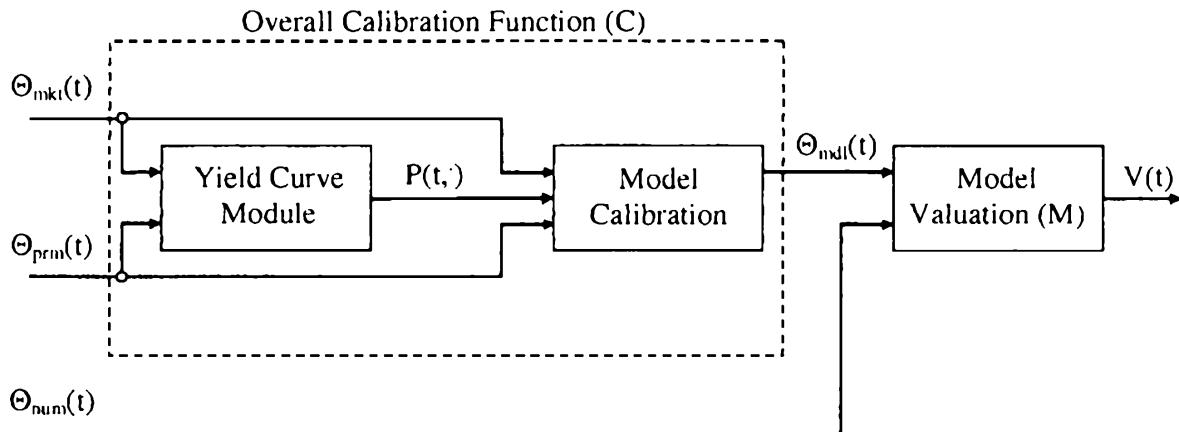
²We use the term vector loosely, since some elements of Θ_{prm} (e.g. the discount curve) may be continuous functions.

those elements of $\Theta_{\text{mkt}}(t)$ and $\Theta_{\text{prm}}(t)$ that are relevant³ for yield curve construction are used to produce a discount bond curve $P(t, T)$, $T \geq t$. Together with the (remaining) elements of $\Theta_{\text{mkt}}(t)$ and $\Theta_{\text{prm}}(t)$, this yield curve is fed to the main model calibration function, which produces $\Theta_{\text{mdl}}(t)$. In any case, we may write

$$\Theta_{\text{mdl}}(t) = C(\Theta_{\text{mkt}}(t); \Theta_{\text{prm}}(t)), \quad (22.1)$$

where C represents the (overall) calibration function.

Fig. 22.1. Information Flow



Notes: Basic information flow in derivatives pricing and model calibration.

Given the time t yield curve and a set of model parameters, we can proceed to use the model to price a given portfolio of derivative contracts. This will require us to load contract data for a specified set of securities, and also to read in additional parameters $\Theta_{\text{num}}(t)$ that control the numerical schemes used in the model. Examples of parameters in $\Theta_{\text{num}}(t)$ would include the number of Monte Carlo paths, the size of discretization steps for finite difference grids and for SDE discretization schemes, and so forth. With $V(t) = V_1(t) + \dots + V_n(t)$ denoting the value of a portfolio of n derivatives, we write (see Figure 22.1),

$$V(t) = M(\Theta_{\text{mdl}}(t); \Theta_{\text{num}}(t)), \quad (22.2)$$

for some function M , originating from the expression of arbitrage-free valuation principles through our chosen model. As $\Theta_{\text{mdl}}(t)$ itself originates from $\Theta_{\text{mkt}}(t)$ and $\Theta_{\text{prm}}(t)$, we may, of course, write

³Recall from Chapter 6 that some yield curve construction algorithms require control parameters (such as tension parameters and precision tolerances), so $\Theta_{\text{prm}}(t)$ may be required for the construction of the discount curve from market inputs.

$$V(t) = H(\Theta_{\text{mkt}}(t); \Theta_{\text{prm}}(t), \Theta_{\text{num}}(t)), \quad (22.3)$$

where H is the overall transfer function that translates market data and control parameters into derivatives values.

Finally, let us quickly note that sometimes the calibration function C will be product-specific, i.e. it will depend on the characteristics of the specific security being valued — recall the discussion of “global” versus “local” calibration in Section 14.5.5. For our purposes here, we ignore this additional level of potential inter-connectivity.

22.1.2 Risk: Theory and Practice

According to basic derivatives pricing theory, the function M in (22.2) assigns value based on dynamic hedging and no-arbitrage principles: the price of a derivative security should equal the cost of hedging the security through its lifetime. In doing so, the model will typically rely on idealized assumptions, e.g. that hedging costs are zero, hedging can take place in continuous time, and so forth. These assumptions are, of course, not true in practice, and will require traders to properly charge⁴ for the cost of running the hedge, as well as for the fact that the hedge is not truly risk-free. A more subtle issue is the fact that the model will compute value based on an assumption of “infallibility” of the parameter estimates $\Theta_{\text{mdl}}(t)$. In particular, once $\Theta_{\text{mdl}}(t)$ has been established, the underlying model will typically assume that, for any $t' > t$,

$$V(t') = L(X(t'); \Theta_{\text{mdl}}(t)), \quad t' > t, \quad (22.4)$$

where $X(t)$ is a random vector of state variables driven by a vector Brownian motion $W(t)$, and L some model-implied map. Our hedging strategy should therefore, as described in Section 1.7, in theory care *only* about neutralizing the effects of movements in X , as caused by W .

In practice, the situation is different. First, actual moves of the yield curve and volatility smiles will inherently deviate from those projected by the model. Second, at time $t' > t$, in reality the model parameter vector $\Theta_{\text{mdl}}(t)$ will be *discarded* and the model calibrated again, followed by an application of (22.2)

⁴For simple securities in simple models, it is possible to derive certain analytical results concerning the costs and risks of actual (discrete-time, costly) hedging strategies. A classical paper in this area is Leland [1985], although the approach has recently received some criticism (see, e.g., Kabanov and Safarian [1997]). Other relevant papers include, among many others, Soner et al. [1995], Barles and Soner [1998], and Derman and Kamal [1999]. As all derivatives traders manage risk at the *portfolio* (or *book*) level, the “net” security owned by traders is far more complicated than those covered by most of the literature. In addition, traders tend to rebalance their trading books according to rules much more complex than those assumed in academic papers. As a result, proper charging for transactions normally requires a heavy element of human judgment, and will depend strongly on the portfolio context.

to establish the new value of the portfolio as $V(t') = M(\Theta_{\text{mdl}}(t'); \Theta_{\text{num}}(t'))$. Of course, the recalibrated model parameters $\Theta_{\text{mdl}}(t')$ will rarely, if ever, be consistent with those used at time t , so equation (22.4) will generally fail. As this equation serves as a fundamental assumption of the underlying model, the practical usage of the model is clearly causing quite profound consistency violations. In particular, we constantly change parameters that are assumed by the model to be invariants.

While at first glance the situation outlined above may seem to strike a death blow to the entire foundation of derivatives pricing and hedging, there are several mitigating factors that make the situation less dire than it appears. First, if the model is fundamentally sound, its dynamics will be close to reality most of the time, and (22.2) and (22.4) will consequently be near-identical on average. Second, the trader can employ several strategies to minimize the risk of model mis-specification. One type of strategy involves the use of robust static or super-replicating hedges, as in Sections 16.6.1, 19.4.5 and 20.2.3. As this is not always possible, a more common strategy involves hedging “too much”, by neutralizing the portfolio to higher-order sensitivities (gamma hedging), and also by hedging against moves in quantities which are assumed by the model to be non-random. A standard example of the latter is the practice of vega hedging with the Black-Scholes model: despite the fact that the model assumes that the volatility is a constant, the dealer will nevertheless put on a hedge against moves in the volatility parameter. As discussed in Hull [2006] or Taleb [1997], there is empirical evidence that such practices considerably improve the hedge robustness and performance in actual markets.

Vega hedging is an example of the common practice of ignoring the theoretical ideal of (22.4), and instead constructing a hedge around (22.3), with the hedge aiming to neutralize (in a standard Taylor-series sense, see Section 22.1.5) as many of the movements in the entire market data vector $\Theta_{\text{mkt}}(t)$ as possible, irrespective of whether a particular model may suggest that this is reasonable or not. As the dimension of the market data vector N_{mkt} can be very high, often 100 or more, it may be too costly and too onerous to hedge against all components of $\Theta_{\text{mkt}}(t)$ individually, so some type of principal components analysis (as in Section 14.3.1, for instance) may be undertaken to guide the level of granularity required in the hedge. Additionally, one would need to contemplate whether hedges against both first- and second-order (or even higher orders) risk are required. The answer to this question would typically be settled by careful analysis of the convexity properties of the portfolio value V as a function of the market data $\Theta_{\text{mkt}}(t)$: whenever there is significant convexity or concavity with respect to a given parameter, it is reasonable to attempt to put on a second-order hedge. As the (Hessian) matrix of second-order derivatives of V with respect to the components of $\Theta_{\text{mkt}}(t)$ has $N_{\text{mkt}}(N_{\text{mkt}} + 1)/2$ distinct elements (which will often reach thousands), again some selectivity will be required in practice.

The “art of derivatives trading” — that is, the practice of cost-efficiently managing of book of derivatives from market data sensitivity reports — requires considerable and detailed market knowledge and is hard, if not impossible, to describe in purely mathematical terms. Consequently we abstain from attempting to do so, but simply notice that the very foundation of the trading exercise are the market data sensitivities themselves. For readers interested in a description of derivative trading practices, the available material is, unfortunately, rather limited. A common reference is Taleb [1997]; Miron and Swannell [1991] could also be consulted.

22.1.3 Example: the Black-Scholes Model

To make the discussion in Sections 22.1.1 and 22.1.2 more concrete, let us assume that our derivatives portfolio is written on a single underlying asset $X(t)$, the risk-neutral dynamics of which are

$$dX(t)/X(t) = r dt + \sigma dW(t), \quad (22.5)$$

where r and σ are constants and $W(t)$ is a one-dimensional Brownian motion. We shall shortly (in Section 22.1.4) make the model a bit more realistic by introducing time-dependence into the parameters, but for now we assume that they are constants. In addition, we assume that r and σ are directly observable in the market (again, we relax this shortly), such that $\Theta_{\text{mkt}}(t) = \Theta_{\text{mdl}}(t) = (X(t), r, \sigma)^{\top}$.

The theoretical hedge for a derivative (or a derivative portfolio) V written on X will take the form of a pure delta hedge, where a position $(-\partial V / \partial X)$ in X is maintained at all times. In practice, a trader will not only be concerned with neutralizing against first-order movements in $X(t)$, but will also manage many of the instantaneous sensitivities listed in Table 22.1.

Common Name	Definition
theta	$\partial V / \partial t$
rho	$\partial V / \partial r$
delta	$\partial V / \partial X$
gamma	$\partial^2 V / \partial X^2$
vega	$\partial V / \partial \sigma$
volga (or vomma)	$\partial^2 V / \partial \sigma^2$
vanna	$\partial^2 V / \partial \sigma \partial X$

Table 22.1. Common greeks in the Black-Scholes model.

Neutralization of rho will typically involve taking positions in interest rate swaps, whereas gamma, vega, volga, vanna can only be eliminated by trading derivative securities with non-linear payouts in X and volatility exposure — typically liquid European options.

We should note that in Table 22.1, the theta (or time decay) has a special status, as the passage of time is both unavoidable and non-random — as such, it makes no sense to try to hedge this “exposure”. We also note that the theta, in a sense, emerges as a combination of other greeks, as can be confirmed from the basic Black-Scholes-Merton valuation PDE from Section 1.9:

$$\frac{\partial V}{\partial t} + rX \frac{\partial V}{\partial X} + \frac{1}{2}\sigma^2 X^2 \frac{\partial^2 V}{\partial X^2} = rV \quad (22.6)$$

or, equivalently,

$$\text{theta} + rX \times \text{delta} + \frac{1}{2}\sigma^2 X^2 \times \text{gamma} = r \times \text{value}.$$

Notice, in particular, that a delta hedged position (delta = 0), will identify theta = $rV - \frac{1}{2}\sigma^2 X^2 \times \text{gamma}$, so for a delta hedger the time decay of his position originates in part from a pure discount effect due to interest accruing on the net present value of the portfolio (the term rV), and in part from a convexity term (the term $-\frac{1}{2}\sigma^2 X^2 \times \text{gamma}$). The latter, of course, represents optionality leaking away over time, and not surprisingly scales with σ^2 .

While traders tend to monitor all entries in Table 22.1, particular attention is typically paid to the delta, gamma and vega exposures. Hedging the vega ensures that inevitable changes to implied volatility will not lead to large moves in the portfolio value, and hedging the gamma helps prevent “slippage” in between dates where the delta hedge is rebalanced. As discussed already in Section 19.7.1, most traders strongly prefer being long gamma exposure (i.e. the net gamma of their portfolio, including hedges, is positive), as a negative gamma position can lose a very substantial amount of money in periods of financial turmoil when the trader cannot adjust his delta hedge quickly enough to track the market.

One would intuitively expect that a portfolio that has low gamma exposure would also have a low vega, a result that in fact can be formalized as follows:

Lemma 22.1.1. *Consider a European-style claim with maturity T and payout function $g(X(T))$, where $X(t)$ satisfies (22.5). With $V(t) = V(t, X(t))$ denoting time t value of this security, we have (employing somewhat loose notation)*

$$\frac{\partial V(t)}{\partial \sigma} = (T - t)X(t)^2 \sigma \frac{\partial^2 V(t)}{\partial X(t)^2}, \quad \text{or} \quad \text{vega} = (T - t)X(t)^2 \sigma \times \text{gamma}. \quad (22.7)$$

Proof. An elementary, if rather inelegant, proof proceeds as follows. First, we let E_t denote time t risk-neutral expectation, such that

$$V(t) = e^{-r(T-t)} E_t(g(X(T))),$$

where

$$X(T) = X(t) \exp \left(\left(r - \frac{1}{2}\sigma^2 \right) (T-t) + \sigma\sqrt{T-t}Z \right),$$

with $Z \sim \mathcal{N}(0, 1)$. Notice that

$$\frac{\partial V(t)}{\partial X(t)} = e^{-r(T-t)} \mathbb{E}_t \left(\frac{\partial g(X(T))}{\partial X(T)} \cdot \frac{X(T)}{X(t)} \right), \quad (22.8)$$

and

$$\frac{\partial X(T)}{\partial t} = -X(T) \left(r - \frac{1}{2}\sigma^2 + \frac{\sigma}{2\sqrt{T-t}}Z \right),$$

such that

$$\begin{aligned} \frac{\partial V(t)}{\partial t} &= r e^{-r(T-t)} \mathbb{E}_t(g(X(T))) + e^{-r(T-t)} \mathbb{E}_t \left(\frac{\partial g(X(T))}{\partial X(T)} \cdot \frac{\partial X(T)}{\partial t} \right) \\ &= rV(t) - e^{-r(T-t)} \mathbb{E}_t \left(\frac{\partial g(X(T))}{\partial X(T)} X(T) \left(r - \frac{1}{2}\sigma^2 + \frac{\sigma}{2\sqrt{T-t}}Z \right) \right) \\ &= rV(t) - X(t) \frac{\partial V(t)}{\partial X(t)} \left(r - \frac{1}{2}\sigma^2 \right) \\ &\quad - \frac{\sigma}{2\sqrt{T-t}} e^{-r(T-t)} \mathbb{E}_t \left(\frac{\partial g(X(T))}{\partial X(T)} \cdot X(T)Z \right), \end{aligned} \quad (22.9)$$

where we have used (22.8) in the last equality. Using the same principles, we get

$$\begin{aligned} \frac{\partial V(t)}{\partial \sigma} &= e^{-r(T-t)} \mathbb{E}_t \left(\frac{\partial g(X(T))}{\partial X(T)} \cdot \frac{\partial X(T)}{\partial \sigma} \right) \\ &= e^{-r(T-t)} \mathbb{E}_t \left(\frac{\partial g(X(T))}{\partial X(T)} \cdot X(T) \left(-\sigma(T-t) + \sqrt{T-t}Z \right) \right) \\ &= -X(t) \frac{\partial V(t)}{\partial X(t)} \sigma(T-t) + e^{-r(T-t)} \sqrt{T-t} \mathbb{E}_t \left(\frac{\partial g(X(T))}{\partial X(T)} X(T)Z \right). \end{aligned} \quad (22.10)$$

Combining (22.9) and (22.10), we get

$$\begin{aligned} \frac{\partial V(t)}{\partial \sigma} &= -X(t) \frac{\partial V(t)}{\partial X(t)} \sigma(T-t) \\ &\quad + \frac{2(T-t)}{\sigma} \left(rV(t) - \frac{\partial V(t)}{\partial t} - X(t) \frac{\partial V(t)}{\partial X(t)} \left(r - \frac{1}{2}\sigma^2 \right) \right) \\ &= \frac{2(T-t)}{\sigma} \left(rV(t) - \frac{\partial V(t)}{\partial t} - rX(t) \frac{\partial V(t)}{\partial X(t)} \right). \end{aligned}$$

The result (22.7) follows after insertion of the pricing PDE (22.6) into the expression above. \square

Remark 22.1.2. A more elegant way to prove Lemma 22.1.1 relies on operator calculus, as in Carr [2000]. As this technique is sometimes quite handy, we demonstrate it in Appendix 22.A, where it is also shown that, for a general European-style claim,

$$\frac{\partial V(t)}{\partial r} = (T - t) \left(X(t) \frac{\partial V(t)}{\partial X(t)} - V(t) \right).$$

While Lemma 22.1.1 only holds for European-style claims in the Black-Scholes model, the observation that there is a close link between vega and gamma holds in general. We shall examine the importance of vega and gamma hedging in more detail shortly, but first let us make the model setting a bit more realistic.

22.1.4 Example: Black-Scholes Model with Time-Dependent Parameters

In actual usage of the Black-Scholes model, one would always allow the short rate and volatility to be time-dependent, in order to match observed term structures of discount bonds and (term) option volatilities. While technically an obvious extension to Section 22.1.3, let us go over the mechanics nevertheless to better illustrate the concepts from Sections 22.1.1 and 22.1.2. We first change (22.5) to

$$dX(t)/X(t) = r(t) dt + \sigma(t) dW(t),$$

where $r(t)$ and $\sigma(t)$ are deterministic functions of time, to be calibrated to the term structure of discount bonds and implied at-the-money (ATM) volatilities observed in the market. We assume that the yield curve is computed from a vector of swap yields $S(t) = (S_1(t), \dots, S_J(t))^\top$, where it is understood that $S_i(t)$ represents the time- t par yield of a swap that matures on some date $T_i > t$, $i = 1, \dots, J$. For simplicity, let us assume, as in Section 6.2.1, that the constructed yield curve is bootstrapped as a piecewise flat curve in forward rate space, with breakpoints located at the swaption maturities $\{T_i\}$. Defining $T_0 \triangleq t$, the resulting time t forward curve may therefore be represented as

$$f(t, u) = -\frac{\partial \ln P(t, u)}{\partial u} = \sum_{i=0}^{J-1} \gamma_i(t) 1_{\{u \in [T_i, T_{i+1})\}},$$

where we use the vector $\gamma(t) = (\gamma_1(t), \dots, \gamma_J(t))^\top$ to store the resulting J forward curve levels. As r is assumed deterministic, we set $r(u) = f(t, u)$, $u \geq t$, and have then completed the interest rate calibration.

To construct $\sigma(t)$, assume that a vector $v(t) = (v_1(t), \dots, v_D(t))^\top$ is observed in the market, where each $v_i(t)$ represents a term volatility to t_i on a specified maturity grid $\{t_i\}_{i=1}^D$, with $t_1 > t$. That is,

$$v_i(t) = \sigma_{\text{ATM}}(t, t_i),$$

where $\sigma_{\text{ATM}}(t, t_i) \triangleq \sigma_{\text{BS}}(t, X(t); t_i, X(t))$ is the ATM implied (Black-Scholes) volatility to time t_i , seen from time t (see Section 7.1.2). If we assume that $\sigma(u)$ is piecewise flat on the $\{t_i\}$ -grid, we can construct $\sigma(u)$ by bootstrapping from the basic relation (see Section 1.9.3)

$$\sigma_{\text{ATM}}(t, u)^2 = (u - t)^{-1} \int_t^u \sigma(s)^2 ds. \quad (22.11)$$

The result of this exercise is a vector $\varsigma(t) = (\varsigma_1(t), \dots, \varsigma_D(t))^{\top}$ of flat volatility levels, such that

$$\sigma(u) = \sum_{i=0}^{D-1} \varsigma_i(t) \mathbf{1}_{\{u \in [t_i, t_{i+1})\}},$$

where we set $t_0 \triangleq t$.

The calibration procedure described above turns $\Theta_{\text{mkt}}(t) = (S(t), v(t))^{\top}$ into $\Theta_{\text{mdl}}(t) = (\gamma(t), \varsigma(t))^{\top}$, with both vectors having dimension $J + D$, i.e. $N_{\text{mkt}} = N_{\text{prm}} = J + D$. The vector $\Theta_{\text{prm}}(t)$ is here empty, but would have contained precision parameters if, say, the more elaborate yield curve construction algorithm of Section 6.3 had been used. The contents of the vector $\Theta_{\text{num}}(t)$ would depend on what numerical method the calibrated model would implement for the purpose of pricing a specific derivative. For instance, if we were using a finite difference grid, $\Theta_{\text{num}}(t)$ might contain a confidence level multiplier (see Section 2.1) to dimension the grid, a θ -parameter to determine the level of implicitness in the solver (see Section 2.2.3), various flags (e.g. whether to use upwinding or not, see Section 2.6.1) and, of course, information to determine the number of grid points in the time (t) and space (X) directions.

As r and σ are now vector-valued, the quantities rho, vega, volga, and vanna in Table 22.1 are no longer scalars⁵ but must be represented as vectors. Also, one issue is how we wish to present this risk information in the first place: for vega, say, do we want to report sensitivities with respect to the market volatilities v (so-called *market vegas*) or the model volatilities ς (so-called *model vegas*)? The former is the most common, but we can here freely translate between the two. Specifically, applying the chain rule to the relationship between v_i 's and ς_j 's given by (22.11) we have that

$$\frac{\partial}{\partial v_i} = \sum_{j=1}^D \frac{\partial}{\partial \varsigma_j} \left(\frac{\partial \varsigma_j}{\partial v_i} \right),$$

⁵If, say, a single vega is nevertheless required, it is often most reasonable to report the sensitivity to a parallel shift of the function $\sigma(t)$ at all values of t . We can think of this as roughly representing the sensitivity with respect to the first principal component of volatility curve moves. A similar principle can be applied to compute a single rho.

where the matrix of partial derivatives $\partial \varsigma / \partial v$ can be obtained by inverting the Jacobian matrix $\partial v / \partial \varsigma$ which can be obtained in closed form from the relation (22.11). A similar translation between sensitivities to swap yields and to forward rate buckets was discussed already, in Section 6.4.

22.1.5 Actual Risk Computations

Getting back to the general representation (22.3), assume that we perturb the market data by a vector-valued amount $\delta = (\delta_1, \dots, \delta_{N_{\text{mkt}}})^\top$. Let us use a Taylor expansion to write (dropping the argument t for clarity)

$$V(\delta) = H(\Theta_{\text{mkt}} + \delta) \approx H(\Theta_{\text{mkt}}) + \nabla^H \cdot \delta + \frac{1}{2} \delta^\top \cdot A^H \cdot \delta, \quad (22.12)$$

where ∇^H is an N_{mkt} -dimensional row vector and A^H an $N_{\text{mkt}} \times N_{\text{mkt}}$ matrix, with elements

$$\nabla_i^H = \left. \frac{\partial H(x)}{\partial x_i} \right|_{x=\Theta_{\text{mkt}}}, \quad A_{i,j}^H = \left. \frac{\partial^2 H(x)}{\partial x_i \partial x_j} \right|_{x=\Theta_{\text{mkt}}}, \quad i, j = 1, \dots, N_{\text{mkt}}.$$

While the situation in interest rate modeling is obviously a little more complicated than for the single-asset Black-Scholes setup in Sections 22.1.3 and 22.1.4, loosely speaking the gradient vector ∇^H will contain deltas (first-order sensitivities with respect to swap yields) and vegas (first-order sensitivities with respect to swaption volatilities), while the Hessian matrix A^H will contain gammas, volgas, and vannas. Notice that the risk measure rho is not used for interest rate derivatives (where, in a sense, delta and rho coincide).

Although not all elements of A are always requested, in a nutshell the main role of an interest rates derivatives risk system is to report⁶ ∇^H and A^H for consumption by the trading desk, risk management, and the middle office. The trading desk will, as we discussed earlier in Section 22.1.2, use the sensitivities to evaluate how much it should rebalance the portfolio to keep it broadly market-neutral and robust to market shocks; the “ideal” configuration for a pure⁷ hedger will obviously be to arrange the portfolio such that $\nabla^H = 0$ and $A^H = 0$. Risk management will use the sensitivities to ensure that the exposures to individual market data components are within given sensitivity *limits*. In addition, risk management will typically

⁶Sometimes the risk system is asked to perform the Taylor expansion around a series of different (perturbed) market data scenarios, not just the current market data. The resulting collections of reports are known as *ladders*.

⁷Most derivatives traders have views on the future market evolution and are allowed to express their views in proprietary (“prop”) positions, meant to make money if the trader’s views turn out to be correct. In this situation, the trader will purposely leave the portfolio open to certain market risk exposures. Banks typically enforce strict limits on the size of these exposures.

compute an overall measure of the portfolio risk, based on a statistical (or historical) model for the perturbation vector δ over a given time horizon, typically one day. We discuss this computation briefly in Section 22.3. In the middle office of a bank, the sensitivities are used for *P&L analysis*, i.e. the process of reconciling observed moves in portfolio value with changes in market data. We review this process in Section 22.2.

Recall that market data input used to compute the value of a derivative security typically goes through a two-step procedure, where first a calibration turns the market data vector Θ_{mkt} into a model data vector Θ_{mdl} , which in turn is used to compute the derivatives price, $V = M(\Theta_{\text{mdl}})$. It is often natural to first compute sensitivities with respect to the model parameters — a process that is independent of the chosen calibration procedure — and then combine these sensitivities with calibration-specific sensitivity information to compute the market data input sensitivities. For instance, we would write, for $j = 1, \dots, N_{\text{mkt}}$,

$$\nabla_j^H = \sum_{i=1}^{N_{\text{mdl}}} \frac{\partial M(y)}{\partial y_i} \frac{\partial C_i(x)}{\partial x_j} \Big|_{y=\Theta_{\text{mdl}}, x=\Theta_{\text{mkt}}},$$

where C_i is the i -th component of the N_{mdl} -dimensional calibration function C in (22.1). The $N_{\text{mdl}} \times N_{\text{mkt}}$ matrix J with elements

$$J_{i,j} = \frac{\partial C_i}{\partial x_j}$$

is known as the *Jacobian* for the map from market to model parameters. The Jacobian can normally be computed in the calibration module, as part of the calibration procedure itself. We saw a simple example of Jacobian matrix usage in Section 22.1.4 above, and will consider the idea in a more realistic setting in Section 26.3.3.

22.1.6 What about Θ_{prm} and Θ_{num} ?

The reader will have noticed that the portfolio value depends not only on market data, but also on various technical parameters that control numerics and the calibration, as well as certain unobservable model parameters. This type of data is fairly static, and sensitivity information is rarely reported on a running basis. As the numbers do affect the official profit-and-loss (P&L) produced by the model, the elements of Θ_{prm} and Θ_{num} are typically supervised by control groups that may impose standards on numerical parameters (e.g. require that the numerical error be within a certain tolerances) and may request that monetary buffers — so-called *reserves* — be set aside to cover the uncertainty of unobservable model parameters. The latter will require some estimates of the uncertainty associated with a given parameter, as well as a computation of the portfolio value sensitivity to the parameters.

While the reserves need to be dynamically updated to reflect changes in the portfolio and in the parameter sensitivities, this is normally done relatively infrequently, e.g. every month or quarter. Given this, from a computational perspective sensitivity generation with respect to market data Θ_{mkt} — which is often done on an inter-day basis — is, by far, the more challenging task. As such, our emphasis in the rest of this book is solely on computation of market data sensitivities.

22.1.7 A Note on Trading P&L and the Computation of Implied Volatility

Before proceeding to discuss applications of the sensitivity analysis of Section 22.1.5, we insert a brief interlude to demonstrate an important result (sometimes known as the *fundamental theorem of derivatives trading*) that provides a link between a portfolio's gamma and expected hedging P&L over a given horizon. The setup is as follows. At time 0 a trader buys a contingent claim on a single non-dividend paying asset X , and chooses to value his position by using a Black-Scholes model with fixed volatility σ_{BS} . Let the trader's mark for his portfolio be $V_{\text{BS}}(t)$ and assume that σ_{BS} is such that the value $V_{\text{BS}}(0)$ coincides with the time 0 market value. We assume that the contingent claim expires at time T with value $g(X(T))$, and pays no cash flows before then. The trader is actively hedging his position, but commits two “sins”: i) he does not gamma or vega hedge his position, but only delta hedges; and ii) he never re-calibrates the model but assumes that σ_{BS} is the correct volatility to use when computing hedge information, even if the volatility of X is observed to change over time.

In analyzing the performance of the hedger's strategy, let us assume that the volatility of X is a random process $\sigma(t)$, i.e. that dynamics in the real-life measure P are of the form

$$dX(t)/X(t) = O(dt) + \sigma(t) dW^P(t),$$

where $W^P(t)$ is a P -Brownian motion. Now, to hedge his long position in the contingent claim, the trader sets up a short position in a portfolio Π with $n_X(t)$ units of X held at time t , along with a cash position $N(t)$. As described above, the trader delta hedges according to the Black-Scholes model, so

$$n_X(t) = \frac{\partial V_{\text{BS}}(t)}{\partial X(t)}.$$

In other words, we have

$$\Pi(t) = n_X(t)X(t) + N(t),$$

where, by construction, $\Pi(0) = V_{\text{BS}}(0)$. Assuming that $N(t)$ is rolled over at the short-term interest rate r (assumed constant for convenience), we therefore get

$$d\Pi(t) = n_X(t) dX(t) + N(t)r dt, \quad (22.13)$$

where the self-financing condition (see (1.10)) justifies ignoring the change in n_X at $t + dt$. The following important result now holds.

Proposition 22.1.3. *The time T terminal value of the delta hedging account in (22.15) is*

$$\Pi(T) = g(X(T)) + \frac{1}{2}e^{rT} \int_0^T e^{-rt} (\sigma_{\text{BS}}^2 - \sigma(t)^2) X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt,$$

where $g(x)$ is the terminal payout function.

Proof. By Ito's lemma, observe that

$$dV_{\text{BS}}(t) = \frac{\partial V_{\text{BS}}(t)}{\partial t} dt + \frac{\partial V_{\text{BS}}(t)}{\partial X(t)} dX(t) + \frac{1}{2} \sigma(t)^2 X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt, \quad (22.14)$$

so combining (22.13) and (22.14) yields

$$\begin{aligned} d(V_{\text{BS}}(t) - \Pi(t)) &= \frac{\partial V_{\text{BS}}(t)}{\partial t} dt + \frac{1}{2} \sigma(t)^2 X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt - N(t)r dt \\ &= \frac{\partial V_{\text{BS}}(t)}{\partial t} dt + \frac{1}{2} \sigma(t)^2 X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt \\ &\quad - r \left(\Pi(t) - \frac{\partial V_{\text{BS}}(t)}{\partial X(t)} X(t) \right) dt. \end{aligned} \quad (22.15)$$

We now recall that $V_{\text{BS}}(t)$ satisfies the Black-Scholes PDE with constant volatility σ_{BS} , wherefore

$$\frac{\partial V_{\text{BS}}(t)}{\partial t} + rX(t) \frac{\partial V_{\text{BS}}(t)}{\partial X(t)} + \frac{1}{2} \sigma_{\text{BS}}^2 X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} = rV_{\text{BS}}(t).$$

Inserting this expression into (22.15) yields, after a little algebra,

$$\begin{aligned} d(V_{\text{BS}}(t) - \Pi(t)) &= \frac{1}{2} (\sigma(t)^2 - \sigma_{\text{BS}}^2) X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt + r(V_{\text{BS}}(t) - \Pi(t)) dt. \end{aligned}$$

We can integrate this equation to yield

$$\begin{aligned} V_{\text{BS}}(T) - \Pi(T) &= e^{rT} (V_{\text{BS}}(0) - \Pi(0)) \\ &\quad + e^{rT} \int_0^T e^{-rt} \frac{1}{2} (\sigma(t)^2 - \sigma_{\text{BS}}^2) X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt, \end{aligned}$$

and the result of the proposition follows from the observation that $V_{\text{BS}}(T) = g(X(T))$ and $V_{\text{BS}}(0) = \Pi(0)$. \square

Proposition 22.1.3 demonstrates that the hedging strategy followed by the trader generally does not work, in the sense that the terminal value $\Pi(T)$ of the self-financing hedging portfolio will fail to equal $g(X(T))$. In certain special cases, however, the hedge will work, e.g. when the hedged claim is gamma-neutral ($\partial^2 V_{\text{BS}} / \partial X^2 = 0$) or when $\sigma(t)$ is close to σ_{BS} “on average”. These observations, while trivial, strongly support the strategy of re-calibrating the model to changing market conditions and to use gamma-hedging to keep portfolio convexity low.

As an aside, we notice that if i) the Black-Scholes gamma is strictly positive, and ii) the realized volatility $\sigma(t)$ is consistently higher than σ_{BS} , we have $\Pi(T) < g(X(T))$ for sure. As the trader is short the hedging portfolio, it follows that the trader keeping a positive gamma benefits from financial turmoil (high volatility), a point we have made several times already.

Finally, let us present an important corollary of Proposition 22.1.3.

Corollary 22.1.4. *Let the claim in Proposition 22.1.3 be a European call or put option with strike K . The time 0 implied volatility $\sigma_{\text{BS}} = \sigma_{\text{BS}}(0, X(0); T, K)$ is given by*

$$\sigma_{\text{BS}}^2 = \frac{\mathbb{E} \left(\int_0^T e^{-rt} \sigma(t)^2 X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt \right)}{\mathbb{E} \left(\int_0^T e^{-rt} X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt \right)}, \quad (22.16)$$

where \mathbb{E} denotes expectation in the risk-neutral measure, and

$$X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} = X(t) \frac{\phi(d_+(X(t)))}{\sigma_{\text{BS}} \sqrt{T-t}}, \quad (22.17)$$

$$d_+(x) = \frac{\ln(x/K) + (r + \frac{1}{2}\sigma_{\text{BS}}^2)(T-t)}{\sigma_{\text{BS}} \sqrt{T-t}},$$

where $\phi(x) = (2\pi)^{-1/2} e^{-x^2/2}$ is the Gaussian density.

Proof. The hedge portfolio generates no cash flows on $[0, T]$, so its time 0 value must equal

$$\begin{aligned} \Pi(0) &= e^{-rT} \mathbb{E}(\Pi(T)) = e^{-rT} \mathbb{E}(g(X(T))) \\ &\quad + \mathbb{E} \left(\int_0^T e^{-rt} \frac{1}{2} (\sigma(t)^2 - \sigma_{\text{BS}}^2) X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt \right), \end{aligned}$$

where $g(x) = (x - K)^+$ for a call and $g(x) = (K - x)^+$ for a put. Here, the term

$$e^{-rT} \mathbb{E}(g(X(T)))$$

equals the time 0 market value of the put or call being hedged, and equals $\Pi(0)$ by assumption. It follows that

$$\mathbb{E} \left(\int_0^T e^{-rt} \frac{1}{2} (\sigma(t)^2 - \sigma_{\text{BS}}^2) X(t)^2 \frac{\partial^2 V_{\text{BS}}(t)}{\partial X(t)^2} dt \right) = 0,$$

which immediately leads to (22.16). The result (22.17) follows from an explicit evaluation of gamma in the Black-Scholes model. \square

To get some insights into the result of the corollary above, assume that $\sigma(t)$ is of the local volatility type, $\sigma(t) = \sigma(X(t))$, in which case (22.16) can be written

$$\sigma_{\text{BS}}(T, K)^2 = \frac{\int_0^\infty \int_0^T e^{-rt} \sigma(x) w(t, x) \psi(t, x) dt dx}{\int_0^\infty \int_0^T e^{-rt} w(t, x) \psi(t, x) dt dx}, \quad (22.18)$$

where $\psi(t, x)$ is the density of $X(t)$ as seen from time 0, and

$$w(t, x) = x \frac{\phi(d_+(x))}{\sigma_{\text{BS}}(T, K) \sqrt{T-t}}.$$

In a sense, (22.18) demonstrates that implied volatility is a weighted average of local volatility, where weights are proportional to the product of gamma and the asset density.

Direct usage of (22.18) is complicated by the fact that the implied volatility figures in both the left- and right-hand sides of the equation, and by the fact that the density of the asset is rarely, if ever, known explicitly. On the other hand, the product $w(t, x) \psi(t, x)$ can be seen to typically form a “ridge” from $x = X(0)$ at time 0 to $x = K$ at time T , a result that holds irrespective of the model specification. This, among other considerations, has inspired some authors to suggest that

$$\sigma_{\text{BS}}(0, X(0); T, K)^2 \approx \frac{1}{T} \int_0^T \mathbb{E}(\sigma(t)^2 | X(t) = x^*(t)) dt,$$

where $x^*(t)$ is to be interpreted as *the most likely path from $X(0)$ to K* . This idea has found applications for both local and stochastic volatility models, see e.g. Gatheral [2006]. While often quite intuitive, approximation techniques based on (22.16) involve a fair amount of heuristics⁸, and precision is often neither impressive nor easy to characterize. As a result, the method — which we believe was originally suggested by Bruno Dupire — is often reserved for qualitative analysis. See Lee [2005] and Gatheral [2006] for additional discussion and applications.

22.2 P&L Analysis

Besides being used by traders to manage the exposure of their books, the sensitivity information contained in Taylor expansions such as (22.12) is

⁸In a pinch, it is often reasonable to simply assume that $x^*(t) = T^{-1}(X(0)(T-t) + Kt)$.

consumed by various support and control functions in a bank. We discuss two such uses: P&L analysis and, in Section 22.3 below, value-at-risk computation.

22.2.1 P&L Predict

The expansion (22.12) may be used in an accounting analysis to analyze and reconcile the realized P&L from one trading day to the next. Although the analysis is carried out at time $t + h$ (where the market data movement δ is known), it is known as a *P&L prediction analysis*, or just a *P&L predict*. Given, at time t , expansion terms $\nabla^H(t)$ and $A^H(t)$, if the observed market data movements over the period $[t, t + h]$ (with h typically equal to one business day) is δ , then, all things equal, we would expect the time $t + h$ portfolio value to be approximately

$$V(t + h) \approx V(t) + \frac{\partial V(t)}{\partial t} h + \nabla^H(t) \cdot \delta + \frac{1}{2} \delta^\top \cdot A^H(t) \cdot \delta. \quad (22.19)$$

Notice the inclusion of the term $\partial V / \partial t$ (theta) in this expansion, to account for the passage of time. The right-hand side of this equation is known as the *second-order P&L predict*; if we omit the convexity term (i.e. set $A^H(t) = 0$), the right-hand side is, naturally, the *first-order P&L predict*. The difference between the right- and left-hand sides of (22.19) may be called the *unpredicted P&L*. If systems and models are working properly, the P&L predict should generally be an accurate and unbiased estimated of actual P&L, so monitoring of the unpredicted P&L serves an important control purpose. Unusually large values of unpredicted P&L may, for instance, hint at problems in the computation of risk sensitivities (and therefore in hedges) or suggest that the portfolio is exposed to large unhedged high-order risks.

We should note that when writing down (22.19), we implicitly made several simplifying assumptions, most notably that the portfolio at time t is the same as at time $t + h$. In reality, trades may expire, get canceled or amended, or entirely new trades may be added to the portfolio on the interval $(t, t + h]$. In addition, cash payments (coupons and settlement amounts) may take place on $(t, t + h]$ and must be added to the left-hand side of (22.19). The function V in (22.19) should therefore really be thought of as representing the part of the portfolio trade population that involves no special events over the period $[t, t + h]$; a full P&L predict analysis will additionally require accounting for a number of adjustments due to cash payments, changes to the portfolio, and rate fixings. Getting all details right is often a fairly complex exercise, and as mentioned earlier is normally handled by dedicated personnel in a bank's middle office.

An important issue in P&L predict (and also in P&L explain, see Section 22.2.2 below) is the computation of the theta term in (22.19). In particular, when advancing time forward, what precisely is it that we should hold fixed?

While one might say that, by definition, all elements of $\Theta_{\text{mkt}}(t)$ should stay at their time t values, this generally causes problems. To demonstrate, assume, as is common, that the yield curve is constructed from a series of Eurodollar contracts, and swap quotes. As discussed in Section 5.4, a Eurodollar futures contract will settle at a fixed point in time — i.e. it has a fixed time *of* maturity — so when advancing time from t to $t + h$, its remaining time *to* maturity will shrink by an amount h . On the other hand, a market-quoted swap always is associated with a standardized time *to* maturity (a fixed swap tenor), so no maturity shrinkage occurs when time is advanced. In total, when advancing time forward, the time to maturity of some, but not all, yield curve instruments will undergo a change. This is not compensated for by a change in the market quote (which is held fixed), which results in an effective move in the forward curve that typically will be highly erratic and entirely unsuitable for a perturbation analysis. To avoid problems of this type, and to properly reflect short-term funding costs, it is natural to compute the expected change in market data

$$\Theta_{\text{mkt}}^f(t) = \mathbb{E}^{t+h} (\Theta_{\text{mkt}}(t+h)|\mathcal{F}_t), \quad (22.20)$$

and then write

$$\frac{\partial V(t)}{\partial t} \approx \frac{V(t+h; \Theta_{\text{mkt}}^f(t)) - V(t)}{h}. \quad (22.21)$$

The term $V(t+h; \Theta_{\text{mkt}}^f(t))$ may be computed by an outright re-valuation of the portfolio after i) advancing calendar time to $t+h$; and ii) moving the market data to $\Theta_{\text{mkt}}^f(t)$. In computation of $\Theta_{\text{mkt}}^f(t)$, the discount curve constructed at time t will simply be “rolled” up to its time $t+h$ forward curve as seen from time t , i.e. for any $T > t+h$ we set

$$P(t+h, T) = \mathbb{E}^{t+h} (P(t, T)|\mathcal{F}_t) = \frac{P(t, T)}{P(t, t+h)} \approx P(t, T) (1 + r(t)h),$$

where $r(t)$ is the short rate. Moving the discount curve in this fashion is consistent with the notion that a risk-free portfolio should earn a rate of $r(t)$ over a short holding period and rationally anchors the P&L predict analysis around forward values of discount bonds, i.e. values that can be locked in by a risk-free trading strategy at time t .

With the choice (22.21) for theta, we have, in effect, moved the expansion point of the Taylor series (22.19) from $\Theta_{\text{mkt}}(t)$ to $\Theta_{\text{mkt}}^f(t)$, which suggest a modified (and improved) expression for the P&L predict:

$$V(t+h) \approx V(t+h; \Theta_{\text{mkt}}^f(t)) + \nabla^H(t) \cdot (\delta - \delta^f) + \frac{1}{2} (\delta - \delta^f)^\top \cdot A^H(t) \cdot (\delta - \delta^f), \quad (22.22)$$

with

$$\delta^f \triangleq \Theta_{\text{mkt}}^f(t) - \Theta_{\text{mkt}}(t). \quad (22.23)$$

In (22.22), the term $V(t+h; \Theta_{\text{mkt}}^f(t))$ represents the value that the portfolio will reach if the market data moves “according to expectations”, and remaining terms add first- and second-order corrections based on the deviation of the time $t+h$ market data away from its time t expectation,

$$\delta - \delta^f = \Theta_{\text{mkt}}(t+h) - \Theta_{\text{mkt}}(t) - (\Theta_{\text{mkt}}^f(t) - \Theta_{\text{mkt}}(t)) = \Theta_{\text{mkt}}(t+h) - \Theta_{\text{mkt}}^f(t).$$

For the reasons explained earlier, (22.22) is typically preferable to (22.19), yet it is not uncommon for P&L analysis systems to implement both.

Finally, let us note that in our description of the P&L predict process we assumed that interest rate risk was captured as sensitivities with respect to market quotes of yield curve instruments, i.e. we start from a par-point report, in the convention of Section 6.4. As discussed in that section, it is, however, not uncommon to instead capture interest rate risk through sensitivities with respect to buckets of the forward curve itself (a forward rate report). This change in approach is easily accommodated by the methodology above, by simply altering the definitions of $\Theta_{\text{mkt}}(t)$, δ , and δ^f accordingly.

22.2.2 P&L Explain

The objective of a *P&L explain*⁹ analysis is to estimate the contribution of each component of the market vector move δ to the overall move in the portfolio value. In a sense, such information is also captured in the P&L predict (through the sensitivities ∇^H and, perhaps, A^H), but the P&L explain analysis does away with Taylor expansions and instead relies on brute-force bumping of market data. As was the case for the P&L predict, the explain analysis is carried out at time $t+h$ when the market data movement δ is known.

22.2.2.1 Waterfall Explain

In one type of P&L explain — a so-called *waterfall explain* — the impact of the i -th component of δ is basically captured as

$$E_i = V(t+h; \Theta_{\text{mkt}}(t) + (\delta_1, \delta_2, \dots, \delta_i, 0, 0, \dots, 0)^\top) - V(t+h; \Theta_{\text{mkt}}(t) + (\delta_1, \delta_2, \dots, \delta_{i-1}, 0, 0, \dots, 0)^\top), \quad (22.24)$$

with $i = 1, \dots, N_{\text{mkt}}$. In other words, the impact of market variable i is recorded as the difference in portfolio values arising from moving the first $i-1$ and i market data variables, respectively, to their time $t+h$ values. The resulting attribution of P&L is often, quite descriptively, termed a “bump-and-do-not-reset” P&L explain.

⁹Also known by the more grammatically sensible names *P&L explanation* or *P&L attribution*.

Notice that

$$\sum_{i=1}^{N_{\text{mkt}}} E_i = V(t+h; \Theta_{\text{mkt}}(t+h)) - V(t+h; \Theta_{\text{mkt}}(t)) \neq V(t+h) - V(t),$$

since $V(t) = V(t; \Theta_{\text{mkt}}(t)) \neq V(t+h; \Theta_{\text{mkt}}(t))$. A complete P&L explain report must therefore add back a theta-type term that measures time decay, i.e. we write

$$V(t+h) - V(t) = \sum_{i=1}^{N_{\text{mkt}}} E_i + \{V(t+h; \Theta_{\text{mkt}}(t)) - V(t; \Theta_{\text{mkt}}(t))\}, \quad (22.25)$$

where the term in the curly brackets accounts for the effect associated with keeping market data fixed and letting time progress from t to $t+h$.

As argued in Section 22.2.1, the time decay definition used in (22.25) is often problematic. A more meaningful definition is given in (22.21), which leads to the following improved accounting for the P&L explain:

$$V(t+h) - V(t) = \sum_{i=1}^{N_{\text{mkt}}} E_i^f + \left\{ V\left(t+h; \Theta_{\text{mkt}}^f(t)\right) - V(t; \Theta_{\text{mkt}}(t)) \right\}, \quad (22.26)$$

where

$$\begin{aligned} E_i^f &= V\left(t+h; \Theta_{\text{mkt}}^f(t) + (\delta_1^f, \delta_2^f, \dots, \delta_i^f, 0, 0, \dots, 0)^T\right) \\ &\quad - V\left(t+h; \Theta_{\text{mkt}}^f(t) + (\delta_1^f, \delta_2^f, \dots, \delta_{i-1}^f, 0, 0, \dots, 0)^T\right), \end{aligned}$$

with δ^f defined in (22.23). Both (22.25) and (22.26) can be found in actual bank systems, but the latter is typically preferable.

22.2.2.2 Bump-and-Reset Explain

By construction, the waterfall P&L explain procedure in Section 22.2.2.1 is always fully able to explain P&L moves, in the sense that both (22.25) and (22.26) are identities, rather than approximations. While this is convenient, one drawback of the method is that the amount (E_i or E_i^f) of the P&L move that is allocated to an individual market data variable depends on how the vector $\Theta_{\text{mkt}}(t)$ happens to be ordered. This lends a certain amount of arbitrariness to the waterfall method, which sometimes can affect the P&L attribution process fairly substantially. To see this, assume that $\Theta_{\text{mkt}}(t)$ consists of interest rate and volatility data, and that interest rates are listed before the volatilities. Consider a position in an out-of-the-money caplet, and a market scenario where both interest rates and volatilities increase over the interval $[t, t+h]$. Further, assume that the shift in interest rates just happens to make the caplet position move from being out-of-the-money

(OTM) to at-the-money (ATM). In the waterfall P&L explain, since our ordering was such that we move interest rates before we move volatilities, when measuring the contribution from volatilities to the P&L move, we will register a decent amount, since ATM options have high vega. On the other hand, had we arbitrarily listed volatilities before interest rates in $\Theta_{\text{mkt}}(t)$, the contribution from volatility would have been computed on an OTM option with little vega, resulting in a much smaller P&L effect.

To avoid the consistency problems of the waterfall method, an alternative approach is to change (22.24) to

$$E_i = V(t + h; \Theta_{\text{mkt}}(t) + (0, 0, \dots, \delta_i, 0, 0, \dots, 0)^T) - V(t + h; \Theta_{\text{mkt}}(t)),$$

which is often called *bump-and-reset P&L explain*. With this definition, however, an exact P&L explain such as (22.25) is not possible, but we must instead content ourselves with an expression of the form

$$V(t + h) - V(t) = \sum_{i=1}^{N_{\text{mkt}}} E_i + \{V(t + h; \Theta_{\text{mkt}}(t)) - V(t; \Theta_{\text{mkt}}(t))\} + U, \quad (22.27)$$

where U represents the unexplained part of the P&L (the “unexplain”), primarily caused by cross-convexity terms in the Hessian matrix A^H , i.e. terms of the type $\partial^2 V / \partial \delta_i \partial \delta_j$, $i \neq j$. We note that (22.27) can be improved to incorporate the same notion of time decay as in (22.26); we leave this straightforward modification to the reader.

If the term U is consistently large, it may be necessary to explicitly add terms that capture cross-convexity exposure. This can be done using terms from A^H , which makes the overall procedure a bit of a hybrid between true P&L predict and explain. Alternatively, if we, say, identify the interaction of δ_i and δ_j as being considerable, we may do a joint bump-and-reset of these market data perturbations to split out a cross-term contribution of

$$\begin{aligned} & V(t + h; \Theta_{\text{mkt}}(t) + (0, 0, \dots, \delta_i, 0, \dots, \delta_j, 0, \dots, 0)^T) \\ & - V(t + h; \Theta_{\text{mkt}}(t) + (0, 0, \dots, \delta_i, 0, 0, \dots, 0)^T) \\ & - V(t + h; \Theta_{\text{mkt}}(t) + (0, 0, \dots, \delta_j, 0, 0, \dots, 0)^T). \end{aligned}$$

Carefully supplementing the basic bump-and-reset P&L explain with cross-term contributions will help ensure that the residual amount of unexplained P&L is small.

22.3 Value-at-Risk

While the P&L predict and explain are largely backward-looking accounting exercises, the risk management team in a bank is primarily focused on

analyzing the distribution of *future* portfolio values, in order to gauge the overall riskiness of a portfolio. Rather than report the entire P&L distribution, it is common to summarize it in a few summary statistics, known as *risk measures*. Many such risk measures exist, but the so-called *value-at-risk* (VaR) is probably the most commonly used in practice. VaR at level α (denoted Λ_α) is simply the $(1 - \alpha)$ -percentile of the distribution of the P&L move $V(t + h) - V(t)$ in the real-life measure P :

$$P(V(t + h) - V(t) \leq \Lambda_\alpha | \mathcal{F}_t) = 1 - \alpha. \quad (22.28)$$

In other words, the probability of losing¹⁰ more than $-\Lambda_\alpha$ over the time interval $[t, t + h]$ is less than $1 - \alpha$. Typically α is set to 99% or 95%, and h to one business day.

Another commonly used risk measure is *conditional value-at-risk* (cVaR), which is defined as the conditional expectation

$$\Xi_\alpha \triangleq E^P(V(t + h) - V(t) | V(t + h) - V(t) \leq \Lambda_\alpha). \quad (22.29)$$

cVaR has certain theoretical advantages to VaR¹¹, but VaR is nevertheless the more common in practice.

To compute Λ_α and Ξ_α , we need a statistical description for the market data increment vector δ . One popular choice uses the historical distribution of δ directly, giving rise to the so-called *historical VaR* risk measure. Here, one takes the actual realizations of δ over the last N_{VaR} trading days (e.g. $N_{\text{VaR}} = 500$, roughly corresponding to two years) and applies them to the current market data, thereby generating the empirical distribution of $V(t + h) - V(t)$. The calculation of VaR then amounts to ranking the impact of the last N_{VaR} market moves on the current portfolio, from worst to best, and using the impact of the market data move on the day with the rank $(1 - \alpha)N_{\text{VaR}}$ as the VaR.

Another VaR methodology uses a parameterized, rather than historical, distribution of market moves. As h is typically a short interval, it is, for instance, often justified to assume that each element in δ is Gaussian with zero mean and standard deviation s_i ,

$$\delta_i \sim \mathcal{N}(0, s_i^2), \quad i = 1, \dots, N_{\text{mkt}}, \quad (22.30)$$

where s_i may be estimated from the annualized basis point volatility σ_i of market element i through the relation

$$s_i = \sigma_i \sqrt{h}.$$

¹⁰Notice that Λ_α virtually always is a negative number. Sometimes it is the absolute value of this number that is reported as the VaR.

¹¹Specifically, VaR is not a *coherent risk measure*, in the sense defined in Artzner et al. [1999].

We capture co-dependence between the elements of δ in a correlation matrix R , typically estimated from historical time series. Notice that even if market data element i has some non-zero drift, the mean of δ_i would be of order $O(h)$ and typically negligible relative to s_i (order $O(\sqrt{h})$); the assumption of zero mean is therefore an innocuous one.

To compute VaR and cVaR in the Gaussian setup, one option is to perform a brute-force simulation of the portfolio value $V(t + h)$, using a full portfolio revaluation for each simulated value of the vector δ . While this is, in fact, sometimes done, it is far easier to rely on the Taylor expansion (22.19). For VaR/cVaR purposes, we may safely ignore the time decay term in (22.19), hence we can take as our starting point the equation

$$V(t + h) = V(t) + \nabla^H(t) \cdot \delta + \frac{1}{2} \delta^\top \cdot A^H(t) \cdot \delta. \quad (22.31)$$

With (22.31) and (22.30), the VaR and cVaR computations are analytically tractable. One simple result is the following.

Proposition 22.3.1. *Let $V(t + h)$ be given as in (22.31), with $A^H(t) = 0$. Also, let the elements of the N_{mkt} -dimensional vector δ have correlation matrix R and satisfy (22.30). Setting $s = (s_1, \dots, s_{N_{\text{mkt}}})^\top$, we have*

$$\Lambda_\alpha = v\Phi^{-1}(1 - \alpha), \quad (22.32)$$

$$\Xi_\alpha = -v(1 - \alpha)^{-1}\phi(\Phi^{-1}(1 - \alpha)), \quad (22.33)$$

where $\phi(x) = (2\pi)^{-1/2}e^{-x^2/2}$ is the Gaussian density, and

$$v^2 = \nabla^H(t) \text{diag}(s) R \text{diag}(s) \nabla^H(t)^\top.$$

Proof. First observe that the covariance matrix C of δ is

$$C = \text{diag}(s) R \text{diag}(s),$$

where $\text{diag}(s)$ is a square matrix with s along the diagonal and zeros elsewhere. Under our assumptions, it is clear that $V(t + h) \sim \mathcal{N}(V(t), v^2)$, where the variance v^2 is given by

$$v^2 = \nabla^H(t) C \nabla^H(t)^\top.$$

Defining $\Delta V = V(t + h) - V(t)$ and writing $\Delta V = vZ$ for $Z \sim \mathcal{N}(0, 1)$, we have

$$P(\Delta V \leq \Lambda_\alpha) = P\left(Z \leq \frac{\Lambda_\alpha}{v}\right) = \Phi\left(\frac{\Lambda_\alpha}{v}\right).$$

Equating this expression to $1 - \alpha$, per (22.28), results in (22.32). To compute the cVaR, we write

$$\begin{aligned} \mathbb{E}^P(\Delta V | \Delta V \leq \Lambda_\alpha) &= P(\Delta V \leq \Lambda_\alpha)^{-1} \times \mathbb{E}^P(\Delta V 1_{\{\Delta V \leq \Lambda_\alpha\}}) \\ &= (1 - \alpha)^{-1} \mathbb{E}^P(v Z 1_{\{Z \leq m\}}) \\ &= (1 - \alpha)^{-1} v \mathbb{E}^P(Z 1_{\{Z \leq m\}}), \end{aligned}$$

where $m = \Lambda_\alpha/v = \Phi^{-1}(1 - \alpha)$. Observe that

$$\mathbb{E}^P(Z 1_{\{Z \leq m\}}) = \int_{-\infty}^m z \phi(z) dz = -\phi(m),$$

and (22.33) follows. \square

The results in Proposition 22.3.1 are often denoted *delta VaR* and *delta cVaR*, respectively, to reflect the fact that we have ignored the Hessian matrix A^H . If we wish to include A^{II} — to compute what is known as *delta-gamma VaR/cVaR* — matters get a bit more complicated, as the distribution of $V(t+h) - V(t)$ is no longer simple. One method, described in Rouvinez [1997], shows that $V(t+h) - V(t)$ can be expressed as a sum of independent non-central chi-square random variables. From this representation, the characteristic function of $V(t+h) - V(t)$ can be constructed and, using a numerical technique, turned into a cumulative distribution function from which VaR and cVaR can be computed. As the topic is somewhat tangential to our needs in this book, we omit the details here but just note that at the end of the day the key to a good delta-gamma VaR/cVaR computation is a reliable and accurate estimate for ∇^H and A^H , a topic that shall occupy the remainder of this book.

22.A Appendix: Alternative Proof of Lemma 22.1.1

Consider a contingent claim with terminal value $g(X(T))$, where $X(t)$ satisfies the Black-Scholes SDE (22.5). Let us write the time t value of this claim as $V(t) = h(T-t, X(t))$, where h satisfies the PDE

$$\frac{\partial h}{\partial \tau} = \mathcal{L}V, \quad \mathcal{L} = -r + rX \frac{\partial}{\partial X} + \frac{1}{2}\sigma^2 X^2 \frac{\partial^2}{\partial X^2}, \quad (22.34)$$

subject to an initial condition $h(0, X) = g(X)$. Operator calculus treats this equation as an ordinary differential equation in $\tau = T-t$. The solution to (22.34) is then given by

$$h(\tau, X) = \exp(\tau \mathcal{L}) g(X), \quad (22.35)$$

where the exponential must be interpreted as

$$\exp(\tau \mathcal{L}) = \sum_{i=0}^{\infty} \frac{(\tau \mathcal{L})^i}{i!}.$$

Differentiating with respect to τ verifies that (22.35) is, indeed, the solution to the initial value problem (22.34).

To form the derivative $\partial V / \partial \sigma = \partial h / \partial \sigma$, we notice that all dependency on σ in the expression (22.35) is in the operator \mathcal{L} . Differentiating with respect to σ , we get

$$\frac{\partial h}{\partial \sigma} = \tau \frac{\partial \mathcal{L}}{\partial \sigma} \exp(\tau \mathcal{L}) g(X) = \tau \sigma X^2 \frac{\partial^2}{\partial X^2} (\exp(\tau \mathcal{L}) g(X)) = \tau \sigma X^2 \frac{\partial^2 h}{\partial X^2},$$

or, equivalently,

$$\frac{\partial V}{\partial \sigma} = (T - t) \sigma X^2 \frac{\partial^2 V}{\partial X^2},$$

which is (22.7).

The operator representation (22.35) makes it easy to compute other parameter derivatives. For instance, we note that

$$\frac{\partial h}{\partial r} = \tau \left(-1 + X \frac{\partial}{\partial X} \right) (\exp(\tau \mathcal{L}) g(X)) = \tau \left(X \frac{\partial h}{\partial X} - g \right),$$

such that

$$\frac{\partial V}{\partial r} = (T - t) \left(X \frac{\partial V}{\partial X} - V \right),$$

as mentioned in Remark 22.1.2.

Payoff Smoothing and Related Methods

As made clear in the previous chapter, practical risk management of a portfolio of interest rate securities revolves around price sensitivities with respect to various valuation inputs, such as market prices and model parameters. These sensitivities are often¹ calculated by applying small perturbations to market and model parameters, followed by a re-pricing of the securities portfolio in question.

Being derivatives of a model price function, price sensitivities (greeks) are inherently less smooth than the price function itself. For instance, it is well-known (see Section 1.10.3) that while the value of a Bermudan security on an exercise date is continuous across the exercise boundary, its delta is not. This lack of smoothness will often put significant stress on a numerical scheme, which effectively is faced with the problem of resolving an irregular boundary condition. As a result, a careless implementation of a numerical sensitivity computation will often produce poor results, with the resulting greeks being less stable than what is expected theoretically. In this chapter we study this problem, with an emphasis on how to adapt numerical schemes to avoid introducing spurious instabilities into the calculation of greeks. Some of the discussion in this chapter builds on previous material, and we suggest that the reader briefly review Sections 2.5 and 3.3 before proceeding.

23.1 Issues with Discretization Schemes

As we saw in Section 3.3.1.2, fixing the random seed when computing deltas (and other greeks) by Monte Carlo methods significantly reduces the standard error. This is a simple example of the general rule that one should ideally attempt to freeze as many aspects of a numerical calculation (such as a Monte Carlo seed, the geometry of a PDE grid, etc.) as possible when doing perturbation analysis in a numerical scheme. In particular, adhering to this

¹But not always — see Chapter 24.

simple rule often ensures that no additional discontinuities are introduced by the numerical method itself.

23.1.1 Problems with Grid Dimensioning

While straightforward in theory, consistently following the rule above can be quite difficult. Sometimes violations are subtle and unintentional, and sometimes they are unavoidable due to systems limitations or computational precision requirements. Let us look at the former case, using as our first example the problem of valuing a simple option by numerically integrating the payout against a Gaussian density. As is frequently done in practice, suppose that the numerical domain of integration is chosen to be $[-5\sigma\sqrt{T}, 5\sigma\sqrt{T}]$, where T is the expiry of the option and σ is the asset volatility (a model input). In addition, let the number of integration nodes, N , be chosen so that the integration grid is uniform with a pre-specified and fixed step δ , i.e.

$$N = \left\lfloor \frac{10\sigma\sqrt{T}}{\delta} \right\rfloor, \quad (23.1)$$

where $\lfloor \cdot \rfloor$ denotes the integer part of a real number. The resulting option valuation scheme appears reasonable, if slightly unconventional. However, imagine now that we wish to compute the option vega by comparing the base value of the option to a value computed after shocking the volatility to a new level of $\sigma + \Delta\sigma$. Since the integration domain depends on σ , it will be slightly larger in the perturbation scenario. Moreover, as the integration step is kept constant, the number of integration nodes may change between the base and the bumped scenarios. As the number of integration points can only move by an integer amount, the change in the grid geometry would not be continuous, introducing a purely artificial contribution of the order $O((\Delta\sigma)^{-1})$. This contribution explodes to infinity as $\Delta\sigma \rightarrow 0$ (as long as the number of steps changes as a result of the volatility perturbation).

The issue that arises in the example above stems from the fact that the number of integration nodes N as given by (23.1) is *not a smooth function* of σ , as the function $x \mapsto \lfloor x \rfloor$ is not differentiable (or even continuous). Since the numerical value of the security is a function of the number of integration nodes N , the value will not be smooth with respect to σ , and the vega, while continuous in the theoretical model, will be discontinuous in the numerical scheme. Of course, the problem is easy to rectify: heed our advice and avoid altering the geometry of the grid when perturbing market inputs.

23.1.2 Grid Shifts Relative to Payout

The example in Section 23.1.1 above is an example of grid geometry changing outright, due to changes in asset moments used for grid position and/or dimensioning. Another problematic case occurs when the grid is frozen in

space, but the nature of the perturbation itself will cause an *effective* shift of the grid relative to the payout. To give an example of this, consider a problem of valuing a European option with a payoff $f(x)$ on an underlying $S(T)$ observed at time T . Assuming zero interest rates, the value of this option is equal to the integral of the payoff against the PDF of the underlying. Let the initial (time 0) spot value of the underlying be denoted by S . Assuming for simplicity that the distribution of the increment $S(T) - S$ is independent of S , we denote the density of $S(T) - S$ by

$$\mathbb{Q}((S(T) - S) \in dx) = \pi(x) dx.$$

The value of the option as a function of the spot S today is then given by two equivalent expressions

$$V(S) = \int_{-\infty}^{\infty} f(x)\pi(x - S) dx \quad (23.2)$$

$$= \int_{-\infty}^{\infty} f(S + x)\pi(x) dx. \quad (23.3)$$

While (23.2) and (23.3) are mathematically equivalent, (23.2) is better suited for numerical computations of sensitivities of V with respect to S . In particular, notice that changes in S here get absorbed into the *density* of S , which in a grid setting simply amounts to changing the weights on individual grid points in a numerical quadrature rule. On the other hand, the formulation (23.3) absorbs changes in S into the *payout*, which effectively causes the grid to move relative to the payout function².

To demonstrate the kind of problems that arise when the discretization grid is not fixed relative to the payoff, let us consider a simple example. We recall that a typical non-adaptive quadrature scheme (including rectangle, trapezoidal and Gaussian quadrature rules) specifies a collection of fixed knot points $\{x_n\}_{n=0}^N$ and weights $\{w_n\}_{n=0}^N$, and approximates $V(S) \approx \tilde{V}(S)$, where

$$\tilde{V}(S) = \sum_{n=0}^N w_n f(x_n + S).$$

Note again that the weights and knots are fixed relative to the density of the process and not the payoff, i.e. contrary to our earlier advice. Let us analyze the behavior of the scheme under shifts of S . For concreteness, we consider a European call option, i.e.

$$f(x) = (x - K)^+$$

for a fixed choice of K . Then

²The observant reader may have noticed a strong connection to material in Chapter 3 on pathwise and likelihood ratio methods for Monte Carlo applications.

$$\tilde{V}(S) = \sum_{n=0}^N w_n (x_n + S - K)^+.$$

The exact derivative of the numerical value $\tilde{V}(S)$ with respect to the initial value of the asset, i.e. the delta, is given by

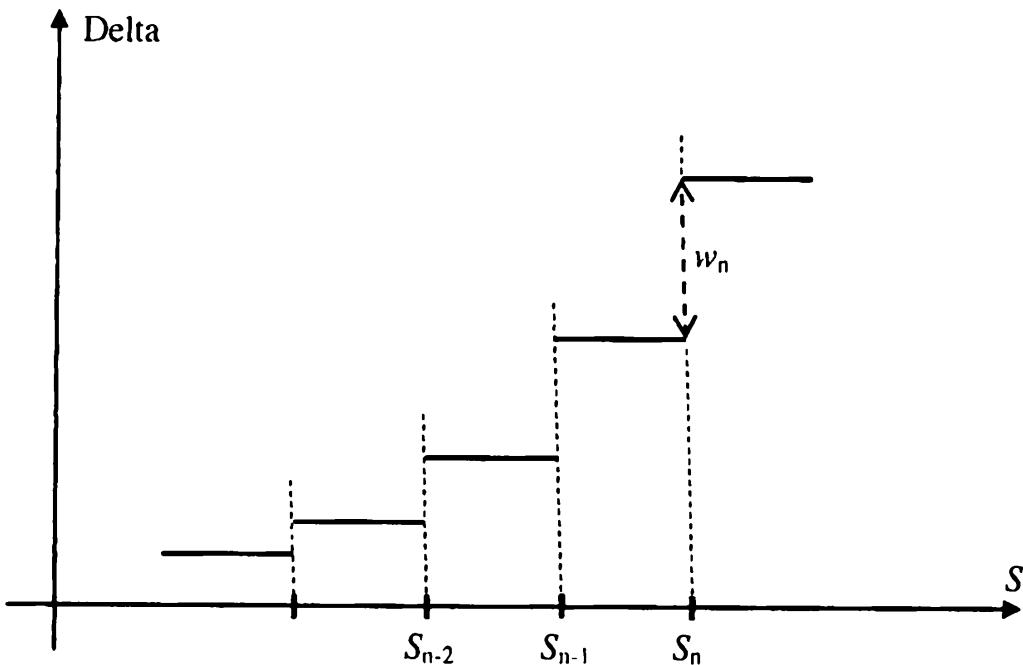
$$\frac{d}{dS} \tilde{V}(S) = \sum_{n=0}^N w_n 1_{\{x_n + S - K > 0\}}.$$

As a function of S , this function has discontinuities of sizes w_n at points $S = S_n$,

$$S_n \triangleq K - x_n,$$

for all $n = 0, \dots, N$. Thus, as the spot S moves, the delta will jump whenever the spot crosses one of the levels S_n . Moreover, in this scheme, the delta does not change as long as S does not cross one of S_n , which is obviously unrealistic. A typical plot of such a “delta” is shown in Figure 23.1.

Fig. 23.1. Discontinuous Delta



Notes: Delta of a derivative evaluated with an integration scheme that is not fixed relative to the payoff.

The irregular deltas in Figure 23.1 are caused by the call option kink crossing over knot points in the grid as a result of parameter perturbation; similar behavior will occur for all payouts with payout discontinuities, kinks, and the like.

23.1.3 Additional Comments

Problems of the types described in Sections 23.1.1 and 23.1.2 are easy to introduce, but often difficult to track down. This is particularly true for complex models that often use sophisticated numerical methods with complicated dependencies on market and model data. When examining numerical routines for problems, we note that one must obviously include both valuation and calibration algorithms, as the computation of stable sensitivities requires that both types of algorithms behave smoothly with respect to moves in market parameters. Local, bootstrap-type calibrations (such as that developed in Section 13.1.7 for quasi-Gaussian models) generally outperform global, best-fit calibrations (such as that from Section 14.5.7 for LM models). We postpone our treatment of calibration effects to Chapter 26, and in this chapter focus on the (post-calibration) problem of building numerically smooth valuation routines.

Even if one is vigilant and tracks down all cases of non-constant (effective) grid geometry, there may, as mentioned earlier, exist violations that are impractical to resolve. For example, the Monte Carlo “grid” is difficult to keep constant by the nature of how it is generated. Also, by not explicitly tailoring a finite difference grid geometry to the market data used in a perturbed scenario, an unacceptable loss of precision may occur. Finally, we note that the organizational setup in many banks may make it hard for valuation code developers to fully control risk sensitivity computations. For all these reasons, it is worthwhile to develop methods that will produce stable sensitivities, *even if* the grid geometry cannot be guaranteed to stay fixed under market data moves.

As a final comment, let us briefly note that some numerical methods, in principle, give us greeks “for free” as part of valuation. For example, in a finite difference grid solution of a vanilla model PDE, the derivative with respect to the asset can be read off the PDE grid by forming a central finite difference of the solution (at time 0) at grid nodes surrounding today’s value of the asset³. This avenue is not available for all greeks, however, and even for deltas and gammas the utility of the method is quite limited in term structure models, since we rarely are interested in theoretical derivatives with respect to the abstract model variables, but instead nearly always wish to compute sensitivities to specific perturbations of the yield curve. In addition, as described in Section 16.1.1, we may often be interested in working with joint moves of interest rates and volatilities that are user-prescribed and incompatible with the theoretical model dynamics. Even for PDE-based

³Sometimes one gets better numerical properties if a spline is fit through all grid values at $t = 0$; deltas and gammas can then be computed by differentiating this spline.

models we therefore typically will need to calculate greeks by applying market data bumps⁴.

23.2 Basic Techniques

An obvious remedy to many of the issues identified in earlier sections is simply to increase the number of grid points. As a rule of thumb, the discretization step of a grid should be significantly smaller than a typical shift applied to inputs when calculating greeks. For instance, in Figure 23.1, if the perturbation size for delta calculation covered a few grid intervals, the delta would vary fairly smoothly with the value of the spot. Of course, the smaller the grid size, the more computationally intensive the numerical scheme will be, and we cannot always trade speed for accuracy with impunity, since greek calculations are often time-sensitive. The next few sections cover several more sophisticated, and less computationally costly, alternatives to brute-force grid refinement.

23.2.1 Adaptive Integration

Increasing the density of grid points uniformly is not the best way to spend a computational budget, as adding extra points away from a discontinuity of a payoff does little to improve the numerical properties of greeks. A reasonable yet relatively simple way of improving numerical stability of an algorithm at a modest computational cost is to first identify the region of the state space where the payoff is likely to have singularities⁵, and then sample this part of the space at a higher resolution than elsewhere. For instance, suppose we know that the function $f(x)$ that we are integrating has a singularity in a particular interval $[x_m, x_{m+1}]$ of the integration grid. Then we can further subdivide $[x_m, x_{m+1}]$ into 10 (or 100) subintervals and apply the trapezoidal rule for the finer grid. This will insure that the singularity is handled accurately, while no extra effort is wasted in the regions where f is smooth.

For any particular payoff function, it is often relatively simple to identify the regions where a denser grid is beneficial — the type of knowledge that can be incorporated directly into the integration routine. For a more generic setup (or for payouts that are hard to analyze directly), a good

⁴In fact, several commonly used risk measures are explicitly meant to be computed by finite-sized shifts. For instance, many swaption traders' definition of gamma is the change in delta, for a 10 basis point move in the yield curve.

⁵Here and elsewhere in this chapter, the term “singularity” refers to a point in the state space at which the payoff function, or one of its derivatives, is discontinuous. Examples include the barrier for a digital option and the strike of a European put or call.

alternative is to rely on adaptive integration routines, often prepackaged in numerical libraries such as IMSL and NAG. In this class of routines, the integral is approximated using ordinary quadrature rules on adaptively refined subintervals of the integration domain until a stopping criterion is met. In effect, the grid points are chosen automatically based directly on properties of the function being integrated. Adaptive algorithms generally work quite well, but care must be taken to keep grid geometry fixed in perturbed scenarios, something that can occasionally be a bit of a challenge if a third-party library routine is used.

23.2.2 Adding Singularities to the Grid

In the situation where perturbations cause a grid shift relative to the payoff, if the position of a payoff singularity is known exactly, numerical properties of the valuation algorithm can be improved substantially by simply adding the singularity to the integration grid. This serves to effectively lock down the grid geometry in the immediate vicinity of the singularity. The method is closely related to the *grid shifting* method used to improve convergence of numerical solutions of PDEs, see Section 2.5.3.

Using the integration problem (23.3) as an example, let us consider a continuous payoff $f(y)$ whose derivative is discontinuous at a single point $y = K$ (multiple singularities can be handled similarly). Let us rewrite the value of the option as

$$V(S) = \int_{-\infty}^{K-S} f(x + S) \pi(x) dx + \int_{K-S}^{\infty} f(x + S) \pi(x) dx.$$

Suppose we proceed to apply a simple numerical quadrature to each integral separately. We start with $N + 1$ integration knots fixed in x -space, and add one more knot at the singularity, i.e. at $x = K - S$. To characterize the location of the additional knot, let the index $\mu(S)$ be defined by

$$x_{\mu(S)} \leq K - S < x_{\mu(S)+1}.$$

Using a trapezoidal integration rule for simplicity, the resulting numerical scheme can formally be written as

$$\begin{aligned} \tilde{V}(S) &= \tilde{V}_1(S) + \tilde{V}_2(S) + \tilde{V}_3(S) + \tilde{V}_4(S), \\ \tilde{V}_1(S) &= \frac{1}{2} \sum_{n=1}^{\mu(S)} (f(x_n + S) \pi(x_n) + f(x_{n-1} + S) \pi(x_{n-1})) \Delta x_n, \\ \tilde{V}_2(S) &= \frac{1}{2} (f(K) \pi(K - S) + f(x_{\mu(S)} + S) \pi(x_{\mu(S)})) (K - S - x_{\mu(S)}), \end{aligned}$$

$$\begin{aligned}\tilde{V}_3(S) &= \frac{1}{2} (f(x_{\mu(S)+1} + S) \pi(x_{\mu(S)+1}) + f(K) \pi(K - S)) \\ &\quad \times (x_{\mu(S)+1} - (K - S)), \\ \tilde{V}_4(S) &= \frac{1}{2} \sum_{n=\mu(S)+2}^N (f(x_n + S) \pi(x_n) + f(x_{n-1} + S) \pi(x_{n-1})) \Delta x_n.\end{aligned}$$

The terms $\tilde{V}_1(S)$ and $\tilde{V}_4(S)$ collect the contributions of integration intervals before and after the singularity, respectively, whereas the terms $\tilde{V}_2(S)$ and $\tilde{V}_3(S)$ represent the contributions of the integration interval containing the singularity.

Let us fix m such that $\mu(S) = m$ for the initial (pre-perturbed) value of S . Clearly there are no issues with the smoothness of $d\tilde{V}(S)/dS$ as we move S in such a way that $K - S \in [x_m, x_{m+1})$, so the only potential discontinuity could arise when $K - S$ crosses one of the integration nodes, i.e. when $\mu(S)$ jumps. To show that the scheme implies smooth behavior across grid points, let us investigate what happens when S crosses $S_m \triangleq K - x_m$. For this analysis, we only need to keep track of the terms $\tilde{U}(S)$ that originate from the intervals adjacent to x_m . For $x_{m+1} > K - S \geq x_m$ we have

$$\begin{aligned}\tilde{U}(S) &= \frac{1}{2} (f(x_m + S) \pi(x_m) + f(x_{m-1} + S) \pi(x_{m-1})) \Delta x_m \\ &\quad + \frac{1}{2} (f(K) \pi(K - S) + f(x_m + S) \pi(x_m)) (K - S - x_m) \\ &\quad + \frac{1}{2} (f(x_{m+1} + S) \pi(x_{m+1}) + f(K) \pi(K - S)) (x_{m+1} - (K - S)),\end{aligned}$$

and for shifted S such that $x_{m-1} \leq K - S < x_m$, i.e. for $\mu(S) = m - 1$,

$$\begin{aligned}\tilde{U}(S) &= \frac{1}{2} (f(K) \pi(K - S) + f(x_{m-1} + S) \pi(x_{m-1})) (K - S - x_{m-1}) \\ &\quad + \frac{1}{2} (f(x_m + S) \pi(x_m) + f(K) \pi(K - S)) (x_m - (K - S)) \\ &\quad + \frac{1}{2} (f(x_{m+1} + S) \pi(x_{m+1}) + f(x_m + S) \pi(x_m)) \Delta x_{m+1}.\end{aligned}$$

Note that $\tilde{U}(S)$ is continuous at $S = S_m$ and

$$\begin{aligned}\tilde{U}(S_m) &= \frac{1}{2} (f(K) \pi(x_m) + f(K - \Delta x_m) \pi(x_{m-1})) \Delta x_m \\ &\quad + \frac{1}{2} (f(K + \Delta x_{m+1}) \pi(x_{m+1}) + f(K) \pi(x_m)) \Delta x_{m+1}.\end{aligned}$$

In fact, the *derivative* of $\tilde{U}(S)$ is also continuous across the grid point, as we prove in Appendix 23.A.

As a final observation, we note that to add the singularity to the grid, the location of the singularity obviously needs to be detected first. In many

alternative is to rely on adaptive integration routines, often prepackaged in numerical libraries such as IMSL and NAG. In this class of routines, the integral is approximated using ordinary quadrature rules on adaptively refined subintervals of the integration domain until a stopping criterion is met. In effect, the grid points are chosen automatically based directly on properties of the function being integrated. Adaptive algorithms generally work quite well, but care must be taken to keep grid geometry fixed in perturbed scenarios, something that can occasionally be a bit of a challenge if a third-party library routine is used.

23.2.2 Adding Singularities to the Grid

In the situation where perturbations cause a grid shift relative to the payoff, if the position of a payoff singularity is known exactly, numerical properties of the valuation algorithm can be improved substantially by simply adding the singularity to the integration grid. This serves to effectively lock down the grid geometry in the immediate vicinity of the singularity. The method is closely related to the *grid shifting* method used to improve convergence of numerical solutions of PDEs, see Section 2.5.3.

Using the integration problem (23.3) as an example, let us consider a continuous payoff $f(y)$ whose derivative is discontinuous at a single point $y = K$ (multiple singularities can be handled similarly). Let us rewrite the value of the option as

$$V(S) = \int_{-\infty}^{K-S} f(x + S) \pi(x) dx + \int_{K-S}^{\infty} f(x + S) \pi(x) dx.$$

Suppose we proceed to apply a simple numerical quadrature to each integral separately. We start with $N + 1$ integration knots fixed in x -space, and add one more knot at the singularity, i.e. at $x = K - S$. To characterize the location of the additional knot, let the index $\mu(S)$ be defined by

$$x_{\mu(S)} \leq K - S < x_{\mu(S)+1}.$$

Using a trapezoidal integration rule for simplicity, the resulting numerical scheme can formally be written as

$$\begin{aligned} \tilde{V}(S) &= \tilde{V}_1(S) + \tilde{V}_2(S) + \tilde{V}_3(S) + \tilde{V}_4(S), \\ \tilde{V}_1(S) &= \frac{1}{2} \sum_{n=1}^{\mu(S)} (f(x_n + S) \pi(x_n) + f(x_{n-1} + S) \pi(x_{n-1})) \Delta x_n, \\ \tilde{V}_2(S) &= \frac{1}{2} (f(K) \pi(K - S) + f(x_{\mu(S)} + S) \pi(x_{\mu(S)})) (K - S - x_{\mu(S)}), \end{aligned}$$

$$\begin{aligned}\tilde{V}_3(S) &= \frac{1}{2} (f(x_{\mu(S)+1} + S) \pi(x_{\mu(S)+1}) + f(K) \pi(K - S)) \\ &\quad \times (x_{\mu(S)+1} - (K - S)),\end{aligned}$$

$$\tilde{V}_4(S) = \frac{1}{2} \sum_{n=\mu(S)+2}^N (f(x_n + S) \pi(x_n) + f(x_{n-1} + S) \pi(x_{n-1})) \Delta x_n.$$

The terms $\tilde{V}_1(S)$ and $\tilde{V}_4(S)$ collect the contributions of integration intervals before and after the singularity, respectively, whereas the terms $\tilde{V}_2(S)$ and $\tilde{V}_3(S)$ represent the contributions of the integration interval containing the singularity.

Let us fix m such that $\mu(S) = m$ for the initial (pre-perturbed) value of S . Clearly there are no issues with the smoothness of $d\tilde{V}(S)/dS$ as we move S in such a way that $K - S \in [x_m, x_{m+1})$, so the only potential discontinuity could arise when $K - S$ crosses one of the integration nodes, i.e. when $\mu(S)$ jumps. To show that the scheme implies smooth behavior across grid points, let us investigate what happens when S crosses $S_m \triangleq K - x_m$. For this analysis, we only need to keep track of the terms $\tilde{U}(S)$ that originate from the intervals adjacent to x_m . For $x_{m+1} > K - S \geq x_m$ we have

$$\begin{aligned}\tilde{U}(S) &= \frac{1}{2} (f(x_m + S) \pi(x_m) + f(x_{m-1} + S) \pi(x_{m-1})) \Delta x_m \\ &\quad + \frac{1}{2} (f(K) \pi(K - S) + f(x_m + S) \pi(x_m)) (K - S - x_m) \\ &\quad + \frac{1}{2} (f(x_{m+1} + S) \pi(x_{m+1}) + f(K) \pi(K - S)) (x_{m+1} - (K - S)),\end{aligned}$$

and for shifted S such that $x_{m-1} \leq K - S < x_m$, i.e. for $\mu(S) = m - 1$,

$$\begin{aligned}\tilde{U}(S) &= \frac{1}{2} (f(K) \pi(K - S) + f(x_{m-1} + S) \pi(x_{m-1})) (K - S - x_{m-1}) \\ &\quad + \frac{1}{2} (f(x_m + S) \pi(x_m) + f(K) \pi(K - S)) (x_m - (K - S)) \\ &\quad + \frac{1}{2} (f(x_{m+1} + S) \pi(x_{m+1}) + f(x_m + S) \pi(x_m)) \Delta x_{m+1}.\end{aligned}$$

Note that $\tilde{U}(S)$ is continuous at $S = S_m$ and

$$\begin{aligned}\tilde{U}(S_m) &= \frac{1}{2} (f(K) \pi(x_m) + f(K - \Delta x_m) \pi(x_{m-1})) \Delta x_m \\ &\quad + \frac{1}{2} (f(K + \Delta x_{m+1}) \pi(x_{m+1}) + f(K) \pi(x_m)) \Delta x_{m+1}.\end{aligned}$$

In fact, the *derivative* of $\tilde{U}(S)$ is also continuous across the grid point, as we prove in Appendix 23.A.

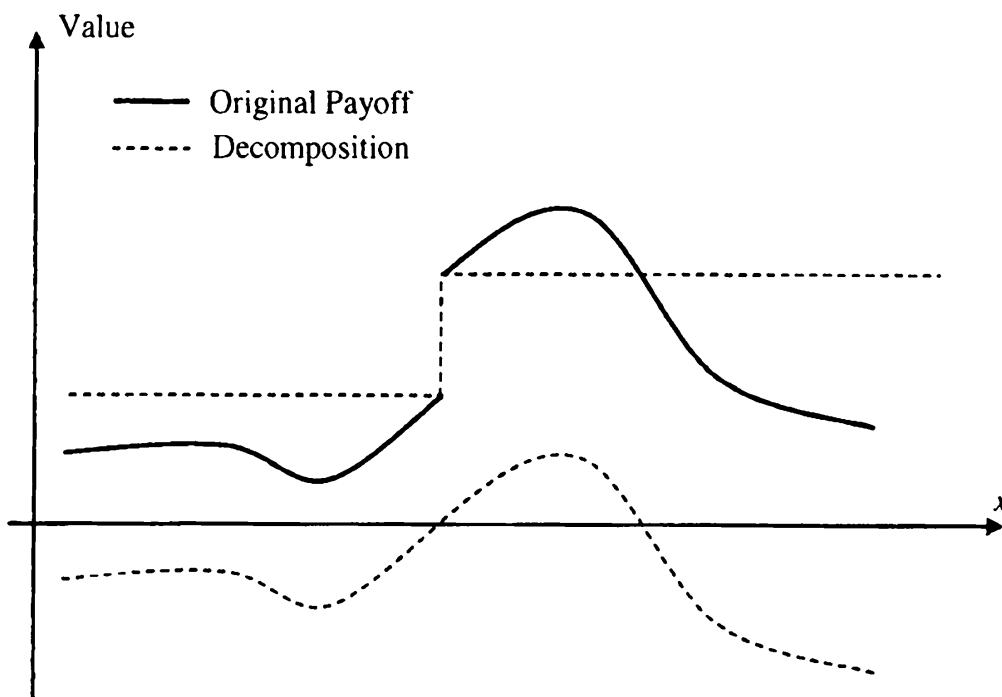
As a final observation, we note that to add the singularity to the grid, the location of the singularity obviously needs to be detected first. In many

cases, this must be done numerically using, for example, the method of Section 23.3.2.1. The numerical improvements to the greeks, however, are usually well worth the extra cost.

23.2.3 Singularity Removal

Most of the noise in greeks comes from the fact that numerical schemes have difficulty handling payoff singularities. It follows that removing these singularities should restore smoothness. The method based on this idea is quite powerful when it works, but is somewhat limited in its scope.

Fig. 23.2. Singularity Removal



Notes: A discontinuous payoff could be decomposed into a continuous one and a step function.

Suppose that, as in Figure 23.2, an otherwise-smooth payoff $f(x)$ has a jump discontinuity at one point $x = K$, with the size of the jump equal to a . The payoff can evidently then be decomposed into the sum of two functions, one being a simple step function $a1_{\{x>K\}}$, and the other equal to $g(x) = f(x) - a1_{\{x>K\}}$. Notice that the function $g(x)$ is smooth, unlike $f(x)$ itself. The integration problem

$$V(S) = \int_{-\infty}^{\infty} f(S + x)\pi(x) dx$$

may now be split in two,

$$V(S) = a \int_{-\infty}^{\infty} 1_{\{x>K\}} \pi(x - S) dx + \int_{-\infty}^{\infty} g(x) \pi(x - S) dx.$$

Let $\Psi(x)$ be the cumulative distribution function corresponding to the density $\pi(x)$, then

$$V(S) = a(1 - \Psi(K - S)) + \int_{-\infty}^{\infty} g(x) \pi(x - S) dx.$$

Suppose the CDF $\Psi(x)$ can be computed analytically, or numerically with high precision (and, say, tabulated). As the only function being integrated numerically, g , is smooth by construction, this scheme produces smooth greeks.

Shifted and scaled step functions $a1_{\{x>K\}}$, and combinations thereof, can be used to remove outright discontinuities in the payoff f . To remove discontinuities in the derivative of f , linear combinations of functions $(ax+b)^+$ for various a, b (i.e., call/put type payoffs) can be used instead. It should be clear, however, that the applicability of the singularity removal method will be limited by the availability of accurate (ideally analytic) methods to compute CDF and call/put option values.

23.2.4 Partial Analytical Integration

For the cases where the CDF or call/put option values are not known analytically, suitable approximations could still be used for smoothing the payoff. In the common case where the density $\pi(x)$ corresponds to a diffusive random variable $S(T)$, one approach is to focus on short times to maturity, where the conditional transition density

$$\pi(x, T; y, T - \delta) dx \triangleq Q(S(T) \in dx | S(T - \delta) = y)$$

can often be approximated by, say, a Gaussian or log-normal distribution for small $\delta > 0$. Then the value of a derivative with the payoff $f(x)$ is given by

$$E(f(S(T))) = E(E_{T-\delta}(f(S(T)))) = \int_{-\infty}^{\infty} \pi(x, T - \delta; S, 0) V(x, T - \delta) dx, \quad (23.4)$$

where $V(x, T - \delta)$ is the value of the derivative at time $T - \delta$,

$$V(x, T - \delta) = \int_{-\infty}^{\infty} \pi(y, T; x, T - \delta) f(y) dy.$$

With the approximation to $\pi(y, T; x, T - \delta)$ by a Gaussian or log-normal density, this integral can often be calculated analytically. This is certainly true for puts/calls and digitals. More complex payoffs can, for each value of x , often be approximated as combinations of simple payoffs in the vicinity

of x , as the width of distribution of $S(T)$ given $S(T - \delta)$ is small for small δ . Moreover, $V(x, T - \delta)$ is often a (much) smoother function of x than $f(x)$; for example the Black-Scholes value of a call option is infinitely differentiable, unlike the payoff itself. Hence, a numerical integration scheme applied to (23.4) should result in smoother and more stable greeks.

The method above (as well as several others reviewed in Section 23.2) is not limited to numerical integration, but can equally well be applied for PDE and Monte Carlo valuation; in fact we already saw a PDE application in Section 2.8. For example, a non-trivial application to TARNs in Monte Carlo is presented in Pietersz and van Regenmortel [2006]. To briefly review this method, let us recall the setup of Section 20.1 and note that

$$V_{\text{tar}_n}(0) = \sum_{n=1}^{N-1} V_{\text{cpn},n}(0),$$

where $V_{\text{cpn},n}(0)$ is the value of the n -th (net) coupon conditional on no early redemption,

$$V_{\text{cpn},n}(0) = E(B(T_{n+1})^{-1} X_n 1_{\{Q_n < R\}}).$$

Here Q_n is the sum of all structured coupons paid on or before T_n , i.e.

$$Q_n = Q_{n-1} + \tau_{n-1} C_{n-1}.$$

Observing that Q_{n-1} is $\mathcal{F}_{T_{n-2}}$ -measurable allows us to write

$$V_{\text{cpn},n}(0) = E(B(T_{n-2})^{-1} V_{\text{cpn},n}(T_{n-2})), \quad (23.5)$$

where

$$V_{\text{cpn},n}(T_{n-2}) = P(T_{n-2}, T_{n+1}) E_{T_{n-2}}^{T_{n+1}} (X_n 1_{\{C_{n-1} < (R - Q_{n-1})/\tau_{n-1}\}}).$$

For Libor-based TARNs, X_n is a function of $L_n(T_n)$ and C_{n-1} is a function of $L_{n-1}(T_{n-1})$, so the calculation of $V_{\text{cpn},n}(T_{n-2})$ involves an integral of a discontinuous function over a joint distribution of $(L_{n-1}(T_{n-1}), L_n(T_n))$, conditioned on $\mathcal{F}_{T_{n-2}}$. As coupon periods are rarely longer than a year, a Gaussian (or log-normal) approximation to this distribution is often accurate enough. Drifts and (co-)variances of Libor rates $(L_{n-1}(T_{n-1}), L_n(T_n))$ can typically be estimated with relative ease from the term structure model used for valuation (see Section 20.1.3 for a relevant discussion), at which point $V_{\text{cpn},n}(T_{n-2})$ would be calculated by an exact or approximate quadrature, perhaps aided by the various methods from Section 17.6. Calculating $V_{\text{cpn},n}(T_{n-2})$ by integration removes the digital discontinuity in the coupon, helping to stabilize Monte Carlo based sensitivity computations for $V_{\text{cpn},n}(0)$ in (23.5).

23.3 Payoff Smoothing For Numerical Integration and PDEs

Upon reflection, it is clear that singularity-removal technique outlined in Section 23.2.4 works by smoothing out an irregular boundary condition by integrating it against a density kernel. A closely related idea involves a direct modification of the payoff, to pre-smooth it before numerical integration or PDE schemes (or even Monte Carlo, as covered in the next section) are applied. We discuss several such payoff smoothing techniques in this section. Let us quickly remind the reader that payoff smoothing has two different, but related benefits. First, payoff smoothing will improve convergence of greeks calculated by PDE or Monte Carlo methods as the number of PDE steps or Monte Carlo paths is increased: the smoother the payoff, the faster the convergence. We have covered this angle in Sections 2.5 and 3.3. Second, payoff smoothing will help alleviate the problems arising if we, for the various reasons mentioned in Section 23.1, are unable to keep discretization grids constant.

23.3.1 Introduction to Payoff Smoothing

In a nutshell, the method of *payoff smoothing* replaces one payoff with a smoother one, to which standard numerical integration or PDE methods are then applied. Payoff smoothing serves to remove points of discontinuity in the payoff and its derivatives which, as we have seen earlier, will help improve the stability and smoothness of greeks.

A simple example of payoff smoothing replaces $f(x)$ with its moving average,

$$f_{\text{smooth}}(x) = \frac{1}{\epsilon} \int_{x-\epsilon/2}^{x+\epsilon/2} f(y) dy, \quad (23.6)$$

for some small $\epsilon > 0$, the choice of which will be discussed later. Payoff smoothing based on (23.6) was already applied in Section 2.5.2 as the *continuity correction* method for improving convergence of numerical solutions of PDEs. An example of moving average payoff smoothing is presented in Figure 23.3.

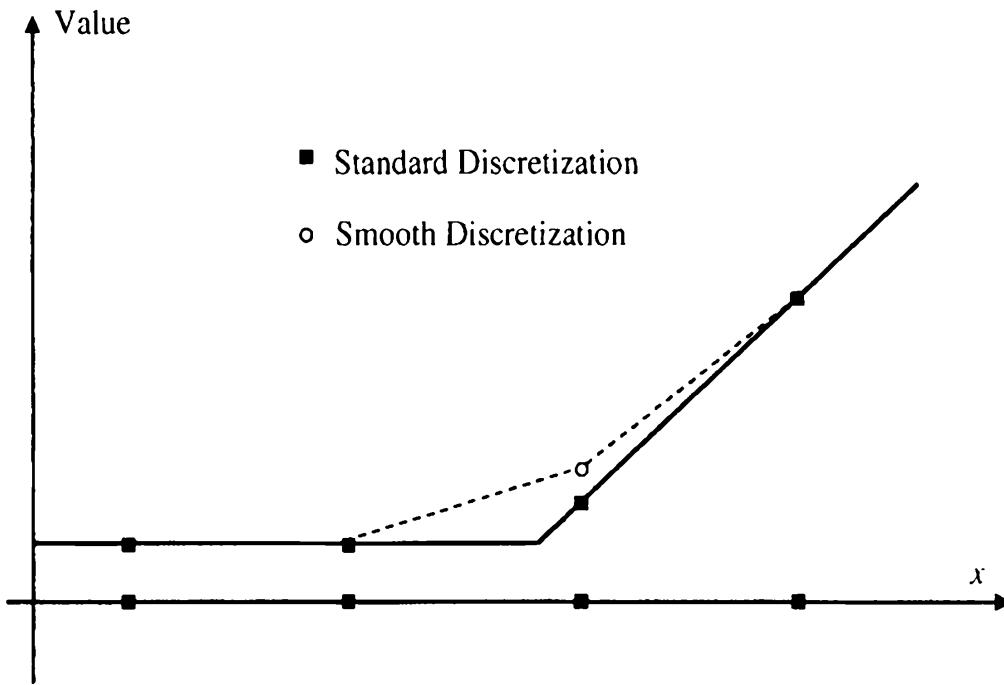
Continuing with the sample setup and the notations of Section 23.1, we recall that the standard numerical quadrature with knots $\{x_n\}$ and weights $\{w_n\}$ specifies that

$$\tilde{V}(S) = \sum_{n=0}^N w_n f(x_n + S). \quad (23.7)$$

The payoff smoothing method replaces this with

$$\tilde{V}(S) = \sum_{n=0}^N w_n f_{\text{smooth}}(x_n + S). \quad (23.8)$$

Fig. 23.3. Payoff Smoothing



Notes: Smoothing a discretized payoff using moving average (23.6).

Because $f_{\text{smooth}}(x)$ has a higher degree of smoothness than $f(x)$, the numerically computed greeks of $\tilde{V}(S)$ behave more smoothly with respect to market inputs.

The function $f_{\text{smooth}}(x)$ usually cannot be computed exactly. However, for small ϵ , various approximations can be made. Here, again, the knowledge of singularities of $f(x)$ is important. If $f(x)$ is known to have no singularities on $[x - \epsilon/2, x + \epsilon/2]$, a simple linear approximation to $f(x)$ on this interval will suffice as the corresponding term in (23.7) will be sufficiently smooth:

$$f(y) \approx f(x - \epsilon/2) + (f(x + \epsilon/2) - f(x - \epsilon/2)) \frac{y - x + \epsilon/2}{\epsilon}$$

for $y \in [x - \epsilon/2, x + \epsilon/2]$, so that

$$f_{\text{smooth}}(x) = \frac{1}{\epsilon} \int_{x-\epsilon/2}^{x+\epsilon/2} f(y) dy \approx f(x). \quad (23.9)$$

If, however, it is known that there is a singularity x^* in $[x - \epsilon/2, x + \epsilon/2]$, then the integral should be handled more carefully by, for example, using separate linear approximations to f in two subintervals, $[x - \epsilon/2, x^*)$ and $(x^*, x + \epsilon/2]$:

$$f(y) \approx \begin{cases} f(x - \epsilon/2) + (f(x^*-) - f(x - \epsilon/2)) \frac{y - x + \epsilon/2}{x^* - x + \epsilon/2}, & y \in [x - \epsilon/2, x^*), \\ f(x^*+) + (f(x + \epsilon/2) - f(x^+)) \frac{y - x^*}{x + \epsilon/2 - x^*}, & y \in (x^*, x + \epsilon/2], \end{cases}$$

so that

$$\begin{aligned} f_{\text{smooth}}(x) &= \frac{1}{\epsilon} \int_{x-\epsilon/2}^{x+\epsilon/2} f(y) dy \approx \frac{x^* - x + \epsilon/2}{\epsilon} f\left(\frac{x - \epsilon/2 + x^*}{2}\right) \\ &\quad + \frac{x + \epsilon/2 - x^*}{\epsilon} f\left(\frac{x + \epsilon/2 + x^*}{2}\right). \end{aligned} \quad (23.10)$$

The method is not dissimilar to the singularity-extended grid method, at least in one dimension, but could be more practical to apply in a PDE schemes, say, if multiple singularities at different locations are introduced at different times.

The performance of moving average smoothing will depend on the choice of ϵ . Higher values of ϵ lead to smoother schemes but makes it more difficult to approximate the required integral, since linear approximations may no longer be accurate enough. More importantly, the introduction of smoothing adds bias to the valuation, as the payoff being integrated becomes increasingly different from the actual one when ϵ is increased. In many cases the choice of ϵ is done semi-empirically, with numerical experiments to determines the highest value of the smoothing window that keeps the bias within acceptable limits. In some cases, the discretization of the grid itself drives the size of ϵ , a case that we consider next.

23.3.2 Payoff Smoothing in One Dimension

To develop the method above in more detail, and to link it more directly to grid-based methods, let us introduce a discrete set of x -values, $\{x_n\}_{n=0}^N$. It is helpful to think of (23.7) as a special case of a more general setup where we define the discretized value f_n as the weighted average of $f(x)$ in a neighborhood of x_n ,

$$f_n = \int_{-\infty}^{\infty} \kappa_n(x) f(x) dx, \quad n = 0, \dots, N, \quad (23.11)$$

with $\{\kappa_n(x)\}_{n=0}^N$ a collection of averaging weights (e.g. we use $\kappa_n(x) = \epsilon^{-1} \mathbf{1}_{\{x \in [x_n - \epsilon/2, x_n + \epsilon/2]\}}$ in the previous section), such that

$$\int_{-\infty}^{\infty} \kappa_n(x) dx = 1, \quad \kappa_n(x) \geq 0, \quad n = 0, \dots, N.$$

Often the weights are taken to be shifted and scaled versions of a common weight function, in the sense that

$$\kappa_n(x) = \frac{1}{\epsilon_n} \kappa\left(\frac{x - x_n}{\epsilon_n}\right),$$

where

$$\int_{-\infty}^{\infty} \kappa(x) dx = 1, \quad \kappa(x) \geq 0.$$

The weights are shifted to center them around the grid points $\{x_n\}_{n=0}^N$, and scaling parameters ϵ_n control the dispersion of the weight around x_n . As the scaling parameters tend to zero, the averages of f tend to the values of f on the grid $\{x_n\}$.

23.3.2.1 Box Smoothing

A particularly simple averaging weight is the indicator function on the interval between middle-points of the grid buckets to the either side of x_n ,

$$\begin{aligned}\kappa_n(x) &= (c_{n+1} - c_n)^{-1} 1_{\{x \in [c_n, c_{n+1}]\}}, \\ c_n &= (x_n + x_{n-1})/2, \quad n = 1, \dots, N-1.\end{aligned}\tag{23.12}$$

Because of the shape of the weight function, the method is sometimes called the *box smoothing method*. The resulting discretization formula is given by

$$f_n = \frac{1}{c_{n+1} - c_n} \int_{c_n}^{c_{n+1}} f(x) dx, \quad c_n = (x_n + x_{n-1})/2, \quad n = 1, \dots, N-1.\tag{23.13}$$

If the function $f(x)$ is known for all x , as is the case for numerical integration, the box smoothing method is easy to apply, using the arguments that lead to (23.9) and (23.10). A more challenging situation arises when only a *discretized* version of the payoff is known, as may happen when f represents a PDE solution rolled back to some intermediate date. While backward induction in a PDE is, in itself, a smoothing operation, singularities may be introduced through the enforcement of jump conditions, as required if the security in question happens to pay a coupon, is exercisable, or has a barrier condition of some kind. Sometimes (e.g., for barrier options) the location of the singularity is known exactly, but often (e.g., for Bermudan style options) it is not. This complicates the application of (23.13), since knowledge of the location of singularities is critical to our ability to compute smooth greeks. We proceed to discuss a scheme to handle the case when the singularity location is not known.

To properly fix our setup, consider a security with terminal payout date T^* whose value is being computed by solving the corresponding PDE numerically, backwards from T^* to time 0. Let $V(t, x)$ represent the true value of the security at time t at state x . As always, let $V(t-, x)$ and $V(t+, x)$ be the value of the security just before and just after time t , respectively. Assuming that a lifecycle event takes place at some intermediate time T (such as an exercise opportunity, a knock-in/knock-out barrier check, a fixing of a structured coupon, and so on), a jump condition will be applied when crossing from $T+$ to $T-$ in the backward recursion scheme, see Section 2.7. Specifically, if $\{V_n^+\}_{n=0}^N$ represents the numerical approximation

to $V(T+, x)$ on the grid $\{x_n\}_{n=0}^N$, the jump condition determines how to compute $\{V_n^-\}_{n=0}^N$, the grid approximation to $V(T-, x)$. Here we make an important observation: most jump conditions can be represented in the following form,

$$V(T-, x) = 1_{\{g(x) \leq h(x)\}} p(x) + 1_{\{g(x) > h(x)\}} q(x), \quad (23.14)$$

where the discretized versions of the *smooth* functions $g(x)$, $h(x)$, $p(x)$ and $q(x)$ are known at $t = T+$.⁶ Some of these functions could be based on $V(T+, x)$, and others are defined by the specifics of the event. Let us give a few examples.

Example 23.3.1. If the security can be canceled at time T , then

$$V(T-, x) = 1_{\{V(T+, x) > 0\}} V(T+, x),$$

i.e. $g(x) = V(T+, x)$, $h(x) = 0$, $p(x) = 0$, $q(x) = V(T+, x)$.

Example 23.3.2. If the security is callable at time T with the exercise value $e(x)$, then

$$V(T-, x) = 1_{\{V(T+, x) \leq e(x)\}} e(x) + 1_{\{V(T+, x) > e(x)\}} V(T+, x),$$

i.e. $g(x) = V(T+, x)$, $h(x) = e(x)$, $p(x) = e(x)$, $q(x) = V(T+, x)$.

Example 23.3.3. Suppose the security knocks out at time T if the rate $r(x)$ is above the barrier b , and the knockout rebate is $a(x)$. Then

$$V(T-, x) = 1_{\{r(x) \leq b\}} V(T+, x) + 1_{\{r(x) > b\}} a(x),$$

i.e. $g(x) = r(x)$, $h(x) = b$, $p(x) = V(T+, x)$, $q(x) = a(x)$.

Example 23.3.4. If a coupon of the form $\max(r(x), s)$ is paid at time T , then

$$V(T-, x) = 1_{\{r(x) \leq s\}} (V(T+, x) + s) + 1_{\{r(x) > s\}} (V(T+, x) + r(x)),$$

i.e. $g(x) = r(x)$, $h(x) = s$, $p(x) = V(T+, x) + s$, $q(x) = V(T+, x) + r(x)$.

Going forward we assume that the event in question can indeed be represented in the form (23.14); let us denote the discretized versions of the smooth functions involved by $\{g_n\}$, $\{h_n\}$, $\{p_n\}$ and $\{q_n\}$.

Turning to the question of localizing singularities in $V(T-, x)$, we notice that the representation of the function in the form (23.14) simplifies our search, as all singularities are given by the solutions to the equation

⁶If there is more than one singularity introduced in the event, the decomposition above holds locally around each singularity. We consider a single singularity case only, with trivial extension to multiple ones.

$$g(x) - h(x) = 0.$$

Assume for simplicity that this equation has only one root⁷, and denote it by x^* . The problem of finding x^* is complicated somewhat by the fact that the values of functions $g(x)$, $h(x)$ are only known at the grid points $\{x_n\}$. However, since $g(x)$ and $h(x)$ are smooth, linear interpolation on each of the intervals $[x_{n-1}, x_n]$, $n = 1, \dots, N$, can be used instead. Specifically, we define

$$\begin{aligned}\hat{h}(x) &= h_n \frac{x - x_{n-1}}{x_n - x_{n-1}} + h_{n-1} \frac{x_n - x}{x_n - x_{n-1}}, \quad x \in [x_{n-1}, x_n], \\ \hat{g}(x) &= g_n \frac{x - x_{n-1}}{x_n - x_{n-1}} + g_{n-1} \frac{x_n - x}{x_n - x_{n-1}}, \quad x \in [x_{n-1}, x_n],\end{aligned}$$

and use the solution \hat{x}^* to $\hat{g}(x) - \hat{h}(x) = 0$ as an approximation to x^* . To locate x^* , first note that an interval $[x_{n-1}, x_n]$ contains x^* (and \hat{x}^*) provided that

$$(g_{n-1} - h_{n-1})(g_n - h_n) \leq 0.$$

If this inequality is satisfied, we pinpoint \hat{x}^* by solving a linear equation $\hat{g}(x) - \hat{h}(x) = 0$ on the interval $[x_{n-1}, x_n]$, so that \hat{x}^* is a solution to

$$h_n \frac{x - x_{n-1}}{x_n - x_{n-1}} + h_{n-1} \frac{x_n - x}{x_n - x_{n-1}} = g_n \frac{x - x_{n-1}}{x_n - x_{n-1}} + g_{n-1} \frac{x_n - x}{x_n - x_{n-1}}.$$

After some trivial algebraic manipulations,

$$\hat{x}^* = \frac{h_n - g_n}{(h_n - g_n) + (g_{n-1} - h_{n-1})} x_{n-1} + \frac{g_{n-1} - h_{n-1}}{(h_n - g_n) + (g_{n-1} - h_{n-1})} x_n.$$

Having established (an approximation of) the singularity in $V(T-, x)$, we can proceed to approximate the integrals

$$\frac{1}{c_{n+1} - c_n} \int_{c_n}^{c_{n+1}} V(T-, x) dx,$$

for all n . There are two cases to consider. When the root is not inside $[c_n, c_{n+1}]$, i.e. $\hat{x}^* \notin [c_n, c_{n+1}]$, (23.9) tells us that we can simply use the value of $V(T-, x)$ at x_n as an approximation to the integral. In other words, for such n we set

$$V_n^- = 1_{\{g_n \leq h_n\}} p_n + 1_{\{g_n > h_n\}} q_n.$$

For the case when $\hat{x}^* \in [c_n, c_{n+1}]$, we split the integration domain into two, $[c_n, \hat{x}^*]$ and $(\hat{x}^*, c_{n+1}]$. By assumption, the function $V(T-, x)$ is smooth on each of the two intervals, so according to (23.10) we can approximate each

⁷In other words, we assume that the discretization is fine enough to guarantee that there is only one singularity per interval.

of the two integrals by the value of the function $V(T-, x)$ in the center of each of the two subintervals,

$$V_n^- = \frac{\hat{x}^* - c_n}{c_{n+1} - c_n} V\left(T-, \frac{\hat{x}^* + c_n}{2}\right) + \frac{c_{n+1} - \hat{x}^*}{c_{n+1} - c_n} V\left(T-, \frac{\hat{x}^* + c_{n+1}}{2}\right). \quad (23.15)$$

As should be clear from (23.14), on the interval $[c_n, \hat{x}^*]$ the function $V(T-, x)$ is equal to one of the functions $p(x)$ or $q(x)$, and on the interval $(\hat{x}^*, c_{n+1}]$ it is equal to the other. More precisely,

$$\begin{aligned} V(T-, x) &= p(x) \quad \text{for } x \in [c_n, \hat{x}^*], \\ V(T-, x) &= q(x) \quad \text{for } x \in (\hat{x}^*, c_{n+1}], \end{aligned} \quad (23.16)$$

if and only if $g_{n-1} < h_{n-1}$. Assume for concreteness that (23.16) in fact holds. Then (23.15) can be rewritten as

$$V_n^- = \frac{\hat{x}^* - c_n}{c_{n+1} - c_n} p\left(\frac{\hat{x}^* + c_n}{2}\right) + \frac{c_{n+1} - \hat{x}^*}{c_{n+1} - c_n} q\left(\frac{\hat{x}^* + c_{n+1}}{2}\right).$$

To find $p((\hat{x}^* + c_n)/2)$, $q((\hat{x}^* + c_{n+1})/2)$, we (once again) use the fact that the functions $p(x)$ and $q(x)$ are smooth, and approximate them linearly. Then, $p((\hat{x}^* + c_n)/2)$, $q((\hat{x}^* + c_{n+1})/2)$ can be computed from the known values of $p(x)$ and $q(x)$ on the grid $\{x_n\}$,

$$p\left(\frac{\hat{x}^* + c_n}{2}\right) \approx \hat{p}, \quad q\left(\frac{\hat{x}^* + c_{n+1}}{2}\right) \approx \hat{q},$$

where

$$\begin{aligned} \hat{p} &= p_{n-1} \frac{x_n - (\hat{x}^* + c_n)/2}{x_n - x_{n-1}} + p_n \frac{(\hat{x}^* + c_n)/2 - x_{n-1}}{x_n - x_{n-1}}, \\ \hat{q} &= q_n \frac{x_{n+1} - (\hat{x}^* + c_{n+1})/2}{x_{n+1} - x_n} + q_{n+1} \frac{(\hat{x}^* + c_{n+1})/2 - x_n}{x_{n+1} - x_n}. \end{aligned}$$

Combining various approximations, we finally obtain that for n such that $\hat{x}^* \in [c_n, c_{n+1}]$,

$$V_n^- = \frac{\hat{x}^* - c_n}{c_{n+1} - c_n} \hat{p} + \frac{c_{n+1} - \hat{x}^*}{c_{n+1} - c_n} \hat{q},$$

with \hat{p} , \hat{q} given just above. This concludes the description of the box smoothing method.

23.3.2.2 Other Smoothing Methods

In (23.11), weight functions other than indicator functions could be considered. A fairly popular alternative to box smoothing uses a weight function based on the linear Lagrange (triangular, or “hat”) weight functions that

we defined in footnote 10 on page 58. The resulting smoothing method is often known as the *hat smoothing method*. Apart from a slightly different functional form of the weight functions, its implementation differs little from the box smoothing method. In particular, for the method to be fully effective, discontinuities still need to be detected (as in the previous section) and incorporated into the calculation of integrals.

Unlike indicator functions, triangular weight functions are continuous, leading to smoother greeks than for box smoothing, especially for higher-order greeks such as gammas. On the other hand, hat smoothing is less “local” than box smoothing, in the sense that f_n computed in hat smoothing will depend on the values of $f(x)$ for $x \in [x_{n-1}, x_{n+1}]$, whereas f_n in box smoothing only depends on the values of $f(x)$ for $x \in [(x_n + x_{n-1})/2, (x_{n+1} + x_n)/2]$. This could be important when, for example, a trade feature (such as the exercise boundary) is close to the initial point of the grid in time and space, where more local functions tend to give better resolution, i.e. lower bias.

Stability of greeks of even higher order can be obtained by using weights that are even more smooth, i.e. Gaussian kernels. This, however, leads to more computationally intensive schemes, as well as schemes of ever more deteriorating locality.

23.3.3 Payoff Smoothing in Multiple Dimensions

The weight-based smoothing methods discussed in the previous sections also apply to multi-dimensional PDEs, although certain challenges quickly become evident, especially when one attempts to split the smoothing integrals into sub-domains around singularities. In one dimension, a singularity is just a single point which splits a single interval into two subintervals. In two dimensions, a singularity is typically a curve which affects multiple rectangles, and splits each affected rectangle into two subdomains of generally irregular geometry. Things get even more complicated in dimensions 3 and higher.

Let us consider the case of box smoothing in two dimensions in more detail, as it presents most of the challenges appearing in higher dimensions. We define $\{x_n\}_{n=0}^N$ and $\{y_m\}_{m=0}^M$ to be grids in x and y dimensions, and denote by (c_n^x, c_m^y) the center point of the rectangle $[x_{n-1}, x_n] \times [y_{m-1}, y_m]$,

$$c_n^x = (x_n + x_{n-1})/2, \quad c_m^y = (y_m + y_{m-1})/2.$$

Furthermore, let us define a rectangle

$$D_{n,m} = [c_n^x, c_{n+1}^x] \times [c_m^y, c_{m+1}^y].$$

Then we can introduce a collection of two-dimensional box weights

$$\kappa_{n,m}(x) = |D_{n,m}|^{-1} \mathbf{1}_{\{(x,y) \in D_{n,m}\}}, \quad n = 1, \dots, N-1, \quad m = 1, \dots, M-1,$$

where $|D|$ is the area of D . To smooth a function $f(x, y)$, an integral

$$f_{n,m} = \frac{1}{|D_{n,m}|} \int_{D_{n,m}} f(x, y) dx dy$$

is calculated for each n, m .

Recall that the one-dimensional box smoothing method we presented in Section 23.3.2.1 was based on the representation (23.14). In a similar vein, and using similar notations, we assume that the time $T-$ value of the security is given by

$$V(T-, x, y) = 1_{\{g(x, y) \leq h(x, y)\}} p(x, y) + 1_{\{g(x, y) > h(x, y)\}} q(x, y). \quad (23.17)$$

Here the discretized versions of the smooth functions $g(x, y)$, $h(x, y)$, $p(x, y)$ and $q(x, y)$ are assumed known at time $T+$ on the 2-dimensional mesh $\{(x_n, y_m)\}$.

The singularity of $V(T-, x, y)$ is given by the one-dimensional curve

$$s \subset \mathbb{R}^2, \quad s = \{(x, y) : g(x, y) = h(x, y)\}.$$

As in one dimension, the value

$$\frac{1}{|D_{n,m}|} \int_{D_{n,m}} V(T-, x, y) dx dy \quad (23.18)$$

for $D_{n,m}$ such that $D_{n,m} \cap s = \emptyset$ can be approximated with the value $V(T-, x_n, y_m)$. If, however, $D_{n,m} \cap s \neq \emptyset$, the domain $D_{n,m}$ needs to be split up into two, and the integral of V computed on each of the subdomains separately. We note that, in general, there will be many rectangles $D_{n,m}$ where $D_{n,m} \cap s \neq \emptyset$.

The box smoothing method naturally splits into the following steps.

1. Approximate the value $g(x, y) - h(x, y)$ at the corners of each of $D_{n,m}$ by linear interpolation. We denote

$$\hat{\xi}_{n,m} \triangleq \hat{g}(c_n^x, c_m^y) - \hat{h}(c_n^x, c_m^y),$$

where \hat{g} and \hat{h} are approximations to g , h computed using bi-linear interpolation off the grid $\{(x_n, y_m)\}$ (on which the values of g and h are known by assumption).

2. Find those rectangles $D_{n,m}$ for which $D_{n,m} \cap s \neq \emptyset$. The search can be conducted efficiently by looking for those $D_{n,m}$ for which the signs of the difference $g(x, y) - h(x, y)$ are not all the same in the corners of $D_{n,m}$. Specifically, we decide that $D_{n,m}$ contains a singularity if not all of $\hat{\xi}_{n,m}, \hat{\xi}_{n+1,m}, \hat{\xi}_{n,m+1}, \hat{\xi}_{n+1,m+1}$ have the same sign.
3. For those $D_{n,m}$ that do not contain a singularity, approximate the integral (23.18) with $V(T-, x_n, y_m)$, available from (23.17).

4. For those $D_{n,m}$ that do contain a singularity, approximate the singularity curve with a straight line through those two edges of the rectangle $D_{n,m}$ that are crossed by s . Note that s crosses the edge $(c_n^x, c_m^y) \rightarrow (c_{n+1}^x, c_m^y)$ if $\hat{\xi}_{n,m} \cdot \hat{\xi}_{n+1,m} \leq 0$, and so on. Find approximations to the positions of the points where s crosses the two edges by linearly interpolating the appropriate values of $\xi_{n,m}, \hat{\xi}_{n+1,m}, \hat{\xi}_{n,m+1}, \hat{\xi}_{n+1,m+1}$. For example, if s crosses the edge $(c_n^x, c_m^y) \rightarrow (c_{n+1}^x, c_m^y)$, then the point where s crosses the edge is given by $(\hat{x}_1^*, \hat{y}_1^*)$, where $\hat{y}_1^* = c_m^y$, and \hat{x}_1^* is a solution to

$$\xi_{n,m} + \frac{\xi_{n+1,m} - \xi_{n,m}}{c_{n+1}^x - c_n^x} (x - c_n^x) = 0.$$

Denote the second crossing point by $(\hat{x}_2^*, \hat{y}_2^*)$ and approximate the part of s inside $D_{n,m}$ by a straight line⁸ segment connecting $(\hat{x}_1^*, \hat{y}_1^*)$ and $(\hat{x}_2^*, \hat{y}_2^*)$.

5. Compute the integral (23.18) on each of the two subdomains of $D_{n,m}$. For that, split $D_{n,m}$ into two parts, $D_{n,m}^1$ and $D_{n,m}^2$, by the line segment connecting $(\hat{x}_1^*, \hat{y}_1^*)$ and $(\hat{x}_2^*, \hat{y}_2^*)$. On one subdomain ($D_{n,m}^1$ for concreteness), $V(T-, x, y)$ is equal to $p(x, y)$ and on the other ($D_{n,m}^2$) to $q(x, y)$. Use the fact that p and q are smooth and approximate them with linear functions \hat{p} and \hat{q} on $D_{n,m}^i$, $i = 1, 2$. To compute the integrals

$$\int_{D_{n,m}^1} \hat{p}(x, y) \, dx \, dy, \quad \int_{D_{n,m}^2} \hat{q}(x, y) \, dx \, dy,$$

use the fact that

$$\int_P l(x, y) \, dx \, dy = l(x_P, y_P) \times |P|$$

for any (integrable) domain P and linear function $l(x, y)$, where (x_P, y_P) is the center of mass of P .

6. Repeat Steps 4 and 5 for all $D_{n,m}$ such that $D_{n,m} \cap s \neq \emptyset$.

In three dimensions, singularities are given by two-dimensional surfaces; within each discretization cube, the singularities can be approximated by planes, and various cases as to how these planes intersect cubes need to be considered. In four dimensions, singularities must be approximated by cubes, and so forth. As the dimensionality of the problem grows, so does the amount of effort required to do smoothing. In a K -dimensional space, a singularity has dimension $K - 1$, so if the K -dimensional grid has an order N discretization for each of the K dimensions, then the number of K -dimensional grid segments (intervals in one-dimension, rectangles in two dimensions, cubes in three dimensions, etc.) that intersect the singularity

⁸In line with the footnote 7 we assume that the rectangles are small enough so that we can assume there is only two crossing points.

is of order N^{K-1} . The amount of work required per segment also generally grows with K , so a direct implementations of the smoothness algorithm in dimension 3 and above can already constitute a significant, if not dominant, proportion of calculation time.

In some multi-dimensional problems, the workload can be reduced by using known features of a product and/or model to understand the structure of singularities. For example, in some cases one of the PDE dimensions in a given model can be identified as being “dominant”, in the sense that the singularity surface will be mostly orthogonal to this direction. In this case, rather than applying the full multi-dimensional smoothing method, a series of one-dimensional smoothing methods in the dominant dimension can be used instead, often at considerable time savings and with good smoothing results. For example, in SV models singularities will typically be present in the asset (S) dimension, while the payoff will be smooth in the variance (z) dimension. This suggests a scheme where one applies, at each discretized value of the stochastic variance, *one-dimensional* payoff smoothing in the direction of the asset state variable.

23.4 Payoff Smoothing for Monte Carlo

While the methods of Section 23.3.1 can be applied in a Monte Carlo setting, there exist more natural payoff smoothing methods for Monte Carlo applications. Starting from a very simple example in Section 23.4.1, we construct one such method here. The method is designed to be applied when calculating sensitivities by direct perturbation (“bump-and-reprice”). Alternative methods for sensitivity computations by Monte Carlo are described in Section 3.3 and, in particular, in Chapter 24).

23.4.1 Tube Monte Carlo for Digital Options

One situation where payoff smoothing is almost universally applied is in valuing digital options. Consider an option that pays

$$1_{\{S(T) > B\}} \tag{23.19}$$

at time T , where $S(t)$ is the process for the underlying that is simulated using Monte Carlo.

Since the payoff is discontinuous, the Monte Carlo estimate of $V = E(1_{\{S(T) > B\}})$, defined by (ω_j are sample paths, $j = 1, \dots, J$)

$$V \approx J^{-1} \sum_{j=1}^J 1_{\{S(T, \omega_j) > B\}},$$

exhibits poor convergence (see Section 3.3.1.1) and unstable greeks. The standard way to remedy this is to replace the digital payoff with a call

spread (or use the likelihood ratio method — but we do not consider it at the moment). Let us choose $K_1 < K_2$ and replace

$$1_{\{x>B\}} \approx f(x), \quad f(x) = \max \left(\min \left(\frac{x - K_1}{K_2 - K_1}, 1 \right), 0 \right).$$

Since $f(x)$ is smoother than $1_{\{x>B\}}$ (at least it is continuous), the stability of greeks is improved. Various choices of K_1, K_2 are possible. If $K_1 = B$, the call spread with the payoff $f(x)$ is often called an “underhedge” (as $f(x) \leq 1_{\{x>B\}}$ for all x). Conversely, if $K_2 = B$, then the call spread with the payoff $f(x)$ is called an “overhedge” ($f(x) \geq 1_{\{x>B\}}$ for all x). A symmetric payoff with $K_1 = B - \epsilon, K_2 = B + \epsilon, \epsilon > 0$, is used most often when the goal is to improve greeks stability while minimizing the bias introduced by the smoothness method. In this case,

$$f_{\text{sym}}(x) = \max \left(\min \left(\frac{x + \epsilon - B}{2\epsilon}, 1 \right), 0 \right). \quad (23.20)$$

The choice of the smoothing window $\epsilon = (K_2 - K_1)/2$ involves a trade-off between a high degree of smoothness (large ϵ) and low bias (low ϵ) and is usually performed experimentally. As we already mentioned, a typical strategy involves formulating a maximum tolerable level of the difference between the values of options with payoffs $1_{\{x>B\}}$ and $f(x)$, and then setting ϵ accordingly.

The method above is (in its symmetric form, at least) a special case of the moving average smoothing approach from Section 23.3.1. As we shall show next, in a Monte Carlo setting the method can also be justified a completely different way. While not particularly useful for the specific case of the digital option, this alternative interpretation allows us to formulate a generic Monte Carlo smoothing strategy applicable to more complicated payoffs.

In the general spirit of payoff smoothing, let us first replace the standard Monte Carlo approximation (23.19) with the following one,

$$V \approx J^{-1} \sum_{j=1}^J V_j, \quad V_j = \mathbb{E} (1_{\{S(T)>B\}} | A_j), \quad (23.21)$$

where \mathbb{E} is the expected value operator that corresponds to the pricing measure Q (whose exact nature is unimportant here, but could be taken to be the risk-neutral measure for concreteness), and where A_j is defined as a small interval centered at $S(T, \omega_j)$,

$$A_j = \{\omega : S(T, \omega) \in [S(T, \omega_j) - \epsilon, S(T, \omega_j) + \epsilon]\}, \quad \epsilon > 0.$$

The difference between (23.21) and the standard estimate (23.19) comes from replacing the “point” sample of the payoff $1_{\{x>B\}}$ at $S(T, \omega_j)$ with an

“average” estimate of the payoff in a small interval around the sample asset value $S(T, \omega_j)$.

To compute $E(1_{\{S(T) > B\}} | A_j)$, we assume that the distribution of $S(T)$ within the interval $[S(T, \omega_j) - \epsilon, S(T, \omega_j) + \epsilon]$ can be approximated with a uniform distribution. If ϵ is small, the error introduced by this approximation is small. Then, if $B \notin A_j$, we have that $E(1_{\{S(T) > B\}} | A_j) = 1_{\{S(T, \omega_j) > B\}}$. If, however, $B \in A_j$, then (using the uniform distribution approximation as discussed above)

$$\begin{aligned} E(1_{\{S(T) > B\}} | A_j) &= Q(S(T) > B | S(T, \omega_j) - \epsilon \leq S(T) \leq S(T, \omega_j) + \epsilon) \\ &= \frac{S(T, \omega_j) + \epsilon - B}{2\epsilon}. \end{aligned}$$

Combining the two cases into one formula, we obtain

$$\begin{aligned} V_j &= E(1_{\{S(T) > B\}} | A_j) \\ &= \frac{S(T, \omega_j) + \epsilon - B}{2\epsilon} \times 1_{\{S(T, \omega_j) - \epsilon \leq B < S(T, \omega_j) + \epsilon\}} \\ &\quad + 0 \times 1_{\{B < S(T, \omega_j) - \epsilon\}} + 1 \times 1_{\{S(T, \omega_j) + \epsilon \leq B\}} \\ &= f_{\text{sym}}(S(T, \omega_j)). \end{aligned}$$

Hence, (23.21) can be rewritten as

$$V \approx J^{-1} \sum_{j=1}^J f_{\text{sym}}(S(T, \omega_j)),$$

and the “call spread” method is motivated from the probabilistic perspective.

The derivation above points to a systematic approach for obtaining Monte Carlo specific payoff smoothing approximations for a wide variety of payoffs. First, we replace point estimates of the payoff along each sample path with averages of the payoff over a suitably defined small neighborhood of the sample path. Then, to compute the required average value over each neighborhood, we use various approximations that can be justified by the fact that each neighborhood is small. We call this method the *tube Monte Carlo* (also sometimes known as *sausage Monte Carlo*, see Piterbarg [2004c]), with the name reflecting the fact that small neighborhoods around sample paths resemble thin, narrow (multi-dimensional) “tubes”. In the next section, we apply the tube Monte Carlo method to a more interesting class of payoffs.

23.4.2 Tube Monte Carlo for Barrier Options

Consider a derivative which is a knock-in barrier into a stream of (net) coupons X_1, \dots, X_{N-1} , with the knock-in feature defined by a stopping time index η : the derivative pays coupons X_i at T_{i+1} for $i = \eta, \dots, N-1$. The value of the security is then given by

$$V_{ki}(0) = E \left(\sum_{i=1}^{N-1} B(T_{i+1})^{-1} X_i 1_{\{i \geq \eta\}} \right), \quad (23.22)$$

where E is the expected value operator for the spot measure Q^B , with $B(t)$ in (23.22) being the discretely compounded money market account. The standard estimate of this value in Monte Carlo simulation is given by

$$V_{ki}(0) \approx \frac{1}{J} \sum_{j=1}^J \left(\sum_{i=1}^{N-1} B(T_{i+1}, \omega_j)^{-1} X_i(\omega_j) 1_{\{i \geq \eta(\omega_j)\}} \right), \quad (23.23)$$

where $\omega_1, \dots, \omega_J$ are Monte Carlo paths and where we use the notation $\xi(\omega)$ for the value of a random variable ξ on path ω .

The indicator functions $1_{\{i \geq \eta(\omega)\}}$ in (23.23) introduce digital discontinuities in the payoff which, as we know, lead to poor stability of risk sensitivities. To improve on this situation, let us consider how to apply the payoff smoothing ideas of Section 23.4.1 here. Let $x(t, \omega)$ be the d -dimensional vector of state variables of the underlying model which we, without practical loss of generality, assume to be Markovian. We further assume that we can write the stopping time index η as the first hitting time of a state-dependent boundary,

$$\eta(\omega) = \min \{n \geq 1 : \psi_n(x(T_n, \omega)) \geq 0\} \wedge N, \quad (23.24)$$

where $\psi_n(x)$ are some functions, $\psi_n : \mathbb{R}^d \rightarrow \mathbb{R}$.

The idea of the tube method is to replace point estimates (23.23) of the payoff with payoff averages over appropriately defined tubes. Let us fix $\epsilon > 0$, the width of the tube. For each j we define the ϵ -tubes in the state space by

$$A_j^\epsilon = \bigcap_{i=1}^{N-1} A_{j,i}^\epsilon, \quad (23.25)$$

$$A_{j,i}^\epsilon = \{\omega : \|x(T_i, \omega) - x(T_i, \omega_j)\| < \epsilon\},$$

where, essentially, A_j^ϵ denotes the set of all sample paths that come within ϵ -distance of $x(T_i, \omega_j)$'s for all T_i , $i = 1, \dots, N-1$. Then we replace (23.23) with the following estimator,

$$V_{ki}(0) \approx J^{-1} \sum_{j=1}^J V_j, \quad V_j \triangleq E \left(\sum_{i=1}^{N-1} [B(T_{i+1}, \omega)^{-1} X_i(\omega) 1_{\{i \geq \eta(\omega)\}}] \middle| A_j^\epsilon \right). \quad (23.26)$$

Since $B(T_{i+1}, \omega)^{-1}$ and $X_i(\omega)$ are, often, smooth⁹ functions of the path ω , we evaluate them just at the sample path,

⁹ $X_i(\omega)$ can, of course, be discontinuous, but this is not our focus at the moment.

$$V_j \approx \sum_{i=1}^{N-1} B(T_{i+1}, \omega_j)^{-1} X_i(\omega_j) E \left(1_{\{i \geq \eta(\omega)\}} \mid A_j^\epsilon \right), \quad (23.27)$$

which approximates (23.26) to order ϵ . To proceed we need the following proposition.

Proposition 23.4.1. *In (23.27) the probabilities*

$$q_i(\omega_j) \triangleq E \left(1_{\{i \geq \eta(\omega)\}} \mid A_j^\epsilon \right),$$

can be approximated as follows:

$$1 - q_i(\omega_j) = \prod_{n=1}^i (1 - p_n(\omega_j)), \quad (23.28)$$

$$p_n(\omega_j) \triangleq \begin{cases} 1, & \psi_{n,j} - \delta_{n,j} \geq 0, \\ \frac{\delta_{n,j} + \psi_{n,j}}{2\delta_{n,j}}, & \psi_{n,j} - \delta_{n,j} < 0 < \psi_{n,j} + \delta_{n,j}, \\ 0, & \psi_{n,j} + \delta_{n,j} \leq 0, \end{cases}$$

where

$$\psi_{n,j} \triangleq \psi_n(x(T_n, \omega_j)), \quad \delta_{n,j} \triangleq \epsilon \|\nabla \psi_{n,j}\|, \quad \nabla \psi_{n,j} \triangleq \nabla \psi_n(x)|_{x=x(T_n, \omega_j)}.$$

Proof. By expressing A_j^ϵ in terms of the functions ψ_n that define the knock-in index time in (23.24), we get

$$1 - q_i(\omega_j) = Q^B \left(\bigcap_{n=1}^i \{\psi_n(x(T_n, \omega)) \leq 0\} \mid A_j^\epsilon \right). \quad (23.29)$$

We claim that, to order ϵ ,

$$1 - q_i(\omega_j) = \prod_{n=1}^i Q^B (\psi_n(x(T_n, \omega)) \leq 0 \mid A_{j,n}^\epsilon). \quad (23.30)$$

The proof follows by repeated applications of Lemma 23.B.1 from Appendix 23.B, although the intuition behind it is rather simple. Conditioning on $A_j^\epsilon = \bigcap_i A_{j,i}^\epsilon$ is essentially equivalent to pinning down the Markov process at times T_i , $i = 1, \dots, N-1$, to known values $\{x(T_i, \omega_j)\}$ with ϵ -accuracy. If a Markov process is conditioned on being at a certain state on a given date, past and future events become conditionally independent. Then, the set intersection on the right-hand side of (23.29) can be unwrapped into the product on the right-hand side of (23.30).

Now, define

$$\begin{aligned} p_n(\omega_j) &= Q^B (\psi_n(x(T_n), \omega) > 0 \mid A_{j,n}^\epsilon) \\ &= Q^B (\psi_n(x(T_n, \omega)) > 0 \mid \|x(T_n, \omega) - x(T_n, \omega_j)\| < \epsilon), \end{aligned}$$

so that (23.30) becomes

$$1 - q_i(\omega_j) = \prod_{n=1}^i (1 - p_n(\omega_j)).$$

To compute the p_n , observe that since we assumed that functions ψ_n are smooth, we may write for x such that $\|x - x(T_n, \omega_j)\| < \epsilon$,

$$\psi_n(x) \approx \psi_{n,j} + \nabla \psi_{n,j} \times (x - x(T_n, \omega_j)), \quad (23.31)$$

(here $\nabla \psi$ is the gradient of ψ , a row vector). Define $O_{n,j} \subset \mathbb{R}$ by

$$O_{n,j} = \psi_n(\{z \in \mathbb{R}^d : \|z - x(T_n, \omega_j)\| < \epsilon\}),$$

i.e. the image of the ball $\|z - x(T_n, \omega_j)\| < \epsilon$ under mapping $\psi_n : \mathbb{R}^d \rightarrow \mathbb{R}$. Then, from (23.31),

$$O_{n,j} \approx [\psi_{n,j} - \|\nabla \psi_{n,j}\| \epsilon, \psi_{n,j} + \|\nabla \psi_{n,j}\| \epsilon],$$

where $\|\nabla \psi_{n,j}\|$ denotes the norm of the linear operator $\nabla \psi_{n,j}$. Under the approximation

$$\begin{aligned} A_{j,n}^\epsilon &= \{\omega : \|x(T_n, \omega) - x(T_n, \omega_j)\| < \epsilon\} \\ &\approx \{\omega : \psi_n(x(T_n, \omega)) \in O_{n,j}\} \end{aligned}$$

we get

$$p_n(\omega_j) \approx Q^B(\psi_n(x(T_n, \omega)) > 0 | \psi_n(x(T_n, \omega)) \in O_{n,j}).$$

Approximating conditional distribution of $\psi_n(x(T_n, \omega))$ by a uniform distribution on the set $O_{n,j}$ we obtain

$$\begin{aligned} p_n(\omega_j) &\approx \frac{|\{\psi_n > 0\} \cap O_{n,j}|}{|O_{n,j}|} \\ &= \frac{|\{\psi_n > 0\} \cap [\psi_{n,j} - \|\nabla \psi_{n,j}\| \epsilon, \psi_{n,j} + \|\nabla \psi_{n,j}\| \epsilon]|}{|[\psi_{n,j} - \|\nabla \psi_{n,j}\| \epsilon, \psi_{n,j} + \|\nabla \psi_{n,j}\| \epsilon]|}. \end{aligned}$$

where we use $|\cdot|$ to denote the length of intervals in \mathbb{R} . Denoting

$$\delta_{n,j} = \epsilon \|\nabla \psi_{n,j}\|,$$

we obtain

$$|[\psi_{n,j} - \|\nabla \psi_{n,j}\| \epsilon, \psi_{n,j} + \|\nabla \psi_{n,j}\| \epsilon]| = 2\delta_{n,j},$$

and

$$p_n(\omega_j) = \begin{cases} 1, & \psi_{n,j} - \delta_{n,j} \geq 0, \\ \frac{\delta_{n,j} + \psi_{n,j}}{2\delta_{n,j}}, & \psi_{n,j} - \delta_{n,j} < 0 < \psi_{n,j} + \delta_{n,j}, \\ 0, & \psi_{n,j} + \delta_{n,j} \leq 0. \end{cases} \quad (23.32)$$

This completes the derivation. \square

Combining the results together, the formula for the tube Monte Carlo of a discrete knock-in barrier is then given by

$$V_{ki}(0) \approx J^{-1} \sum_{j=1}^J V_j, \quad V_j = \sum_{i=1}^{N-1} B(T_{i+1}, \omega_j)^{-1} X_i(\omega_j) q_i(\omega_j), \quad (23.33)$$

with q_i 's given by Proposition 23.4.1.

Let us analyze this formula in some detail. The quantity $\psi_{n,j}$ in (23.28) measures how far into the knock-in region the state process went, so we call it the “overshoot” function. The quantity $\delta_{n,j}$ is the “window” over which the overshoot function is smoothed out. It is equal to the universal constant ϵ (smoothing window for the state variables $x(\cdot)$) times the size of the gradient of the overshoot function. This provides consistent scaling of smoothing windows across different times/simulated paths. If the overshoot function is high (above $\delta_{n,j}$) then the knock-in barrier is deemed completely breached, and we set $p_n(\omega_j) = 1$. If the overshoot function is low (below $-\delta_{n,j}$), knock-in region is deemed to not have been reached at all. And for cases in between, the knock-in barrier is considered “partially” breached¹⁰, and a weight of $(\delta_{n,j} + \psi_{n,j})/2\delta_{n,j}$ is used to measure the extent of the barrier breach. Another analogy uses the idea of a partial knock-in: if the path ω_j is near the knock-in boundary, relevant coupons get included in the derivative value only partially, with the weights $q_i(\omega_j)$ defining the fractions of the coupons that count. This is in contrast to the standard Monte Carlo formula (23.23) in which coupons get included in the value either completely or not at all.

Critically, the weights $p_n(\omega_j)$ change smoothly with ω_j as do therefore the V_j 's in (23.33) (unlike those in (23.23) with digital discontinuities). The tube Monte Carlo formula (23.33) converges to the standard formula (23.23) as ϵ gets small. Clearly, the larger the smoothing window ϵ is, the smoother the payoff becomes, resulting in more stable risk sensitivities. With larger ϵ , however, the bias of the approximation becomes larger. In practice, to balance smoothness versus accuracy, one would start with a small ϵ and then keep increasing it for as long as the observed bias in the price remains within pre-set tolerances. Once the upper acceptable bound on ϵ is established, it can be used in risk sensitivity calculations.

¹⁰The concept of “partial” membership in a set should be familiar to those schooled in *fuzzy logic* (see Zadeh [1965]), and tube Monte Carlo can, in fact, be considered a probabilistically motivated fuzzy logic algorithm. For more discussion of fuzzy logic applications to Monte Carlo sensitivity computations in finance, see Withington and Lucic [2009].

23.4.3 Tube Monte Carlo for Callable Libor Exotics

The method of Section 23.4.2 can be applied directly to callable Libor exotics in Monte Carlo (see Section 18.3) whose valuation often relies on representing them as knock-in discrete barriers with the knock-in defined by an estimate of the exercise index — see e.g. (18.27), (18.28). Interestingly, the exact value of a CLE is a *smooth* function of the underlying path (as we establish later in Chapter 24), yet the representation such as (18.28) introduces digital discontinuities in the payoff. Therefore, for CLEs it is more advisable to use risk calculation methods that are specifically adopted to the CLE structure and its smoothness; such methods are developed in Chapter 24. Tube Monte Carlo, however, still has its place in the arsenal of valuation methods for CLEs as it often integrates better with standard risk system designs, compared to the more specialized methods of Chapter 24. In terms of performance, the effectiveness of the tube method compared to the alternatives depends on many underlying factors, but it is shown in Piterbarg [2005a] that to achieve comparable risk stability, the tube Monte Carlo method typically requires only about 1/4 of the path count needed for the standard simulation. The pathwise differentiation method of Chapter 24 reduces the required number of paths by another factor of 3 to 4.

Most of the mechanics required to apply the tube Monte Carlo method to CLEs have already been developed in Section 23.4.2. In fact, we only need to describe the functions ψ_n that define knock-in (or, in the context of CLEs, exercise) regions for each exercise date. This is straightforward to do; with (18.27) in mind, we just set

$$\psi_n(x(T_n, \omega)) = \mathcal{C} \left(\tilde{U}_n(T_n) \right)^T \zeta(T_n, \omega) - \mathcal{C} \left(\tilde{H}_n(T_n) \right)^T \zeta(T_n, \omega),$$

where we treat the right-hand-side — that is, the difference between exercise and hold values as measured by exercise and hold regression polynomials applied to explanatory variables of the regression — as a function of the model state variable vector $x(T_n, \omega)$. The method of Section 23.4.2 now carries over unchanged.

23.4.4 Tube Monte Carlo for TARNs

A TARN (see Section 20.1) can be represented as a derivative that pays a stream of (net) coupons until a knock-out event takes place when a sum of structured coupons exceeds a certain target. As knock-out derivatives are closely related to knock-in's, it is no surprise that a tube Monte Carlo method similar to that of Section 23.4.2 can be developed for TARNs (see Piterbarg [2004c]).

Let us recall the main TARN valuation formula (20.2), which we rewrite in a form similar to (23.22),

$$V_{\text{tarn}}(0) = \mathbb{E} \left(\sum_{i=1}^{N-1} B(T_{i+1})^{-1} X_i 1_{\{i < \eta\}} \right),$$

where

$$\eta(\omega) = \min \{n \geq 1 : Q_n(\omega) - R \geq 0\} \wedge N.$$

Then, in a close analogy to (23.33), we can write the approximation formula for the tube Monte Carlo method,

$$V_{\text{tarn}}(0) \approx J^{-1} \sum_{j=1}^J V_j, \quad V_j = \sum_{i=1}^{N-1} B(T_{i+1}, \omega_j)^{-1} X_i(\omega_j)(1 - q_i(\omega_j)),$$

where

$$1 - q_i(\omega_j) = \prod_{n=1}^i (1 - p_n(\omega_j)),$$

$$p_n(\omega_j) = \min \left(\max \left(\frac{Q_n(\omega_j) - R + \delta_{n,j}}{2\delta_{n,j}}, 0 \right), 1 \right),$$

and

$$\delta_{n,j} = \epsilon \|\nabla Q_n(\omega_j)\|, \quad (23.34)$$

with $\nabla Q_n(\omega_j)$ understood to be the gradient of Q_n expressed as a function of the model state vector.

High level of accuracy is not really required when calculating scaling constants $\delta_{n,j}$ in (23.34). In particular, for efficiency reasons we may use a simpler, deterministic scaling

$$\delta_{n,j} = \epsilon_n$$

for a collection $\{\epsilon_n\}$ or even a time-independent deterministic scaling

$$\delta_{n,j} = \epsilon.$$

The same simplifications could, of course, be adopted for knock-in barrier options and callable Libor exotics.

23.A Appendix: Delta Continuity of Singularity-Enlarged Grid Method

To show that the derivative of $\tilde{U}(S)$ is continuous across the grid point, it is sufficient to show that the left derivative of $\tilde{U}(s)$ at S_m equals the right derivative (at S_m). To simplify notations, let us assume that $\Delta x_n \equiv \Delta$, $n = 1, \dots, N$, $\psi(x) \equiv 1$ for x in some neighborhood of x_m , and redefine

$$\overline{U}(S) = \frac{2}{\Delta} \tilde{U}(S).$$

We then have, for $\epsilon > 0$, that

$$\begin{aligned}\overline{U}(S_m + \epsilon) &= (f(K) + f(K - \Delta + \epsilon)) \left(1 - \frac{\epsilon}{\Delta}\right) \\ &\quad + (f(K + \epsilon) + f(K)) \frac{\epsilon}{\Delta} \\ &\quad + (f(K + \Delta + \epsilon) + f(K + \epsilon)), \\ \overline{U}(S_m) &= 2f(K) + f(K - \Delta) + f(K + \Delta), \\ \overline{U}(S_m - \epsilon) &= (f(K - \epsilon) + f(K - \Delta - \epsilon)) \\ &\quad + (f(K) + f(K - \epsilon)) \frac{\epsilon}{\Delta} \\ &\quad + (f(K + \Delta - \epsilon) + f(K)) \left(1 - \frac{\epsilon}{\Delta}\right).\end{aligned}$$

In particular

$$\begin{aligned}\overline{U}(S_m + \epsilon) - \overline{U}(S_m) &= (f(K + \Delta + \epsilon) + f(K + \epsilon)) - (f(K + \Delta) + f(K)) \\ &\quad + (f(K) + f(K - \Delta + \epsilon)) - (f(K) + f(K - \Delta)) \\ &\quad + \frac{\epsilon}{\Delta} (f(K + \epsilon) - f(K - \Delta + \epsilon)),\end{aligned}$$

and

$$\begin{aligned}D^+ \overline{U}(S_m) &\triangleq \lim_{\epsilon \downarrow 0} \epsilon^{-1} (\overline{U}(S_m + \epsilon) - \overline{U}(S_m)) \\ &= D^+ f(K) + \frac{1}{\Delta} (f(K) - f(K - \Delta)) \\ &\quad + D^+ f(K + \Delta) + D^+ f(K - \Delta).\end{aligned}\tag{23.35}$$

Likewise,

$$\begin{aligned}\overline{U}(S_m) - \overline{U}(S_m - \epsilon) &= (f(K) + f(K - \Delta)) - (f(K - \epsilon) + f(K - \Delta - \epsilon)) \\ &\quad + (f(K + \Delta) - f(K + \Delta - \epsilon)) \\ &\quad + \frac{\epsilon}{\Delta} (f(K + \Delta - \epsilon) - f(K - \epsilon)),\end{aligned}$$

and

$$\begin{aligned}D^- \overline{U}(S_m) &\triangleq \lim_{\epsilon \uparrow 0} \epsilon^{-1} (\overline{U}(S_m) - \overline{U}(S_m - \epsilon)) \\ &= D^- f(K) + \frac{1}{\Delta} (f(K + \Delta) - f(K)) \\ &\quad + D^- f(K + \Delta) + D^- f(K - \Delta).\end{aligned}\tag{23.36}$$

Combining (23.35), (23.36) together and using the fact that the derivative of $f(x)$ is continuous everywhere except at K , we obtain

$$\begin{aligned} D^+ \bar{U}(S_m) - D^- \bar{U}(S_m) &= (D^+ f(K) - D^- f(K)) \\ &\quad + \left(\frac{1}{\Delta} (f(K) - f(K - \Delta)) - \frac{1}{\Delta} (f(K + \Delta) - f(K)) \right). \end{aligned}$$

We note that, to the second order,

$$\begin{aligned} \frac{1}{\Delta} (f(K) - f(K - \Delta)) &\approx D^- f(K), \\ \frac{1}{\Delta} (f(K + \Delta) - f(K)) &\approx D^+ f(K), \end{aligned}$$

and thus, to the second order,

$$D^+ \bar{U}(S_m) \approx D^- \bar{U}(S_m).$$

We conclude that the quadrature method produces smooth deltas.

23.B Appendix: Proof of Approximate Conditional Independence for Tube Monte Carlo

Here we prove a lemma needed in the proof of Proposition 23.4.1. Let $x(t, \omega)$, $t \in [0, T]$, be a Markov process with a state space \mathbb{R}^d for some $d \geq 1$. Assume its transition density

$$Q(x(t) = z | x(s) = y)$$

is differentiable in z and y for all $t, s > 0$. Let $T_1 < T_2$ and define, for some $\epsilon > 0$,

$$U_i^\epsilon = \{\omega : \|x(T_i, \omega) - x_i\| < \epsilon\}, \quad i = 1, 2,$$

for some x_1, x_2 .

Lemma 23.B.1. *Let X_1 and X_2 be two subsets of the state space \mathbb{R}^d and define*

$$Z_i = \{\omega : x(T_i, \omega) \in X_i\}, \quad i = 1, 2.$$

Then

$$Q(Z_1 \cap Z_2 | U_1^\epsilon \cap U_2^\epsilon) = Q(Z_1 | U_1^\epsilon) Q(Z_2 | U_2^\epsilon) (1 + O(\epsilon))$$

as $\epsilon \rightarrow 0$.

Proof. We have,

$$Q(Z_1 \cap Z_2 | U_1^\epsilon \cap U_2^\epsilon) = \frac{Q(Z_1 \cap Z_2 \cap U_1^\epsilon \cap U_2^\epsilon)}{Q(U_1^\epsilon \cap U_2^\epsilon)}.$$

For the expression in the numerator we have

$$\begin{aligned} Q(Z_1 \cap Z_2 \cap U_1^\epsilon \cap U_2^\epsilon) &= E(1_{\{Z_1\}} 1_{\{Z_2\}} 1_{\{U_1^\epsilon\}} 1_{\{U_2^\epsilon\}}) \\ &= \int dy \int dz Q(x(T_1) = y) 1_{\{y \in X_1\}} 1_{\{\|y - x_1\| < \epsilon\}} \\ &\quad \times Q(x(T_2) = z | x(T_1) = y) 1_{\{z \in X_2\}} 1_{\{\|z - x_2\| < \epsilon\}}. \end{aligned}$$

As the transition density is differentiable, for y such that $1_{\{U_1^\epsilon\}}(y) \neq 0$ we have that

$$Q(x(T_2) = z | x(T_1) = y) = Q(x(T_2) = z | x(T_1) = x_1)(1 + O(\epsilon)),$$

so we can write

$$\begin{aligned} Q((Z_1 \cap Z_2 \cap U_1^\epsilon \cap U_2^\epsilon)) &= \int dy \int dz Q(x(T_1) = y) 1_{\{y \in X_1\}} 1_{\{\|y - x_1\| < \epsilon\}} \\ &\quad \times Q(x(T_2) = z | x(T_1) = x_1) 1_{\{z \in X_2\}} 1_{\{\|z - x_2\| < \epsilon\}} (1 + O(\epsilon)) \\ &= E(1_{\{Z_1\}} 1_{\{U_1^\epsilon\}}) E(1_{\{Z_2\}} 1_{\{U_2^\epsilon\}} | x(T_1) = x_1) (1 + O(\epsilon)). \end{aligned}$$

Now

$$\begin{aligned} E(1_{\{Z_2\}} 1_{\{U_2^\epsilon\}} | x(T_1) = x_1) &= E(1_{\{Z_2\}} | x(T_1) = x_1, U_2^\epsilon) E(1_{\{U_2^\epsilon\}} | x(T_1) = x_1) \\ &= E(1_{\{Z_2\}} | x(T_1) = x_1, x(T_2) = x_2) (1 + O(\epsilon)) E(1_{\{U_2^\epsilon\}} | x(T_1) = x_1), \end{aligned}$$

where again we used the regularity properties of the transition density. By the Markovian property and the regularity of density,

$$\begin{aligned} E(1_{\{Z_2\}} | x(T_1) = x_1, x(T_2) = x_2) &= E(1_{\{Z_2\}} | x(T_2) = x_2) \\ &= E(1_{\{Z_2\}} | U_2^\epsilon) (1 + O(\epsilon)). \end{aligned}$$

Therefore, up to $O(\epsilon)$,

$$\begin{aligned} Q(Z_1 \cap Z_2 \cap U_1^\epsilon \cap U_2^\epsilon) &= E(1_{\{Z_1\}} 1_{\{U_1^\epsilon\}}) \\ &\quad \times E(1_{\{Z_2\}} | x(T_2) = x_2) E(1_{\{U_2^\epsilon\}} | x(T_1) = x_1) (1 + O(\epsilon)). \end{aligned}$$

Hence

$$\begin{aligned} Q(Z_1 \cap Z_2 | U_1^\epsilon \cap U_2^\epsilon) &= \frac{E(1_{\{Z_1\}} 1_{\{U_1^\epsilon\}}) E(1_{\{Z_2\}} | U_2^\epsilon) E(1_{\{U_2^\epsilon\}} | x(T_1) = x_1) (1 + O(\epsilon))}{E(1_{\{U_1^\epsilon\}}) E(1 | U_2^\epsilon) E(1_{\{U_2^\epsilon\}} | x(T_1) = x_1) (1 + O(\epsilon))} \\ &= \frac{E(1_{\{Z_1\}} 1_{\{U_1^\epsilon\}})}{E(1_{\{U_1^\epsilon\}})} \times \frac{E(1_{\{Z_2\}} | U_2^\epsilon)}{E(1 | U_2^\epsilon)} (1 + O(\epsilon)) \\ &= Q(Z_1 | U_1^\epsilon) Q(Z_2 | U_2^\epsilon) (1 + O(\epsilon)), \end{aligned}$$

as claimed. \square

Pathwise Differentiation

The various payoff smoothing methods of Chapter 23 primarily target greeks computed through outright repricing with perturbed market data. However, as we have already seen in Section 3.3, there exist methods for risk calculations that avoid brute-force repricing entirely. In this chapter, we concentrate on the convenient *pathwise differentiation method*, paying particular attention to applications involving securities with barriers or early exercise rights.

24.1 Pathwise Differentiation: Foundations

24.1.1 Callable Libor Exotics

The pathwise differentiation method for European-style derivatives has been considered (in a Monte Carlo setting) in Section 3.3.2. As it turns out, Bermudan-style callable derivatives are also quite amendable to the pathwise differentiation method, as shown in Piterbarg [2004b]. Let us outline the basic ideas.

Using the notations from Chapter 18, we recall that the main valuation recursion for a CLE is given by (see (18.7))

$$H_{n-1}(T_{n-1}) = B(T_{n-1}) \mathbb{E}_{T_{n-1}} \left(B(T_n)^{-1} \max(H_n(T_n), U_n(T_n)) \right), \quad (24.1)$$

for $n = N - 1, \dots, 1$, with the starting condition $H_{N-1} \equiv 0$. Here, \mathbb{E} is the expected value operator for the spot measure Q^B , $H_n(t)$ is the n -th hold value, and $U_n(t)$ the n -th exercise value, that is, the value of all future cash flows received upon exercise at time T_n :

$$U_n(t) = B(t) \sum_{i=n}^{N-1} \mathbb{E}_t \left(B(T_{i+1})^{-1} X_i \right).$$

Here X_i are net coupons, $X_i = \tau_i(C_i - L_i(T_i))$, with C_i being the structured coupons and L_i the Libor rates.

Let Δ_α represent a pathwise differentiation operator with respect to a given parameter α . In this section we derive the main representation result that allows us to write a pathwise derivative of a callable Libor exotic as an expectation of a function of the optimal exercise time.

24.1.1.1 CLE Greeks and the Optimal Exercise Time

In order for the pathwise differentiation method to be applicable, we always assume that all coupons X_n , $n = 1, \dots, N - 1$, and the inverse numeraire $B(t)^{-1}$, are Lipschitz continuous functions of the parameter α . It follows then that the pathwise derivative $\Delta_\alpha X_n$ exists almost surely for each $n = 1, \dots, N - 1$.

From (24.1), carrying out the differentiation under the expectation operator, we obtain our first result for pathwise derivatives of CLEs.

Proposition 24.1.1. *Provided the coupons and inverse numeraire are Lipschitz continuous, then, for any n , $n = 1, \dots, N - 1$,*

$$\begin{aligned} \Delta_\alpha (B(T_{n-1})^{-1} H_{n-1}(T_{n-1})) \\ = \mathbb{E}_{T_{n-1}} (1_{\{U_n(T_n) > H_n(T_n)\}} \Delta_\alpha (B(T_n)^{-1} U_n(T_n))) \\ + \mathbb{E}_{T_{n-1}} (1_{\{H_n(T_n) > U_n(T_n)\}} \Delta_\alpha (B(T_n)^{-1} H_n(T_n))). \end{aligned}$$

Proof. The assumption of Lipschitz continuity of the coupons and the inverse numeraire implies that $U_n(T_n)$ for each n , $n = 1, \dots, N - 1$, is Lipschitz continuous in α , as is (by assumption) the inverse numeraire $B(t)^{-1}$. Since the function $\max(x, y)$ is Lipschitz continuous in x (and y), it can be shown recursively from (24.1) that $H_n(T_n)$ for each n , $n = 0, \dots, N - 1$, is Lipschitz continuous in α as well. Hence, Proposition 3.3.1 applies and we have,

$$\begin{aligned} \Delta_\alpha (B(T_{n-1})^{-1} H_{n-1}(T_{n-1})) \\ = \mathbb{E}_{T_{n-1}} (\Delta_\alpha (\max (B(T_n)^{-1} H_n(T_n), B(T_n)^{-1} U_n(T_n)))) . \end{aligned}$$

The function $x \mapsto \max(x, c)$ is absolutely continuous with a derivative that is equal to $1_{\{x>c\}}$. Hence, we can differentiate $\max(H_n(T_n), U_n(T_n))$ inside the expected value to obtain

$$\begin{aligned} \mathbb{E}_{T_{n-1}} (\Delta_\alpha \max (B(T_n)^{-1} H_n(T_n), B(T_n)^{-1} U_n(T_n))) \\ = \mathbb{E}_{T_{n-1}} (1_{\{U_n(T_n) > H_n(T_n)\}} \Delta_\alpha (B(T_n)^{-1} U_n(T_n))) \\ + \mathbb{E}_{T_{n-1}} (1_{\{H_n(T_n) > U_n(T_n)\}} \Delta_\alpha (B(T_n)^{-1} H_n(T_n))) . \end{aligned}$$

Combining equations we obtain the statement of the proposition. \square

Proposition 24.1.1 provides us with a recursive relationship (in n , the exercise date index) between $\Delta_\alpha(B(T_{n-1})^{-1} H_{n-1}(T_{n-1}))$ and

$\Delta_\alpha(B(T_n)^{-1}H_n(T_n))$. The next proposition “unwraps” this recursion to give us the formula for $\Delta_\alpha H_0$.

Proposition 24.1.2. *Let η be the optimal exercise time index (see Section 18.2.2). Then*

$$\Delta_\alpha H_0(0) = E \left(\sum_{n=\eta}^{N-1} \Delta_\alpha (B(T_n)^{-1} X_n) \right). \quad (24.2)$$

Proof. Unwrapping the recursive statement of Proposition 24.1.1, we find that

$$\begin{aligned} \Delta_\alpha H_0(0) &= \\ &\sum_{n=1}^{N-1} E \left(\left(\prod_{i=1}^{n-1} 1_{\{H_i(T_i) > U_i(T_i)\}} \right) 1_{\{U_n(T_n) > H_n(T_n)\}} \Delta_\alpha (B(T_n)^{-1} U_n(T_n)) \right). \end{aligned}$$

As η is the optimal exercise time index,

$$1_{\{\eta=n\}} = \left(\prod_{i=1}^{n-1} 1_{\{H_i(T_i) > U_i(T_i)\}} \right) 1_{\{U_n(T_n) > H_n(T_n)\}},$$

from which it follows that

$$\Delta_\alpha H_0(0) = \sum_{n=1}^{N-1} E \left(1_{\{\eta=n\}} \Delta_\alpha (B(T_n)^{-1} U_n(T_n)) \right).$$

From Proposition 3.3.1 and the fact that

$$B(T_n)^{-1} U_n(T_n) = \sum_{i=n}^{N-1} E_{T_n} (B(T_{i+1})^{-1} X_i)$$

we obtain

$$\Delta_\alpha (B(T_n)^{-1} U_n(T_n)) = \sum_{i=n}^{N-1} E_{T_n} (\Delta_\alpha (B(T_{i+1})^{-1} X_i)),$$

and therefore

$$\Delta_\alpha H_0(0) = \sum_{n=1}^{N-1} E \left(1_{\{\eta=n\}} \sum_{i=n}^{N-1} E_{T_n} (\Delta_\alpha (B(T_{i+1})^{-1} X_i)) \right).$$

The event $\{\eta = n\}$ is in the sigma-algebra \mathcal{F}_{T_n} because η is a stopping time. Thus we may carry the indicator $1_{\{\eta=n\}}$ inside the expectation E_{T_n} , to get

$$\Delta_\alpha H_0(0) = E \left(\sum_{n=1}^{N-1} \sum_{i=n}^{N-1} 1_{\{\eta=n\}} (\Delta_\alpha (B(T_{i+1})^{-1} X_i)) \right).$$

Changing the order of summation we obtain

$$\Delta_\alpha H_0(0) = E \left(\sum_{i=\eta}^{N-1} (\Delta_\alpha (B(T_{i+1})^{-1} X_i)) \right),$$

and the proposition follows. \square

The result of Proposition 24.1.2 provides the foundation for computing pathwise derivatives of callable Libor exotics, and relates the derivative of a CLE to derivatives of coupons (that can typically be computed easily) and to the optional stopping time, a quantity that is computed during normal CLE valuation anyway. To proceed further, we need to specialize the setup to either PDE or Monte Carlo based models. First, however, we study some important implications, as well as some generalizations, of Proposition 24.1.2.

24.1.1.2 Keeping the Exercise Time Constant

It is instructive to compare the expression for the value of a callable Libor exotic (18.6) with the one for its pathwise derivative in Proposition 24.1.2:

$$\begin{aligned} H_0(0) &= E \left(\sum_{n=\eta}^{N-1} B(T_{n+1})^{-1} X_n \right), \\ \Delta_\alpha H_0(0) &= E \left(\sum_{n=\eta}^{N-1} \Delta_\alpha (B(T_{n+1})^{-1} X_n) \right). \end{aligned} \tag{24.3}$$

Somewhat surprisingly, it appears that one can compute the derivative Δ_α by differentiating the sum in (24.3) and pretending that the optimal exercise time index η *does not depend on α* . But, paradoxically, in most cases the distribution of η *does depend on α* .

The seeming contradiction above can be resolved with the help of the following argument, known to economists as the *envelope theorem* (see Sydsæter and Hammond [2008]). For an arbitrary stopping time index ζ , define $V_{ki}(\zeta, X)$ by

$$V_{ki}(\zeta, X) = E \left(\sum_{n=\zeta}^{N-1} B(T_{n+1})^{-1} X_n \right),$$

where, in somewhat loose notation, X in the argument of $V_{ki}(\zeta, X)$ represents all coupons X_n and all numeraire factors $B(T_n)^{-1}$. We can think of $V_{ki}(\zeta, X)$

as the value of a knock-in barrier option with the barrier defined by ζ . Note that $V_{\text{ki}}(\zeta, X)$ is equal to $H_0(0)$ for $\zeta = \eta$. Formally differentiating with respect to α ,

$$\Delta_\alpha V_{\text{ki}}(\zeta, X) = \frac{\partial}{\partial \zeta} V_{\text{ki}}(\zeta, X) \Delta_\alpha \zeta + \frac{\partial}{\partial X} V_{\text{ki}}(\zeta, X) \Delta_\alpha X. \quad (24.4)$$

Substituting $\zeta = \eta$ into the last equation, we make a critical observation that

$$\left. \frac{\partial}{\partial \zeta} V_{\text{ki}}(\zeta, X) \right|_{\zeta=\eta} = 0, \quad (24.5)$$

because η by definition is the *optimal* stopping time index that maximizes the value of a callable Libor exotics over all stopping times (and (24.5) is the necessary first-order optimality condition). Due to (24.5), the first term in (24.4) drops out and we are left with

$$\Delta_\alpha H_0 = \Delta_\alpha V_{\text{ki}}(\eta, X) = \frac{\partial}{\partial X} V_{\text{ki}}(\eta, X) \times \Delta_\alpha X.$$

The expression on the right hand side can be interpreted as the partial derivative of the sum in (24.3) with η held constant.

The effective insensitivity of the stopping time with respect to parameter changes has some significant practical applications, even in situations where the pathwise differentiation method cannot be used (or is not used for some other reason). Recall that often a valuation of a callable security in Monte Carlo involves two steps (see Section 18.3.6): first, an optimal exercise boundary is estimated; and second, the value of the callable security is computed as a knock-in option, using the estimated exercise boundary as the barrier. In an implementation where the greeks are computed by shocking the inputs and revaluing the security, the result above states that the exercise time from the base scenario (which is, in a Monte Carlo simulation, just an integer index for each path) could be reused in the shocked scenario — i.e. we would force the exercise on a given path in a shocked scenario at exactly the same index as on the same path in the base scenario¹. Besides obvious savings in computational time (there is now no need to re-estimate the exercise boundary in the shocked scenario), this scheme improves stability of the greeks, as we explain in the next section.

We should note that if the exercise boundary being used in computations is not truly optimal (which is, of course, nearly always the case in practice), freezing the stopping times in the manner described above will change the meaning of the greeks slightly, in a manner described in Section 24.3.4. Unless the exercise rule is truly inaccurate, these differences are typically small enough to ignore. Also, we point out that a theoretically valid alternative technique involves freezing the exercise *boundary*, rather than the exercise *index*. As explained below, the latter has superior numerical properties.

¹Of course, heeding advice from Chapter 23, we should use the same seed and the same number of paths in the base and shocked scenarios.

24.1.1.3 Noise in CLE Greeks

To expand on the discussion above, and to tie it to greeks computations, let us for concreteness consider a Monte Carlo application where we attempt to evaluate CLE greeks by brute-force perturbation methods. From the results above, we have three valid alternatives when deciding how to treat the exercise decision in the perturbed market data scenario: i) re-estimate the exercise boundary (by regression, say); ii) re-use the base scenario exercise boundary; and iii) re-use the base scenario stopping times. While theoretically equivalent in the large-sample limit (due to the envelope theorem), for a realistic number of Monte Carlo paths the numerical properties of these three alternatives will differ substantially. For instance, it should be intuitively clear that re-estimating the exercise boundary will induce a large amount of spurious noise, so most practitioners have traditionally worked with a frozen boundary, as in alternative ii). Even with this approach, however, the derivatives of Bermudan-callable security prices will typically exhibit much higher levels of simulation error than prices of European options. Let us examine why this is the case.

Armed with the estimate of the optimal exercise index $\tilde{\eta}$, the Monte Carlo estimate of the value of a callable Libor exotic is given by (see Section 18.3.6)

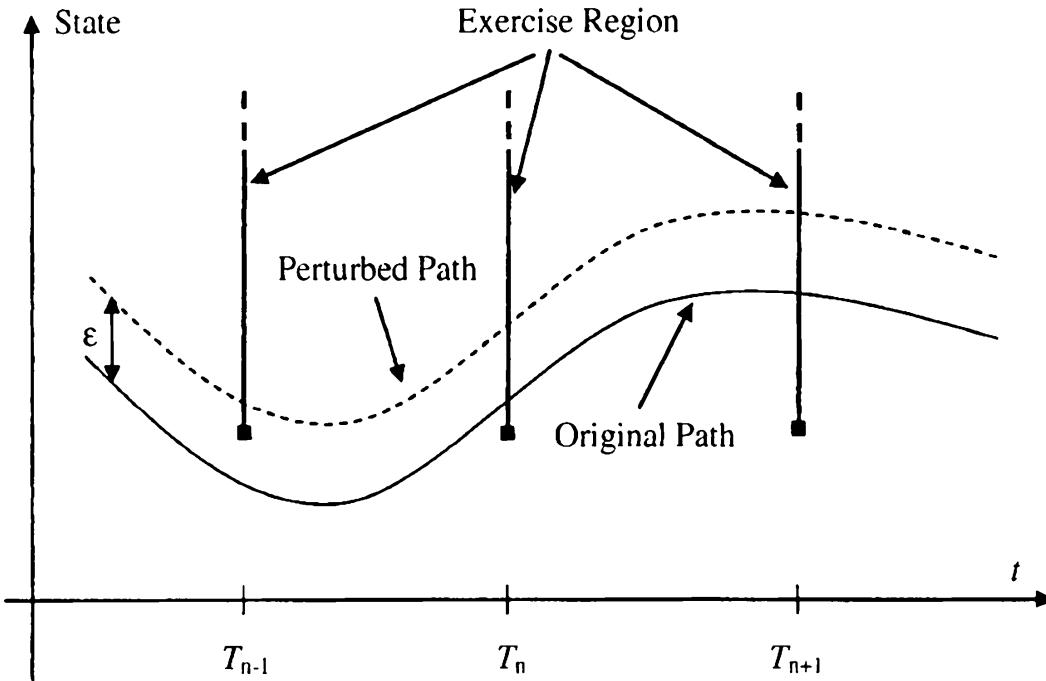
$$\tilde{H}_0 = J^{-1} \sum_{j=1}^J \sum_{i=1}^{N-1} [B(T_{i+1}, \omega_j)^{-1} X_i(\omega_j) 1_{\{i \geq \tilde{\eta}(\omega_j)\}}],$$

where simulated paths are denoted by ω_j , $j = 1, \dots, J$. Clearly, the valuation formula involves exercise indicators $1_{\{i \geq \tilde{\eta}(\omega_j)\}}$. Importantly for our analysis, these indicators are discontinuous functions of the simulated path ω . Figure 24.1 demonstrates the problem visually, for the case where the exercise boundary is frozen (our alternative ii) above).

Notice from Figure 24.1 that if a simulated path passes sufficiently close to the exercise boundary, then a small change in the parameter α can push the path outside of the exercise region for one of the exercise dates, losing a whole coupon as a result. Such a digital-type discontinuity — which is not present in European call/put or other continuous-payoff securities — leads to poorer stability and larger simulation errors for risk sensitivities in Bermudan-callable securities, compared to their European-call counterparts.

One way to improve stability and accuracy of the greeks is to use the payoff smoothing method from Section 23.4.3. However, it is much easier to use alternative iii) above, i.e. to re-use the estimate of the optimal exercise index $\tilde{\eta}(\omega)$ from the base scenario. In practice it means that for each simulated path, we just force the exercise of a CLE at exactly the same time in calculations with the shocked market data as with the base market data. In this approach, no discontinuities are introduced.

Fig. 24.1. Discontinuity of CLE Value in Monte Carlo



Notes: A whole coupon could be added to/subtracted from the value of a CLE valued in Monte Carlo under small, order- ϵ , perturbations of Monte Carlo paths if the exercise time is not kept constant.

24.1.2 Barrier Options

A CLE can be interpreted as a type of barrier option where the barrier condition is defined by the optimal exercise rule; one might therefore speculate that the pathwise differentiation method could be extended to general barrier options. This is, indeed, the case, although the presence of discontinuities in barrier options requires some additional care. As a warm-up exercise, recall the example of a pathwise derivative of digital option in Section 3.3.2.1 and consider a T -maturity European payoff

$$X = 1_{\{G > h\}} R, \quad (24.6)$$

where G and R are \mathcal{F}_T -measurable random variables and h is a particular strike. Differentiating the payoff with respect to α , we obtain

$$\Delta_\alpha (1_{\{G > h\}} R) = 1_{\{G > h\}} \Delta_\alpha R + \left(\frac{\partial 1_{\{G > h\}}}{\partial G} \right) R \Delta_\alpha G.$$

Formally,

$$\frac{\partial 1_{\{G > h\}}}{\partial G} = \delta(G - h),$$

where $\delta(x)$ is the delta function at zero. Assuming we can exchange the order of differentiation and expectation, we have

$$\begin{aligned}
\Delta_\alpha \mathbb{E}(B(T)^{-1}X) &= \mathbb{E}(\Delta_\alpha(B(T)^{-1}X)) \\
&= \mathbb{E}(\Delta_\alpha(B(T)^{-1})1_{\{G>h\}}R) \\
&\quad + \mathbb{E}(B(T)^{-1}1_{\{G>h\}}\Delta_\alpha R) \\
&\quad + \mathbb{E}(B(T)^{-1}\delta(G-h)R\Delta_\alpha G).
\end{aligned}$$

Rewriting the second term in the last equality, we find that the sensitivity of this digital option is given by

$$\begin{aligned}
\Delta_\alpha \mathbb{E}(B(T)^{-1}1_{\{G>h\}}R) &= \mathbb{E}(\Delta_\alpha(B(T)^{-1})1_{\{G>h\}}R) \\
&\quad + \mathbb{E}(B(T)^{-1}1_{\{G>h\}}\Delta_\alpha R) + \gamma_G(h)\mathbb{E}(B(T)^{-1}R\Delta_\alpha G|G=h), \quad (24.7)
\end{aligned}$$

where $\gamma_G(h)$ is the density of G , at $G = h$. While the conditions of Proposition 3.3.1 do not hold for the payoff (24.6) and we cannot rely on Proposition 3.3.1 to justify differentiation inside the expected value operator, the formula (24.7) is, nevertheless, correct, and can be justified by Malliavin calculus, see Fournie et al. [1999].

The expected values on the right-hand side of (24.7) can in general be computed in a numerical scheme such as Monte Carlo, as long as the density of G is known, and the conditional expected value $\mathbb{E}(B(T)^{-1}R\Delta_\alpha G|G=h)$ can be evaluated. Both can, in principle, be computed in Monte Carlo using Malliavin calculus techniques (see Fournie et al. [1999]). However, the application of this method is much more practical if both of these quantities can be computed (or approximated) in closed form, which is the case in some models such as, say, Gaussian models.

Proceeding to the case of barrier options, let us introduce a barrier schedule $\{T_n\}_{n=1}^{N-1}$, to which we associate knockout variables G_n and barrier levels h_n , for $n = 1, \dots, N - 1$. We consider an option that pays the value R_n , $n = 1, \dots, N - 1$, on the first T_n where $G_n > h_n$; if this event never takes place, the option pays nothing. Formally, we define η , the knockout index, by

$$\eta = \min\{k \geq 1 : G_k > h_k\} \wedge N.$$

For notational convenience, set

$$R_N \equiv 0.$$

The time 0 value of our barrier option is then given by

$$V_{ki}(0) = \mathbb{E}(B(T_\eta)^{-1}R_\eta).$$

Note that this is the same knock-in option as considered in Section 23.4.2, if we define $R_n = B(T_n) \sum_{i=n}^{N-1} B(T_{i+1})^{-1}X_i$.

More generally, let us denote

$$\eta_n = \min\{k \geq n + 1 : G_k > h_k\} \wedge N$$

and

$$V_{\text{ki},n}(t) = B(t) \mathbb{E}_t \left(B(T_{\eta_n})^{-1} R_{\eta_n} \right), \quad (24.8)$$

with the convention

$$V_{\text{ki},N}(t) \equiv 0.$$

Here, $V_{\text{ki},n}(t)$ can be seen as the value of the option with the barrier condition checked at times T_{n+1}, \dots, T_{N-1} only. In particular,

$$V_{\text{ki}}(0) = V_{\text{ki},0}(0).$$

We denote by $\gamma_n(x)$ the density of G_n at time T_n .

Proposition 24.1.3. *For the barrier option paying R_n on the first T_n where $G_n > h_n$, $n = 1, \dots, N - 1$, the pathwise derivative with respect to a parameter α is given by*

$$\begin{aligned} \Delta_\alpha V_{\text{ki}}(0) &= \mathbb{E} \left(B(T_\eta)^{-1} \Delta_\alpha R_n \Big|_{n=\eta} \right) \\ &+ \mathbb{E} \left(B(T_\eta)^{-1} \gamma_\eta(h_\eta) (R_\eta - V_{\text{ki},\eta}(T_\eta)) \Delta_\alpha G_n \Big|_{n=\eta} \Big| G_\eta = h_\eta \right). \end{aligned} \quad (24.9)$$

Proof. (Sketch) The values of the family of knock-in options defined by (24.8) satisfy the following recursive relationship,

$$\begin{aligned} B(T_n)^{-1} V_{\text{ki},n}(T_n) &= \mathbb{E}_{T_n} \left(B(T_{n+1})^{-1} R_{n+1} \mathbf{1}_{\{G_{n+1} > h_{n+1}\}} \right) \\ &+ \mathbb{E}_{T_n} \left(B(T_{n+1})^{-1} V_{\text{ki},n+1}(T_{n+1}) \mathbf{1}_{\{G_{n+1} \leq h_{n+1}\}} \right). \end{aligned}$$

Differentiating formally,

$$\begin{aligned} \Delta_\alpha (B(T_n)^{-1} V_{\text{ki},n}(T_n)) &= \mathbb{E}_{T_n} \left(\Delta_\alpha (B(T_{n+1})^{-1} R_{n+1}) \mathbf{1}_{\{G_{n+1} > h_{n+1}\}} \right) \\ &+ \mathbb{E}_{T_n} \left(\Delta_\alpha (B(T_{n+1})^{-1} V_{\text{ki},n+1}(T_{n+1})) \mathbf{1}_{\{G_{n+1} \leq h_{n+1}\}} \right) \\ &+ \mathbb{E}_{T_n} \left(B(T_{n+1})^{-1} (R_{n+1} - V_{\text{ki},n+1}(T_{n+1})) \delta(G_{n+1} - h_{n+1}) \right). \end{aligned}$$

This defines an expression of $\Delta_\alpha V_{\text{ki},n}(T_n)$ in terms of $\Delta_\alpha V_{\text{ki},n+1}(T_{n+1})$ and other quantities. Unwrapping the recursion, as in Proposition 24.1.2 earlier, proves the proposition. \square

Example 24.1.4. Callable Libor exotics are a special case with $R_n = U_n$, $G_n = U_n - H_n$, $h_n = 0$. Proposition 24.1.2 follows from Proposition 24.1.3 once continuity condition for CLEs,

$$(R_n - V_{\text{ki},n}(T_n))|_{G_n=h_n} = 0,$$

is taken into account.

Example 24.1.5. A TARN (See Chapter 20 and Section 23.4.4) can be represented in barrier form. In particular, using the notations of Section 23.4.4, we define

$$G_n \triangleq Q_n, \quad h_n \triangleq R_n, \quad R_n \triangleq \sum_{k=n}^{N-1} B(T_{k+1})^{-1} X_k(T_k),$$

so that

$$V_{\text{ki},0}(0) + V_{\text{tarn}}(0) = \mathbb{E} \left(\sum_{k=1}^{N-1} B(T_{k+1})^{-1} X_k(T_k) \right),$$

where the right-hand side equals the price of a straight (exotic) swap.

Although Proposition 24.1.3 extends the pathwise differentiation method to barrier options, the complexity of the result limits the practicality of the method. In particular, the transition densities and conditional probabilities in (24.9) are often difficult to compute, and it may ultimately be more fruitful to use methods in Chapters 23 and 25 to smooth or integrate out any discontinuities before applying pathwise differentiation techniques.

24.2 Pathwise Differentiation for PDE Based Models

The pathwise differentiation method can be applied to both PDE and Monte Carlo based models. In this section we consider PDE applications, mostly following Piterbarg [2004a]; we address Monte Carlo applications in Section 24.3.

24.2.1 Model and Setup

The treatment of European-style options in the Section 3.3.2 is rather generic and can easily be implemented in PDE-based numerical schemes. Callable Libor exotics, on the other hand, require more effort, to be undertaken in this section. While the method is developed for, and can be applied to, rather general CLEs, for a number of reasons Bermudan swaptions are probably the most natural target of the techniques described here. Indeed, not only have we shown (in Section 19.2) that low-dimensional, PDE-based Markovian models are appropriate for Bermudan swaptions, but Bermudan swaptions often constitute a dominant part of portfolios of interest rate exotics and are therefore subject to high demands for stable and accurate risk reporting.

To focus on the main features of the method without distraction from minor details, let us consider a Gaussian interest rate model as developed in Section 10.1.2.2, parameterized in terms of the short rate state $x(t)$ as in Proposition 10.1.7. We denote the infinitesimal generator associated with the dynamics of $x(t)$ by \mathcal{L} ,

$$\mathcal{L} = (y(t) - \kappa(t)x) \frac{\partial}{\partial x} + \frac{1}{2} \sigma_r(t)^2 \frac{\partial^2}{\partial x^2}. \quad (24.10)$$

If $V = V(t, x)$ is the value of a contingent claim at time t given $x(t) = x$, then $V(t, x)$ satisfies the equation (see (10.29))

$$\frac{\partial V}{\partial t}(t, x) + \mathcal{L}V(t, x) = (f(0, t) + x)V(t, x).$$

Note that this valuation expression is associated with the risk-neutral measure Q , induced by the continuous money market account $\beta(t)$. Previous material in this chapter (and in Chapter 18) used a spot Libor measure, but results carry over unchanged to the risk-neutral measure. Recycling notations, we now denote by E the corresponding, i.e. risk-neutral, expected value operator; while E was earlier used as the expected value operator in spot Libor measure, there should be no confusion which measure is used going forward.

Let us consider a CLE with net coupons $\{X_n\}$, as in Section 24.1.1; as we work in a PDE setting, we assume that the value of the net coupon X_n does not depend on the state of the yield curve prior to T_n , $n = 1, \dots, N-1$. Recall Proposition 24.1.2 which states that the pathwise derivative Δ_α of a CLE is given by

$$\Delta_\alpha H_0(0) = E \left(\sum_{n=\eta}^{N-1} \Delta_\alpha (\beta(T_{n+1})^{-1} X_n) \right),$$

where η is the optimal exercise date index. The hold and exercise values are now deterministic functions of x , so we use the notations $H_n(t, x)$, $U_n(t, x)$ where appropriate.

24.2.2 Bucketed Deltas

Arguably, the most important risk measures for an interest rate security are the so-called *bucketed interest rate deltas*, see Section 6.4, that measure the sensitivity of the value of the security to changes in various parts of the yield curve. For the CLE in question, the most natural bucketing² of deltas is induced by the tenor structure $\{T_n\}_{n=0}^N$. Specifically, we define the m -th (continuously compounded) forward rate by

$$y_m(0) \triangleq y(0, T_m, T_{m+1}) = -\frac{1}{\tau_m} \ln(P(0, T_m, T_{m+1})), \quad m = 0, \dots, N-1,$$

and denote by Δ_m the pathwise derivative with respect to $y_m(0)$,

²Naturally, the sensitivities to these rates can be projected into any other “basis”, i.e. a set of rates used to define and aggregate curve sensitivities. For the relevant techniques, see Section 6.4.3.

$$\Delta_m \triangleq \Delta_{y_m(0)}, \quad m = 1, \dots, N - 1.$$

To establish Δ_m , let us start by rewriting the pathwise derivative in a more convenient form

$$\begin{aligned} \Delta_m H_0(0) &= \sum_{n=1}^{N-1} \mathbb{E} \left(1_{\{\eta \geq n\}} \Delta_m (\beta(T_{n+1})^{-1} X_n) \right) \\ &= \sum_{n=1}^{N-1} \mathbb{E} \left(1_{\{\eta \geq n\}} \beta(T_{n+1})^{-1} \frac{\Delta_m (\beta(T_{n+1})^{-1})}{\beta(T_{n+1})^{-1}} X_n \right) \\ &\quad + \sum_{n=1}^{N-1} \mathbb{E} \left(1_{\{\eta \geq n\}} \beta(T_{n+1})^{-1} \Delta_m (X_n) \right). \end{aligned} \quad (24.11)$$

We shall also need the following lemma.

Lemma 24.2.1. *In the Gaussian model the following holds,*

$$\frac{\Delta_m (\beta(T_{n+1})^{-1})}{\beta(T_{n+1})^{-1}} = -1_{\{m \leq n\}} \tau_m.$$

Proof. We have

$$\begin{aligned} \beta(T_{n+1})^{-1} &= \exp \left(- \int_0^{T_{n+1}} r(t) dt \right) = \exp \left(- \int_0^{T_{n+1}} (f(0, t) + x(t)) dt \right), \\ \Delta_m (\beta(T_{n+1})^{-1}) &= \beta(T_{n+1})^{-1} \Delta_m \left(- \int_0^{T_{n+1}} (f(0, t) + x(t)) dt \right). \end{aligned}$$

Since the dynamics of $x(t)$ do not depend on the initial yield curve $P(0, \cdot)$,

$$\frac{\Delta_m (\beta(T_{n+1})^{-1})}{\beta(T_{n+1})^{-1}} = -\Delta_m \left(\int_0^{T_{n+1}} f(0, t) dt \right).$$

Moreover, by definition of y_k 's,

$$\int_0^{T_{n+1}} f(0, t) dt = \sum_{k=0}^n \tau_k y_k(0).$$

Hence,

$$\Delta_m \left(\int_0^{T_{n+1}} f(0, t) dt \right) = \Delta_m \left(\sum_{k=0}^n \tau_k y_k(0) \right) = 1_{\{m \leq n\}} \tau_m.$$

□

For the next result we need the following definition.

Definition 24.2.2. The time t survival measure $\Psi(\cdot; t)$ is defined for $\Gamma \subset \mathbb{R}$ by the formula

$$\Psi(\Gamma; t) = \mathbb{E} (\beta(t)^{-1} 1_{\{\eta \geq q(t)\}} 1_{\{x(t) \in \Gamma\}}),$$

where the index function $q(t)$ for the tenor structure $\{T_n\}$ is defined in (14.2). The survival density, the density $\psi(x; t)$ of the survival measure with respect to the Lebesgue measure dx , is defined by

$$\Psi(\Gamma; t) = \int_{\Gamma} \psi(y; t) dy,$$

and is assumed to exist.

Combined with this definition of the survival density and the representation (24.11), Lemma 24.2.1 allows us to derive the following representation of bucketed deltas of a CLE.

Proposition 24.2.3. In the Gaussian model, the m -th bucketed delta of a CLE is given by

$$\begin{aligned} \Delta_m H_0(0) &= - \sum_{n=m}^{N-1} \tau_m \int_{\mathbb{R}} V_{cpn,n}(x) \psi(x; T_n) dx \\ &\quad + \sum_{n=1}^m \int_{\mathbb{R}} D_{cpn,n,m}(x) \psi(x; T_n) dx, \end{aligned}$$

where $V_{cpn,n}(x)$ and $D_{cpn,n,m}(x)$ are the conditional expectations of the discounted value and the discounted derivative of the n -th coupon,

$$\begin{aligned} V_{cpn,n}(x) &= \mathbb{E} (\beta(T_n) \beta(T_{n+1})^{-1} X_n | x(T_n) = x), \\ D_{cpn,n,m}(x) &= \mathbb{E} (\beta(T_n) \beta(T_{n+1})^{-1} \Delta_m X_n | x(T_n) = x). \end{aligned}$$

Proof. From (24.11),

$$\begin{aligned} \Delta_m H_0(0) &= - \sum_{n=1}^{N-1} 1_{\{m \leq n\}} \tau_m \mathbb{E} (1_{\{\eta \geq n\}} \beta(T_{n+1})^{-1} X_n) \\ &\quad + \sum_{n=1}^{N-1} \mathbb{E} (1_{\{\eta \geq n\}} \beta(T_{n+1})^{-1} \Delta_m X_n). \end{aligned}$$

By definition of the survival density and from the fact that $q(T_n-) = n$ (see (14.2)) we obtain

$$\begin{aligned} \mathbb{E} (1_{\{\eta \geq n\}} \beta(T_{n+1})^{-1} X_n) &= \mathbb{E} (1_{\{\eta \geq q(T_n-) = n\}} \beta(T_n)^{-1} V_{cpn,n}(x(T_n))) \\ &= \int V_{cpn,n}(x) \Psi(dx; T_n) \end{aligned}$$

and

$$\begin{aligned} \mathbb{E}(1_{\{\eta \geq n\}} \beta(T_{n+1})^{-1} \Delta_m X_n) &= \mathbb{E}(1_{\{\eta \geq q(T_n, -)\}} \beta(T_n)^{-1} D_{\text{cpn},n,m}(x(T_n))) \\ &= \int D_{\text{cpn},n,m}(x) \Psi(dx; T_n). \end{aligned}$$

Since the net coupon X_n depends on the yield curve on or after the fixing time T_n ,

$$D_{\text{cpn},n,m}(x) \equiv 0$$

for $m < n$. The result follows. \square

Remark 24.2.4. The functions $V_{\text{cpn},n}(x)$ and $D_{\text{cpn},n,m}(x)$ are usually easy to calculate, as the net coupon X_n is typically a function of discount factors observed at time T_n . The reader may want to consult Section 24.3.2 where related calculations are performed.

Proposition 24.2.3 represents bucketed deltas in terms of the integrals of known (or easily computed) functions $V_{\text{cpn},n}(x)$ and $D_{\text{cpn},n,m}(x)$ against the survival density. Note that the survival density is *universal*, i.e. it does not depend on a particular delta index m . Hence, if we can calculate it efficiently, all pathwise bucketed deltas can be computed quickly, as only simple integrals are required for their calculations. This should be compared to the standard way of computing deltas, where the relevant forward rate is perturbed and the value of a CLE is recomputed by solving a full-blown PDE. We discuss the computation of the survival density in the next section.

24.2.3 Survival Density

As a reference point, let us consider the following family of measures defined on \mathbb{R} . We fix time s and position x and define, for $t \geq s$,

$$\begin{aligned} \Pi_{s,x}(t, \Gamma) &\triangleq \mathbb{E}_s(\beta(t)^{-1} 1_{\{x(t) \in \Gamma\}} \mid x(s) = x) \\ &= \mathbb{E}_s\left(e^{-\int_s^t r(u) du} 1_{\{x(t) \in \Gamma\}} \mid x(s) = x\right), \quad \Gamma \subset \mathbb{R}. \end{aligned}$$

For each s, x we can define the density³ $\pi_{s,x}(t, y)$ by

$$\Pi_{s,x}(t, \Gamma) = \int_{\Gamma} \pi_{s,x}(t, y) dy,$$

where $\pi_{s,x}(t, y)$ can be recognized as values of the Arrow-Debreu securities we introduced in Section 11.3.2.1. For fixed s, x these satisfy an analog to forward Kolmogorov equations, see (11.30), which we rewrite in our notations as

³Note that this measure density is *not* a probability density, as it does not integrate to 1.

$$\frac{\partial}{\partial t} \pi_{s,x}(t, y) = \mathcal{L}^* \pi_{s,x}(t, y) - r(t) \pi_{s,x}(t, y), \quad (s, x) \text{ fixed}, \quad (24.12)$$

for $t \geq s$. Here \mathcal{L}^* is the operator adjoint to \mathcal{L} (see (24.10)),

$$\mathcal{L}^* q(t, y) = -\frac{\partial}{\partial y} ((y(t) - \kappa(t)x) q(t, y)) + \frac{\partial^2}{\partial y^2} \left(\frac{1}{2} \sigma_r(t)^2 q(t, y) \right),$$

which is applied to $\pi_{s,x}(t, y)$ as a function of y .

The following proposition outlines an efficient procedure for computing the survival density ψ . The idea of the theorem is that in between the “interesting” times $\{T_n\}$, the density ψ behaves just like the density π in the proposition above. When the time crosses an exercise time T_n , the density ψ gets multiplied by an extra “survival” indicator function $1_{\{H_n(y, T_n) > U_n(y, T_n)\}}$.

Proposition 24.2.5. *For each n , $0 \leq n \leq N - 1$, the survival density $\psi(y; t)$ satisfies the forward PDE*

$$\frac{\partial}{\partial t} \psi(y; t) = (\mathcal{L}^* \psi)(y; t) - r(t) \psi(y; t),$$

on the time interval

$$t \in (T_n, T_{n+1}),$$

with the boundary condition

$$\psi(y; T_n) = \psi(y; T_n-) \times 1_{\{H_n(T_n, y) > U_n(T_n, y)\}}. \quad (24.13)$$

The initial condition for the first interval, $(T_0, T_1) = (0, T_1)$, is given by the delta function

$$\psi(y; 0) = \delta(y - x(0)).$$

Proof. Assume

$$T_n < t < T_{n+1},$$

so that $q(t) = n + 1$. Then

$$\begin{aligned} \Psi(\Gamma; t) &= E \left(e^{- \int_0^t r(u) du} 1_{\{\eta \geq n+1\}} 1_{\{x(t) \in \Gamma\}} \right) \\ &= E \left(e^{- \int_0^{T_n} r(u) du} 1_{\{\eta > n\}} E_{T_n} \left(e^{- \int_{T_n}^t r(u) du} 1_{\{x(t) \in \Gamma\}} \right) \right) \\ &= E \left(e^{- \int_0^{T_n} r(u) du} 1_{\{\eta > n\}} \Pi_{T_n, x(T_n)}(t, \Gamma) \right). \end{aligned}$$

From this formula we obtain

$$\psi(y; t) = E \left(e^{- \int_0^{T_n} r(u) du} 1_{\{\eta > n\}} \pi_{T_n, x(T_n)}(t, y) \right).$$

Differentiating this equality with respect to t , exchanging the order of differentiation and taking the expectation, applying (24.12) and exchanging

the order of the linear operator $\mathcal{L}^* - r(t)$ and the expectation operator, we obtain that the same equation as (24.12) holds for $\psi(y; t)$,

$$\frac{\partial}{\partial t} \psi(y; t) = \mathcal{L}^* \psi(y; t) - r(t) \psi(y; t)$$

for $t \in (T_n, T_{n+1})$. To derive boundary conditions we notice that

$$\begin{aligned}\Psi(\Gamma; T_n) &= E(\beta(T_n)^{-1} 1_{\{\eta \geq n+1\}} 1_{\{x(T_n) \in \Gamma\}}) \\ &= E(\beta(T_n)^{-1} 1_{\{\eta \geq n\}} 1_{\{H_n(T_n, x(T_n)) > U_n(T_n, x(T_n))\}} 1_{\{x(T_n) \in \Gamma\}}).\end{aligned}$$

As $x(t)$ and $\beta(t)$ are continuous at $t = T_n$ we have

$$\Psi(\Gamma; T_n) = \Psi(\Gamma \cap \{H_n(T_n, x(T_n)) > U_n(T_n, x(T_n))\}; T_n-)$$

and, calculating the densities of both sides, we obtain (24.13). For more details see Piterbarg [2004a]. \square

Remark 24.2.6. The time- T_n conditions (24.13) require knowledge of the “hold” regions

$$\{x \in \mathbb{R} : H_n(T_n, x) > U_n(T_n, x)\}.$$

These are computed as a by-product of the CLE valuation, since on each exercise date T_n , the hold values $H_n(T_n, x)$ are determined as functions of the state process $x(t)$ evaluated at time $t = T_n$.

Proposition 24.2.5 outlines a procedure for computing the survival density in one forward PDE “sweep”, starting at $t = 0$ with a delta function. The solution is computed forward using an appropriate PDE scheme (see e.g. Sections 11.3.2.1, 11.3.2.2) until the first exercise time T_1 . At this point, the solution (i.e. the survival density) is multiplied by the indicator function of the no-exercise condition. The density is then rolled forward again until the next exercise date where it is multiplied by another no-exercise indicator function, and so on.

The pathwise differentiation method for calculating deltas handily outperforms the standard approach of re-evaluation of a derivative under shocked scenarios. For the pathwise differentiation method, to calculate all N bucketed deltas, we need to calculate *one* survival density at a cost comparable to one PDE valuation of the derivative, and $2N$ integrals of Proposition 24.2.3, at a combined cost of about twice the PDE valuation of the derivative. In contrast, the perturb-and-revalue method would require N PDE valuations, one for each bucketed delta. As N is typically significantly larger than 3, the cost savings therefore can be quite significant. Nor is the pathwise method limited to deltas only; as shown in Piterbarg [2004a], one can handle vegas and gammas in the same way. As we can reuse much of the calculations (the survival density) among all these greeks, performance improvements are even more dramatic.

24.3 Pathwise Differentiation for Monte Carlo Based Models

Let us now consider applications of pathwise differentiation to Monte Carlo based models. For concreteness, we develop the technique for the LM model (14.13)–(14.14) with separable deterministic local volatility:

$$dL_n(t) = \varphi(L_n(t)) \lambda_n(t)^\top (\mu_n(t) dt + dW(t)), \quad (24.14)$$

$$\mu_n(t) = \sum_{j=q(t)}^n \frac{\tau_j \varphi(L_j(t)) \lambda_j(t)}{1 + \tau_j L_j(t)},$$

The basic principles are, however, quite generic and straightforward to apply to other models. With LM models more naturally presented in the spot measure Q^B , we use a setup where the numeraire is chosen to be the rolling money market $B(t)$, see (14.8). For notational convenience, we assume that the LM model and the security to be priced share the same tenor structure $\{T_n\}$; in practice, of course, this need not be the case.

24.3.1 Pathwise Derivatives of Forward Libor Rates

The discrete and spanning nature of forward Libor rates makes the definition of bucketed deltas⁴ easy, and we define Δ_m to be the pathwise derivative with respect to $L_m(0)$, $m = 0, \dots, N - 1$,

$$\Delta_m X \triangleq \frac{\partial X}{\partial L_m(0)}$$

for any random variable X . Note that in order to keep notation light, we reuse the definition Δ_m for pathwise derivatives from Section 24.2, but redefine their meaning slightly, as we here calculate derivatives with respect to simply compounded rates, rather than to the continuously compounded rates used in Section 24.2.

As should be clear from the basic discussion in Section 3.3.2, to successfully apply the pathwise differentiation method to a Libor market model, we need to be able to simulate the pathwise derivatives of the forward Libor rates $\Delta_m L_n(t)$, $n, m = 0, \dots, N - 1$. To determine the Q^B -dynamics of $\Delta_m L_n(t)$, we use the standard technique of differentiating the SDEs for $L_n(t)$. From (14.13)–(14.14), differentiating with respect to $L_m(0)$, we get

$$d(\Delta_m L_n(t)) = \varphi'(L_n(t)) \lambda_n(t)^\top \Delta_m L_n(t) (\mu_n(t) dt + dW(t)) \\ + \varphi(L_n(t)) \lambda_n(t)^\top \sum_j \frac{\partial \mu_n(t)}{\partial L_j(t)} \Delta_m L_j(t) dt. \quad (24.15)$$

⁴As we mentioned before in Section 6.4, once the deltas in a particular “basis” are computed, it is a matter of simple linear algebra to re-express them in any other basis, e.g. the one used by a risk management system.

The initial conditions for these SDEs are found by differentiating the initial conditions for $L_n(t)$'s, resulting in

$$\Delta_m L_n(0) = 1_{\{n=m\}}. \quad (24.16)$$

The system of SDEs given by (24.15) and (14.13)–(14.14) fully specifies the dynamics of the forward Libor rates and their pathwise derivatives through time.

There are N equations in the system (14.13) and N^2 equations in the system (24.15), and simulating all is computationally expensive, even for relatively low values of N . A significant part of the numerical effort originates with drift computations, so it is natural to investigate whether simplifications of the drift term in (24.15) can lighten the computational burden. Glasserman and Zhao [1999] propose to use the following simplified system of SDEs for simulating values and pathwise deltas of forward Libor rates,

$$\begin{aligned} dL_n(t) &= \varphi(L_n(t)) \lambda_n(t)^\top (\mu_n(t) dt + dW(t)), \\ d(\Delta_m L_n(t)) &= \varphi(L_n(t)) \lambda_n(t)^\top \Delta_m L_n(t) (\mu_n(0) dt + dW(t)) \\ &\quad + \varphi'(L_n(t)) \lambda_n(t)^\top \sum_j \frac{\partial \mu_n(0)}{\partial L_j(0)} \Delta_m L_j(t) dt. \end{aligned} \quad (24.17)$$

Notice that we here retain the original Libor rate dynamics, but have applied the standard “freezing” technique to the drifts when calculating the dynamics of pathwise derivatives of Libor rates. This allows for a considerable speed-up, as the drifts in the equations for deltas of forward Libor rates can be pre-computed before the simulation. Glasserman and Zhao [1999] show numerically that the loss of accuracy in (24.17) is typically quite small.

The cost of propagating pathwise derivatives of Libor rates often dominates the computations of pathwise deltas, so let us consider computational complexity in more detail. We denote by $\Delta\mathbf{L}(t)$ an $N \times N$ matrix with the (n, m) -th element equal to $\Delta_m L_n(t)$,

$$\Delta\mathbf{L}(t)_{n,m} = \frac{\partial L_n(t)}{\partial L_m(0)}, \quad n, m = 0, \dots, N - 1. \quad (24.18)$$

To fix ideas, we assume that we need to propagate $\Delta\mathbf{L}(t)$ for $0 \leq t \leq T = T_k$ for some k , and that we discretize (24.15) using the Euler scheme over the time grid $\{T_i\}_{i=0}^k$ of the LMM tenor structure. Further assume that a path of the Brownian motion $W(t)$ has been drawn, and we have denoted $Z_{i-1} = (W(T_i) - W(T_{i-1})) / \sqrt{\tau_{i-1}}$. Then, for any time step i , $i = 1, \dots, k$, we can rewrite (24.15) in matrix form⁵

$$\Delta\mathbf{L}(T_i) = \mathbf{D}(T_{i-1}) \Delta\mathbf{L}(T_{i-1}), \quad i = 1, \dots, k, \quad (24.19)$$

⁵We could also use the faster approximation (24.17) here; we leave relevant modifications for the reader to explore.

where the matrix $\mathbf{D}(T_{i-1})$ has elements

$$\begin{aligned}\mathbf{D}(T_{i-1})_{n,m} &= 1_{\{n=m\}} + \varphi(L_n(T_{i-1}))\lambda_n(T_{i-1})^\top \frac{\partial \mu_n(T_{i-1})}{\partial L_m(T_{i-1})}\tau_{i-1} \\ &+ 1_{\{n=m\}}\varphi'(L_n(T_{i-1}))\lambda_n(T_{i-1})^\top (\mu_n(T_{i-1})\tau_{i-1} + Z_{i-1}\sqrt{\tau_{i-1}}).\end{aligned}\quad (24.20)$$

We see that propagating $\Delta \mathbf{L}(t)$ over one time step requires a matrix-matrix multiplication of order $O(N^3)$, so the calculation of $\Delta \mathbf{L}(T_k)$ (which requires k steps) has total computation effort of order $O(kN^3)$.

It is interesting to compare the computational complexity of this algorithm to a brute-force perturbation method. Let us estimate the cost of calculating N deltas to $L_m(0)$, $m = 0, \dots, N-1$, by shocking each of these forward Libor rates and re-running the simulation (24.14). Stepping one Libor rate in one perturbation scenario over one time step costs $O(1)$. Hence, stepping all N Libor rates in all N scenarios over all k time steps has complexity $O(kN^2)$, i.e. it is *faster*, by a factor of $O(N)$, than propagating $\Delta \mathbf{L}(t)$. On the other hand, once $\Delta \mathbf{L}(t)$ is simulated, it could be reused for calculating deltas of multiple payoffs, a point we return to in Section 24.3.3. It follows that a naive implementation of the pathwise differentiation method only becomes competitive speed-wise when there are more than N payoffs to differentiate in the same simulation. This seemingly limits the usability of the pathwise method, as simultaneous calculation of risk sensitivities for multiple payoffs is often difficult to achieve in practice, since most risk systems treat each trade as a separate work unit⁶. In Section 24.3.3 below, we show that by suitably rearranging the order of calculations in the pathwise differentiation method, the computational cost can be brought down to $O(kN^2)$, making pathwise differentiation computationally competitive with the revaluation method.

Given that the computational effort is no better than for the significantly simpler perturbation-based methods, the reader may wonder whether the pathwise differentiation method is ultimately worth the effort. The answer to this question is not always entirely obvious. On one hand, pathwise differentiation produces a true derivative estimate without a difference coefficient bias (see Section 3.3.1) and, in a sense, can be seen as the ultimate way of “geometry fixing” for Monte Carlo (or PDEs), since greeks are calculated in exactly the same simulation as the base value. Recalling the analysis in Section 23.1, it is therefore not surprising that pathwise differentiation often produces greeks of superior quality to those produced by run-of-the-mill perturbation methods. On the other hand, by scrutinizing a given product in detail and carefully “locking” all relevant computational details, it is often possible to construct perturbation methods that produce

⁶Also, as calibration and time-discretization is normally set up in a product-specific manner, it can be awkward (and even sub-optimal) to attempt to price many securities simultaneously in a single Monte Carlo loop.

greeks of comparable quality to those of the pathwise differentiation — CLEs are good examples of this, as discussed in Section 24.1.1.2.

24.3.2 Pathwise Deltas of European Options

As in Section 14.6.2.1, let $\mathbf{L}(t)$ be the vector of all Libor rates, and consider a European-style option with time T payoff $V(\mathbf{L}(T))$, for a deterministic function $V(\mathbf{x})$, $\mathbf{x} = (x_0, \dots, x_{N-1})$. The option value at time 0 equals

$$V(\mathbf{L}(0)) = \mathbb{E}(B(T)^{-1}V(\mathbf{L}(T))),$$

where \mathbb{E} denotes expectations in the spot measure Q^B . As required by Proposition 3.3.1, we suppose that $V(\mathbf{x})$ is a Lipschitz continuous function of \mathbf{x} . Then the pathwise delta Δ_m can be carried under the expectation operator,

$$\Delta_m \mathbb{E}(B(T)^{-1}V(\mathbf{L}(T))) = \mathbb{E}(\Delta_m(B(T)^{-1}V(\mathbf{L}(T)))) ,$$

so that

$$\begin{aligned} \Delta_m V(\mathbf{L}(0)) &= \mathbb{E}(\Delta_m(B(T)^{-1})V(\mathbf{L}(T))) \\ &\quad + \mathbb{E}\left(B(T)^{-1} \sum_{i=0}^{N-1} \frac{\partial V(\mathbf{x})}{\partial x_i} \Big|_{\mathbf{x}=\mathbf{L}(T)} \Delta_m L_i(T)\right). \end{aligned}$$

To compute deltas of the option, we need to be able to compute the deltas of the numeraire $\Delta_m(B(T)^{-1})$, as well as the derivatives of the payoff $\partial V(\mathbf{x})/\partial x_i$. We start with the numeraire.

24.3.2.1 Pathwise Deltas of the Numeraire

Recall that the discrete money market account $B(T)$ is given by

$$B(T) = \left(\prod_{i=0}^n (1 + \tau_i L_i(T_i)) \right) P(T, T_{n+1}),$$

where we have assumed that $T_n \leq T < T_{n+1}$. The pathwise derivative of the stub bond $P(T, T_{n+1})$ will depend (mildly) on the interpolation scheme used in the model, see Section 15.1. To keep the exposition simple, we choose the zero-volatility interpolation for the front stub $P(T, T_{n+1})$ of Section 15.1.4, whereby

$$P(T, T_{n+1}) = P(T_n, T, T_{n+1}).$$

Applying constant interpolation of simply compounded rates, see (15.4), we arrive at

$$P(T, T_{n+1}) = P(T_n, T, T_{n+1}) = \frac{1}{1 + (T_{n+1} - T) L_n(T_n)},$$

so that

$$B(T) = \left(\prod_{i=0}^{n-1} (1 + \tau_i L_i(T_i)) \right) \frac{1 + \tau_n L_n(T_n)}{1 + (T_{n+1} - T) L_n(T_n)}.$$

Differentiating, we obtain

$$\Delta_m (B(T)^{-1}) \quad (24.21)$$

$$= \sum_{j=0}^n \frac{\partial (B(T)^{-1})}{\partial L_j(T_j)} \Delta_m L_j(T_j) \quad (24.22)$$

$$= -B(T)^{-1} \sum_{j=0}^{n-1} \frac{\tau_j}{1 + \tau_j L_j(T_j)} \Delta_m L_j(T_j) \\ - B(T)^{-1} \frac{T - T_n}{(1 + (T_{n+1} - T) L_n(T_n)) (1 + \tau_n L_n(T_n))} \Delta_m L_n(T_n).$$

24.3.2.2 Pathwise Deltas of the Payoff

A typical (Lipschitz continuous) interest rate payoff $V(\mathbf{x})$ can be represented as an absolutely continuous function, say $f(\cdot)$, of one (or more) Libor or CMS rates. Therefore, the pathwise derivatives of this payoff with respect to initial Libor rates $L_m(0)$, $m = 0, \dots, N - 1$, can be computed by a chain rule, as long as we know how to differentiate market rates with respect to $L_m(0)$, $m = 0, \dots, N - 1$. For instance, for some swap rate $S(t)$ (note that t is not necessarily equal to T), we get

$$\Delta_m V(\mathbf{L}(t)) = \Delta_m f(S(t)) = f'(S(t)) \Delta_m S(t).$$

The derivatives $\Delta_m S(t)$ are determined by the way the yield curve at future time t is constructed from simulated primary Libor rates $\mathbf{L}(t)$, as discussed in detail in Section 15.1. In particular, the rate $S(t)$ is always a known function of zero-coupon bonds $P(t, s)$ for various s , so $\partial S(t)/\partial P(t, s)$ is easily computed. Finally, as we have an algorithm to construct all $P(t, s)$ from $\mathbf{L}(t)$ per Section 15.1, we can calculate $\partial P(t, s)/\partial L_m(t)$ along the same lines as in (24.22).

For rates $S(t)$ that are aligned with the tenor structure $\{T_n\}$ of the model (i.e. $S(t) = S_{i,j}(t)$ for some i, j as defined in (4.10)), calculations simplify significantly, and we have already derived relevant derivatives in Section 14.4.2, see (14.31). Other methods from Section 14.4.2 also apply for general rates $S(t)$; in particular, we can recycle ideas from Section 14.4.2 on swap rate volatility approximations used for calibrating LM models. Recall the freezing idea of Proposition 14.4.3,

$$\frac{\partial S(t)}{\partial L_m(t)} \frac{\varphi(L_m(t))}{\varphi(S(t))} \approx \frac{\partial S(0)}{\partial L_m(0)} \frac{\varphi(L_m(0))}{\varphi(S(0))}.$$

By simple algebraic manipulations we obtain that

$$\Delta_m S(t) = \frac{\partial S(T)}{\partial L_m(T)} \approx \Delta_m S(0) \frac{\varphi(L_m(0))}{\varphi(S(0))} \frac{\varphi(S(T))}{\varphi(L_m(T))}. \quad (24.23)$$

Numerical errors arising from the approximation (24.23) are typically small, and performance gains are significant as the quantities $\partial S(0)/\partial L_m(0)$ can now be pre-computed before the simulation starts.

It is worth mentioning at this point that, while a good part of the discussion above was about deltas, other first-order risk sensitivities such as vegas and, even, second-order sensitivities such as gammas could be computed in a pathwise method as well, as briefly discussed in Section 3.3.2.

24.3.3 Adjoint Method For Greeks Calculation

Let us continue contemplating hedge computations for European options paying at time T some function V of the Libor vector $\mathbf{L}(T)$; for notational convenience, define $U \triangleq B(T)^{-1}V(\mathbf{L}(T))$. Once the values of pathwise derivatives of all forward Libor rates $\Delta_m L_n(T)$, $n, m = 0, \dots, N - 1$, are simulated for a given path (using, for example, (24.15) or (24.17)), then the full set of pathwise deltas $\Delta U \triangleq (\Delta_0 U, \dots, \Delta_{N-1} U)$ can be calculated at a small cost of multiplying a vector by a matrix (note the row-vector form of left- and right-hand sides),

$$\Delta U = \frac{\partial U}{\partial \mathbf{L}(T)} \Delta \mathbf{L}(T), \quad (24.24)$$

where

$$\frac{\partial U}{\partial \mathbf{L}(T)} = \left(\frac{\partial U}{\partial L_0(T)}, \dots, \frac{\partial U}{\partial L_{N-1}(T)} \right)$$

is payoff-specific (but often easy to calculate, see Sections 24.3.2.1 and 24.3.2.2), and $\Delta \mathbf{L}(t)$ is an $N \times N$ matrix with the (n, m) -th element equal to $\Delta_m L_n(t)$, see (24.18).

As we already mentioned in Section 24.3.1, the representation (24.24) is convenient if we want to calculate pathwise deltas of multiple payoffs simultaneously, as the matrix $\Delta \mathbf{L}(T)$ can be reused for each payoff. On the other hand, the calculation of the matrix $\Delta \mathbf{L}(T)$ is computationally costly — as we showed in Section 24.3.1, it is of the order $O(N)$ slower than just calculating deltas by revaluation, and if we only need to calculate pathwise deltas of a *single* payoff, it is not clear why one would ever want to use the pathwise differentiation method. However, it turns out that by rearranging the order of calculations in what is known as the *adjoint method* (see Giles

and Glasserman [2006]), the speed of calculations in the pathwise method can be significantly improved.

It follows from (24.24), (24.19) and (24.16) that

$$\Delta U = \frac{\partial U}{\partial \mathbf{L}(T_k)} \mathbf{D}(T_{k-1}) \dots \mathbf{D}(T_0), \quad (24.25)$$

where the matrices $\mathbf{D}(T_i)$ are defined in (24.20). The standard pathwise differentiation method calculates matrices $\mathbf{D}(T_0)$, $\mathbf{D}(T_1)\mathbf{D}(T_0)$, and so on, using a matrix-matrix multiplication on each step and ultimately multiplying the final matrix by the vector $\partial U / \partial \mathbf{L}(T_k)$. We can, however, rearrange the order of calculations so that on each step we have a *vector-matrix* multiplication. To accomplish this, we just need to group the terms in (24.25) “from the left”:

$$\Delta U = \left(\dots \left(\left(\frac{\partial U}{\partial \mathbf{L}(T_k)} \mathbf{D}(T_{k-1}) \right) \mathbf{D}(T_{k-2}) \right) \dots \right) \mathbf{D}(T_0).$$

In particular, let us define

$$Y^k = \frac{\partial U}{\partial \mathbf{L}(T_k)}$$

and then, recursively,

$$Y^{i-1} = Y^i \mathbf{D}(T_{i-1}), \quad i = k, \dots, 1. \quad (24.26)$$

Then Y^0 gives the final solution,

$$Y^0 = \Delta U,$$

after applying the recursion (24.26) k times. Each step involves a vector-matrix multiplication and requires only $O(N^2)$ operations — i.e. savings of a factor of N compared to the standard pathwise scheme (see Section 24.3.1) — as is clear from both (24.26) and from the following explicit representation obtained from (24.20):

$$\begin{aligned} Y_m^{i-1} &= Y_m^i + \left(\sum_n Y_n^i \varphi(L_n(T_{i-1})) \lambda_n(T_{i-1})^\top \frac{\partial \mu_n(T_{i-1})}{\partial L_m(T_{i-1})} \right) \tau_{i-1} \\ &\quad + Y_m^i \varphi'(L_m(T_{i-1})) \lambda_m(T_{i-1})^\top (\mu_m(T_{i-1}) \tau_{i-1} + Z_{i-1} \sqrt{\tau_{i-1}}) \end{aligned} \quad (24.27)$$

for $m = 0, \dots, N - 1$. The computational effort is further reduced by noting that this expression simplifies significantly for some combinations of the indices i, m, n . For instance,

$$\lambda_m(T_{i-1}) = 0, \quad \mu_m(T_{i-1}) = 0 \quad \text{for } m \leq i - 1$$

in line with our conventions $L_m(t) \equiv L_m(T_m)$ for $t \geq T_m$. Also, in the spot Libor measure, the drift derivatives $\partial\mu_n(T_{i-1})/\partial L_m(T_{i-1})$ are non-zero only for $i \leq m \leq n$, and similar conditions exist for drifts in other measures. All these facts could (and should) be used to obtain an efficient numerical implementation.

The recursion (24.26) proceeds backward in time, but as is clear from (24.27) the i -th step requires the (simulated) value of the Libor vector $\mathbf{L}(T_{i-1})$, which can only be obtained in a forward simulation. This is not much of a problem, however, as we can always save the required values of the Libor rates when calculating the *value* of the option in the (forward) simulation (14.13)–(14.14) and then use these rates in the backward recursion (24.26), (24.27) when calculating deltas. The extra memory requirements are modest as this is done path-by-path.

From the discussion in this section it should be clear that when the pathwise differentiation method is used, there is limited downside to using the adjoint method to arrange the calculation order. An exception occurs if one is able to compute risk on more than $O(N)$ derivatives in the same model, on a time line shared by all products in the same simulation. In this case, the Libor delta matrix $\Delta \mathbf{L}$ should be pre-computed and applied to each payoff via (24.24).

24.3.4 Pathwise Delta Approximation for Callable Libor Exotics

Calculations of pathwise deltas for CLEs can be based on the fundamental result of Proposition 24.1.2 that expresses the pathwise derivative of a CLE in terms of pathwise derivatives of the (net) coupons. Conveniently, as the coupons can be regarded as European options, the results of the previous section can be used to compute pathwise deltas of coupons. Per Proposition 24.1.2, we additionally require an estimate $\tilde{\eta}$ of the optimal exercise index, which fortunately is almost always found as a by-product of a typical Monte Carlo valuation of a callable security, see Section 18.3. Once $\tilde{\eta}$ is obtained, the (lower bound) estimate of the value of the CLE is given by

$$\tilde{H}_0(0) = \mathbb{E} \left(\sum_{n=\tilde{\eta}}^{N-1} B(T_{n+1})^{-1} X_n \right), \quad (24.28)$$

as computed in a Monte Carlo simulation.

Replacing the true exercise index η with its estimator $\tilde{\eta}$ gives an approximation of the value of a callable Libor exotic. In the same vein, replacing in Proposition 24.1.1 the true optimal exercise index η with $\tilde{\eta}$ gives an approximation of the pathwise delta,

$$\tilde{\Delta}_m H_0(0) \triangleq \mathbb{E} \left(\sum_{n=\tilde{\eta}}^{N-1} \Delta_m (B(T_{n+1})^{-1} X_n) \right), \quad m = 0, \dots, N-1. \quad (24.29)$$

It is shown in Piterbarg [2004b] that, as the exercise policy estimate converges to the optimal policy, the estimate in (24.29) approaches the true pathwise delta,

$$\tilde{\Delta}_m H_0(0) \rightarrow \Delta_m H_0(0).$$

This gives rise to an elegant formula for estimating deltas of a callable Libor exotic that is easy to implement in practice. With the estimate of the optimal exercise time, $\tilde{\eta}$, coming “for free” from the estimation step of the Monte Carlo valuation, we approximate the pathwise delta by

1. Running a forward simulation, for each path ω determining the exercise time index $\tilde{\eta}(\omega)$.
2. For each path, computing pathwise deltas of all coupons X_n , $n = 1, \dots, N - 1$ (as well as the deltas of the inverse numeraire $B(t)^{-1}$), per Section 24.3.2.
3. Adding up deltas $\Delta_m(B(T_{n+1})^{-1}X_n)$ for those coupons that occur after the exercise index $\tilde{\eta}(\omega)$.
4. Averaging the result over all paths.

As the pathwise delta of a CLE is given by the sum of pathwise deltas of European-style options, it follows trivially that the adjoint method of Section 24.3.3 could fruitfully be used here as well — an idea discussed at length in Leclerc et al. [2009].

We call the values $\tilde{\Delta}_m H_0(0)$, $m = 0, \dots, N - 1$, *pathwise delta approximations*. These should not be confused with the true deltas of the lower bound CLE price estimate. To state this more succinctly, recall the definition of $\tilde{H}_0(0)$ in (24.28), which can be interpreted as the value of a barrier-style Libor exotic that knocks into $U_n(T_n)$ for the first n for which the approximate exercise region (characterized by $\tilde{\eta}$) is hit. Then, we generally have,

$$\tilde{\Delta}_m H_0(0) \neq \Delta_m \tilde{H}_0(0),$$

where on the right-hand side we have the true delta of the lower bound CLE price estimate.

It can easily be shown that under mild regulatory conditions,

$$\Delta_m \tilde{H}_0(0) \rightarrow \Delta_m H_0(0),$$

as the exercise boundary converges to the optimal one. Hence, both approximations $\tilde{\Delta}_m H_0(0)$ and $\Delta_m \tilde{H}_0(0)$ provide converging approximations to the true delta $\Delta_m H_0(0)$. We note that it is normally $\Delta_m \tilde{H}_0(0)$ that is typically computed in the standard perturbation method. Piterbarg [2004b] compares deltas computed by perturbations and by pathwise differentiation, and finds the latter both more stable and significantly faster to compute: in the tests performed, pathwise delta approximations required about 15 times

less computational effort than delta computations by direct perturbation methods⁷.

As both $\tilde{\Delta}_m H_0(0)$ and $\Delta_m \tilde{H}_0(0)$ converge to the same value when the exercise policy approaches optimality, we can use the difference between the two as an informal measure of the quality of our exercise decision approximation (or, equivalently, the gap between the true value and the lower bound value calculated in Monte Carlo). In practice this works best if we aggregate all deltas together, and monitor the difference

$$\sum_{m=0}^{N-1} \tilde{\Delta}_m H_0(0) - \sum_{m=0}^{N-1} \Delta_m \tilde{H}_0(0)$$

for significant deviations from 0.

24.4 Notes on Likelihood Ratio and Hybrid Methods

Section 3.3.3 introduced another non-perturbative differentiation method, the likelihood ratio method. The method shifts differentiation from the payoff to the density of the process and is not limited to smooth (Lipschitz continuous) payoffs. Practical applications of likelihood ratio methods in interest rate modeling are typically limited to fairly special situations, so we do not here expand much on our introduction in Section 3.3.3. Still, it is instructive to understand *why* the likelihood ratio method in its basic form is not particularly useful for our purposes⁸.

We start by recalling the expression for the log-likelihood ratio (3.80) in the Black-Scholes model. Of particular relevance to our discussion is the presence of \sqrt{T} term in the denominator of the expression for the log-likelihood ratio in (3.80). Clearly, with T approaching zero, the log-likelihood ratio grows to infinity, resulting in exploding variance of the estimate of the likelihood ratio derivative (3.79). In general, it is, in fact, not the time to option expiry that determines how fast the variance of the estimate explodes, but the earliest observation date of the underlying asset process. To demonstrate, we consider a security with payoff $g(Y(T_1), \dots, Y(T_N))$ for

⁷Note that we are here comparing against a brute-force perturbation method where the exercise boundary — rather than the exercise *time* — is kept fixed under perturbations. Had we instead kept the exercise time fixed in perturbed scenarios, we would likely have obtained greeks of quality comparable to those produced by the pathwise method. Recall our comments at the end of Section 24.3.1.

⁸Another potential drawback is the need to know the transition density of the underlying process, although one always has the option of using a Gaussian approximation based on an Euler discretization of the true process. As discussed in Chen and Glasserman [2007b], the limit of this procedure for small time steps is deeply connected to the Malliavin calculus.

some $0 < T_1 < \dots < T_N$, with the process $Y(t)$ defined in Section 3.3.3.1. Then, clearly

$$\mathbb{E}(g(Y(T_1), \dots, Y(T_N))) = \mathbb{E}(\tilde{g}(Y(T_1))),$$

where

$$\tilde{g}(y) = \mathbb{E}(g(Y(T_1), \dots, Y(T_N)) | Y(T_1) = y).$$

Hence,

$$\begin{aligned} \frac{d}{dS_0} \mathbb{E}(g(Y(T_1), \dots, Y(T_N))) &= \frac{d}{dS_0} \mathbb{E}(\tilde{g}(Y(T_1))) \\ &= \mathbb{E}(l(Y(T_1)) \mathbb{E}(g(Y(T_1), \dots, Y(T_N)) | Y(T_1))) \\ &= \mathbb{E}(l(Y(T_1)) g(Y(T_1), \dots, Y(T_N))) \\ &= \frac{1}{S_0 \sigma \sqrt{T_1}} \mathbb{E}(Z_1 g(Y(T_1), \dots, Y(T_N))), \end{aligned}$$

where $Z_1 = W(T_1)/\sqrt{T_1} \sim \mathcal{N}(0, 1)$. The time T_1 here could, for example, be the time to the first coupon fixing date, to the first exercise date of a CLE, or to the first knockout date of a barrier. Because of the regular structure of most interest rate derivatives, the time T_1 will in most cases be rather short, resulting in high variance of the estimate.

The fact that the likelihood ratio method does not work for many interest rate derivatives is unfortunate, since, as described in Section 3.3.3, likelihood ratio methods have the potential of handling irregular (e.g., discontinuous) payouts that are outside the scope of pathwise differentiation and perturbation methods. For such payouts, we will often have to apply one of the payoff smoothing methods from Chapter 23, and *then* apply pathwise differentiation or a perturbation method. An alternative approach involves invoking a *hybrid* method, that attempts to combine features of both the pathwise differentiation and likelihood ratio methods. In a series of works Fries and Kampen [2006], Fries and Joshi [2008a], Fries [2007], the authors have introduced successively more elaborate hybrid schemes that, roughly speaking, attempt to choose the right combination of a pathwise and a likelihood ratio derivative for each Monte Carlo path, depending on the relationship between the path and “interesting” product features such as strikes or barriers. We cannot possibly do justice to all the nuances involved in developing these schemes, so we simply refer interested readers to the source papers. Many of these methods are both fairly involved and rather specialized, so their deployment in generic risk systems will often be challenging. It is fair to say that the jury is still out when it comes to the practicality of these schemes in actual trading systems.

Importance Sampling and Control Variates

Even if sophisticated payoff smoothing and pathwise derivative schemes are employed, obtaining high-quality Monte Carlo greeks will always require the statistical simulation error to be kept low. Several generic variance reduction techniques were already introduced in Chapter 3; here, we expand on certain applications that are of particular relevance in interest rate modeling. As it turns out, some techniques, such as importance sampling techniques of Section 25.2, produce benefits for greeks estimation that go beyond mere variance reduction. On the other hand, other techniques are less impressive for the greeks than for basic value estimation, as described in Section 25.6. Nevertheless, all variance reduction techniques in this chapter are useful to know.

In our discussion here, we first study a number of applications of the importance sampling technique originally introduced in Section 3.4.4, with a particular emphasis on barrier and TARN products. Subsequently, we turn our attention to the control variate method initially considered in Section 3.4.3, discussing a variety of model- and instrument-based strategies for finding good controls.

25.1 Importance Sampling In Short Rate Models

We first look at a classic importance sampling application in simulation of short rate models. For concreteness, let us consider the pricing of a zero-coupon bond maturing at time T in the generic model (11.54); i.e. we are interested in evaluating

$$X(0, T) = \mathbb{E} \left(\exp \left(- \int_0^T x(u) du \right) \right) \triangleq \mathbb{E}(Y(T)) \quad (25.1)$$

by Monte Carlo methods. Notice that in (25.1) the expectation \mathbb{E} is assumed taken under the risk-neutral probability measure Q . We consider now chang-

ing probability measure, from Q to some other measure \tilde{Q} , with the measure change characterized by a density $\varsigma(t) = E_t(d\tilde{Q}/dQ)$ with

$$d\varsigma(t) = -\varsigma(t)q(t, x(t)) dW(t), \quad \varsigma(0) = 1,$$

for some function $q(t, x(t))$ sufficiently regular for $\varsigma(t)$ to be a Q -martingale. By the Radon-Nikodym theorem

$$X(0, T) = \tilde{E}(Y(T)/\varsigma(T)),$$

where \tilde{E} is the expected value operator for measure \tilde{Q} .

In measure \tilde{Q} , Girsanov's theorem tells us that the joint process for $x(t)$, $1/\varsigma(t)$, and $Y(t)$ becomes

$$\begin{aligned} d \begin{pmatrix} x(t) \\ 1/\varsigma(t) \\ Y(t) \end{pmatrix} &= \begin{pmatrix} \mu_x(t, x(t)) - q(t, x(t))\sigma_x(t, x(t)) \\ 0 \\ -x(t)Y(t) \end{pmatrix} dt \\ &\quad + \begin{pmatrix} \sigma_x(t, x(t)) \\ q(t, x(t))/\varsigma(t) \\ 0 \end{pmatrix} d\tilde{W}(t), \end{aligned}$$

where $d\tilde{W}(t) = dW(t) + q(t, x(t))dt$ is a \tilde{Q} -Brownian motion.

As shown in Section 3.4.4.3, we can arrange for the random variable $Y(T)/\varsigma(T)$ to have zero variance in measure \tilde{Q} , provided that we use (3.96) to set

$$\begin{aligned} \varsigma(t) &= \exp \left(- \int_0^t x(u) du \right) X(t, T, x(t))/X(0, T), \\ q(t, x(t)) &= -X(t, T, x(t))^{-1} \sigma_x(t, x(t)) \frac{\partial X(t, T, x(t))}{\partial x(t)}, \end{aligned} \quad (25.2)$$

where

$$X(t, T, x) = E_{t,x} \left(\exp \left(- \int_t^T x(u) du \right) \right).$$

For the SDE (11.54) we generally do not have an analytical (reconstitution) expression for $X(t, T, x)$, but we are free to provide a guess for it. While doing so will most likely not reduce the variance of $Y(T)/\varsigma(T)$ to zero, if the guess is at all reasonable we can still expect a significant variance reduction effect. One route to an estimate for the function $X(t, T, x)$ is to assume that the SDE for $x(t)$ can be approximated by a simpler SDE for which a closed-form bond reconstitution formula exists; possible candidates would be, say, the affine class of short rate models or the quadratic Gaussian model. For instance suppose that we feel that the SDE for $x(t)$ can be approximated with a mean-reverting Gaussian model

$$dx(t) \approx (m - \kappa_G x(t)) dt + \sigma dW(t),$$

then we would obtain

$$q(t, x) = \sigma_x(t, x) \frac{1 - e^{-\kappa_G(T-t)}}{\kappa_G}. \quad (25.3)$$

In practice, most models will have a linear mean-reverting drift term, so the estimate of the “best” choice of κ_G in (25.3) is often straightforward. If $\mu_x(t, x)$ is non-linear, we could simply linearize it around $x = 0$ for the purpose of estimating κ_G .

Andersen [1996] (see also Andersen and Boyle [2000]) tests the efficiency of the choice (25.3) when applied to the problem of computing discount bond prices for the CIR process; the results are far superior to those obtained by traditional variance reduction techniques. Andersen [1996] also notes that the quality of the measure transformation method improves significantly as the number of time steps in the simulation path is increased; this behavior is not surprising given that the method has been designed around the continuous-time limit of the discretized process for $x(t)$. The tendency of the measure transform method to improve with increasing number of discretization steps is quite attractive as it complements the behavior of the bias in the SDE discretization scheme: increasing the number of time steps will lower *both* the systematic bias and the random Monte Carlo error.

Finally, we note that the principles at play in the method above are general and can be applied to more complicated securities than discount bonds; all that is required is some decent estimate of the expectation value as a function of t and x . Often such estimates can be derived — either exactly or at least approximately — in a Gaussian model, for instance. We note that knowledge of an expectation in a closely related model can also form the basis for an application of the control variate method, an idea that we discuss in more detail starting from Section 25.3 below.

25.2 Payoff Smoothing by Importance Sampling in TARNs and General Barrier Options

Let us now take a different tack and demonstrate that the importance sampling method may also be used to produce payoff smoothing, in the vein of Chapter 23.

25.2.1 Binary Options

We first study a simple example that clearly illustrates the connection between importance sampling and payoff smoothing. Let X be a Gaussian random variable with mean μ and variance σ^2 , and consider an option that pays $g(X)$ (for some smooth $g(x)$) if X is below a certain barrier b , so that the value of the security is given by

$$V = \mathbb{E}(g(X)1_{\{X < b\}}), \quad (25.4)$$

where \mathbb{E} is an expected value operator for some pricing measure P (note that we do not include discounting in this illustrative example for clarity). Valuing this security by Monte Carlo requires simulating independent Gaussian samples, discarding those that end up above the barrier b , and averaging the payoff values over non-discarded samples. If b is low, then the proportion of paths that contribute to the average is small, which leads to a large simulation error (see related discussion in Section 3.4.4.5). Also, the digital feature in the payoff reduces the accuracy and stability of Monte Carlo estimates of greeks. In light of this, it seems natural to change the probability measure to increase the proportion of “interesting” samples, as we did in Section 3.4.4.5. Alternatively, it is tempting to integrate the digital option analytically. As we shall show, the importance sampling method can be set up to implement ‘both’ strategies.

Let us rewrite the value by conditioning on the survival,

$$V = \mathbb{E}(g(X)|X < b) P(X < b). \quad (25.5)$$

The probability of survival in our simple example is known in closed form,

$$P(X < b) = \Phi\left(\frac{b - \mu}{\sigma}\right),$$

where $\Phi(z)$ is the standard Gaussian CDF. Calculating the remaining term $\mathbb{E}(g(X)|X < b)$ by Monte Carlo simulation requires us to draw random samples of X , conditioned on the event $\{X < b\}$. In order to do this, let us briefly recall from Section 3.1.1 how Gaussian random variables are typically simulated. If U is a uniform random variable on $[0, 1]$, then X is obtained by

$$X = \Phi^{-1}(U). \quad (25.6)$$

Therefore, X conditioned on $\{X < b\}$ can be sampled by simply drawing a random variable U' uniformly distributed on the interval $[0, \Phi(b)]$, followed by an application of the mapping (25.6):

$$X| \{X < b\} = \Phi^{-1}(U'), \quad U' \sim \mathcal{U}(0, \Phi(b)). \quad (25.7)$$

From (25.5), we may write our option value as

$$V = \mathbb{E}(g(\Phi^{-1}(U'))) \Phi\left(\frac{b - \mu}{\sigma}\right), \quad U' \sim \mathcal{U}(0, \Phi(b)). \quad (25.8)$$

Here, the function $g(x)$ is smooth by assumption, and thus a Monte Carlo evaluation of $\mathbb{E}(g(\Phi(U')))$ will have good convergence and exhibit stable greek estimates. By conditioning, we have, in effect, managed to integrate out the discontinuity analytically, and used the Monte Carlo method for the smooth part of the payoff only.

While a close connection of the method above to the payoff smoothing methods of Chapter 23 is obvious, the method can also be interpreted as a particular case of importance sampling, since in (25.8) all drawn samples come from the “interesting” part of the sample space where survival is guaranteed. The measure effectively used for sampling in (25.8) is often known as the *survival measure*. To characterize this measure further, notice first that in (25.4) the variable $1_{\{X < b\}}$ is not strictly positive, which requires some additional considerations before using this variable to define a measure shift. Indeed, starting from Section 1.3, we so far have only considered *equivalent* measures defined by strictly non-zero random variables as Radon-Nikodym derivatives. The definition of measure change can, however, be extended to Radon-Nikodym derivatives which can hit zero, but in this case the new measure \tilde{P} is not equivalent to the original measure P ; instead it is *absolutely continuous* with respect to the original measure:

$$P(A) = 0 \Rightarrow \tilde{P}(A) = 0,$$

but not necessarily the other way around. Notice that the two measures are equivalent when restricted to the set on which the Radon-Nikodym derivative is strictly positive. This set is what we are interested in here. The specific Radon-Nikodym derivative we need is given by

$$\Lambda = \frac{1_{\{X < b\}}}{P(X < b)};$$

note the normalization factor so that $E(\Lambda) = 1$. With this definition, we may write

$$V = E(g(X)1_{\{X < b\}}) = E(g(X)\Lambda) P(X < b) = \tilde{E}(g(X)) P(X < b),$$

where \tilde{E} denotes expectation in the survival measure \tilde{P} , defined by

$$\frac{d\tilde{P}}{dP} = \Lambda.$$

\tilde{P} is called a survival measure because it assigns zero probability to all events in the “no-survival” region, i.e. for any event A such that $A \subset \{X \geq b\}$, we have $\tilde{P}(A) = 0$. The distribution of X under \tilde{P} coincides with the distribution of X conditioned on $\{X < b\}$, and is given by (25.7).

The example above is simple, but it demonstrates a general approach to smoothing via importance sampling. Even for more complex barrier-style options, conditioning on survival will often allow us to handle discontinuities analytically, and evaluate the smooth part of the payoff by sampling under the survival measure. These ideas are fully developed in Glasserman and Staum [2001], where the authors observe that conditioning on full survival for a general barrier option is usually not analytically tractable, and instead

propose to condition on one-step survival, from one barrier observation date to the next; at each time step, the measure is changed locally to allow the process to survive until the next time step. Since the behavior of most processes is much simpler on shorter time scales than on longer ones, this strategy will often lead to analytical tractability and efficient Monte Carlo implementation. Our treatment of TARNs in the next section is based on these ideas.

25.2.2 TARNs

TARNs and their valuation by Monte Carlo have been introduced in Chapter 20. We recall that the main TARN valuation formula (20.2) under the spot measure Q^B reads

$$V_{\text{tarn}}(0) = \mathbb{E} \left(\sum_{n=1}^{N-1} B(T_{n+1})^{-1} X_n(T_n) 1_{\{Q_n < R\}} \right), \quad (25.9)$$

where $B(t)$ is as the discretely rolled money market numeraire, X_n 's are net coupons, $Q_n = \sum_{i=1}^{n-1} \tau_i C_i$ are accumulated structured coupons, and R is the total return. Here \mathbb{E} is (re-)defined to be the expected value operator for measure Q^B . As described in detail elsewhere, the net coupon X_n is paid only if the process survives up to time T_n ; by analogy to the simple example in Section 25.2.1, we expect that conditioning on survival may reduce variance and improve risk stability. We have already seen a payoff smoothing method applied to TARNs based on partial analytical integration in Section 23.2.4. Here we approach the problem from a different angle.

25.2.3 Removing the First Digital

In many cases, the biggest contributor to the simulation noise in a TARN is the first embedded digital option, i.e. the contract feature that specifies a knock-out event at date T_2 if $C_1(T_1)$ is above a certain barrier¹. The variance of the estimate can be reduced if we could handle this digital option explicitly, outside of the Monte Carlo simulation.

To develop the idea in detail, let us for concreteness focus on a TARN of the inverse floating type, where the structured coupon is as in (20.1),

$$C_n = (s - g \times L_n(T_n))^+. \quad (25.10)$$

We also introduce a sequence of random variables

$$b_n = (s - (R - Q_n) / \tau_n) / g, \quad (25.11)$$

¹Sometimes a TARN is structured so that the first digital is virtually worthless, but the second one is important. The discussion that follows should then be modified accordingly.

with b_n being $\mathcal{F}_{T_{n-1}}$ -measurable. The first variable in the sequence,

$$b_1 = (s - R/\tau_1)/g$$

(see (20.4)) is deterministic, and we have

$$\{Q_2 < R\} \Leftrightarrow \{L_1(T_1) > b_1\}.$$

Let us denote by \mathcal{V} the path value of the coupons that depend on the first knockout event (the first coupon $X_1(T_1)$ is paid always and is easy to handle separately),

$$\mathcal{V} = \sum_{n=2}^{N-1} B(T_{n+1})^{-1} X_n(T_n) 1_{\{Q_n < R\}}.$$

Then

$$\begin{aligned} E(\mathcal{V}) &= E(\mathcal{V} | L_1(T_1) > b_1) Q^B(L_1(T_1) > b_1) \\ &\quad + E(\mathcal{V} | L_1(T_1) \leq b_1) Q^B(L_1(T_1) \leq b_1). \end{aligned}$$

Clearly

$$E(\mathcal{V} | L_1(T_1) \leq b_1) = 0$$

so that

$$E(\mathcal{V}) = E(\mathcal{V} | L_1(T_1) > b_1) Q^B(L_1(T_1) > b_1). \quad (25.12)$$

In (25.12), since T_1 is typically small (less than one year), the probability $Q^B(L_1(T_1) > b_1)$ of not knocking out can nearly always be approximated analytically with a high degree of precision. For instance, since time to expiry is short, the issue of non-deterministic drift of $L_1(T_1)$ under the spot measure can be easily dealt with by, say, freezing the drift along the forward yield curve (see related discussion in Section 23.2.4).

The value $E(\mathcal{V}|L_1(T_1) > b_1)$ can be interpreted as the value of the TARN under the condition that it does not knock out on the date T_1 . This value can be computed in a Monte Carlo simulation by either sampling Gaussian variates that generate simulation steps by a scheme similar to (25.7), or by adjusting the drifts of the forward Libor model in such a way as to move the Libor rate L_1 away from the knockout region. We do not go into details as we will present a more general scheme in the next section.

25.2.4 Smoothing All Digitals by One-Step Survival Conditioning

Removing the first discontinuity from the payoff being calculated by Monte Carlo often reduces the simulation error substantially. However, we can go further. Typically, given the information available on the coupon date T_n , we can evaluate the probability of knockout on the next day (quasi)-analytically. Following Piterbarg [2004c], we can use this information to develop a scheme where *all* discontinuities are integrated outside of Monte Carlo.

Proposition 25.2.1. *A TARN with structured coupon (25.10) can be valued as follows,*

$$V_{\text{tarn}}(0) = \sum_{n=1}^{N-1} \tilde{\mathbb{E}} (\psi_n \mathbb{E}_{T_{n-1}} (B(T_{n+1})^{-1} X_n(T_n))) , \quad (25.13)$$

$$\psi_n = \prod_{k=1}^{n-1} Q_{T_{k-1}}^B (L_k(T_k) > b_k) . \quad (25.14)$$

Here the measure \tilde{Q}^B is defined by its Radon-Nikodym derivative with respect to Q^B ,

$$\Lambda(t) = \mathbb{E}_t \left(\frac{d\tilde{Q}^B}{dQ^B} \right) , \quad (25.15)$$

where $\Lambda(t)$ is a non-negative, normalized Q^B -martingale such that

$$\Lambda(t) = \frac{Q_t^B (L_{m+1}(T_{m+1}) > b_{m+1})}{Q_{T_m}^B (L_{m+1}(T_{m+1}) > b_{m+1})} \prod_{k=1}^{m-1} \frac{1_{\{L_{k+1}(T_{k+1}) > b_{k+1}\}}}{Q_{T_k}^B (L_{k+1}(T_{k+1}) > b_{k+1})}$$

for $t \in [T_m, T_{m+1})$.

Proof. We observe that due to non-negativity of the structured coupon C_{n-1} , the following equality holds Q^B -almost surely,

$$\begin{aligned} \{Q_n < R\} &\Leftrightarrow \left\{ Q_{n-1} < R, (s - g L_{n-1}(T_{n-1}))^+ < (R - Q_{n-1}) / \tau_{n-1} \right\} \\ &\Leftrightarrow \{Q_{n-1} < R, L_{n-1}(T_{n-1}) > b_{n-1}\} . \end{aligned}$$

Likewise, using non-negativity of all C_i 's,

$$\begin{aligned} \{Q_n < R\} &\Leftrightarrow \{Q_1 < R, Q_2 < R, \dots, Q_n < R\} \\ &\Leftrightarrow \{L_1(T_1) > b_1, \dots, L_{n-1}(T_{n-1}) > b_{n-1}\} . \end{aligned}$$

Hence

$$1_{\{Q_n < R\}} = \prod_{k=1}^{n-1} 1_{\{L_k(T_k) > b_k\}} .$$

Define

$$\Lambda_n(t) = \begin{cases} \frac{Q_t''(L_{n+1}(T_{n+1}) > b_{n+1})}{Q_{T_n}''(L_{n+1}(T_{n+1}) > b_{n+1})}, & t \in [T_n, T_{n+1}), \\ \frac{1_{\{L_{n+1}(T_{n+1}) > b_{n+1}\}}}{Q_{T_n}''(L_{n+1}(T_{n+1}) > b_{n+1})}, & t \geq T_{n+1}, \\ 1, & t < T_n . \end{cases}$$

We note that $\Lambda_n(t)$ is a non-negative Q^B -martingale. Moreover, $\Lambda_n(t)$ is constant on $[0, T_n]$ and $[T_{n+1}, \infty)$. In addition, $\Lambda_n(t)$ is $\mathcal{F}_{T_{n+1}}$ -measurable for $t \geq T_{n+1}$.

We define $\Lambda(t)$ by

$$\Lambda(t) = \prod_{n=0}^{N-2} \Lambda_n(t).$$

It is not hard to show that $\Lambda(t)$ is a Q^B -martingale as well. Let us denote the value of the n -th coupon, contingent on survival to time T_n , by

$$V_{\text{cpn},n}(0) = E(B(T_{n+1})^{-1} X_n(T_n) 1_{\{Q_n < R\}}).$$

As

$$\begin{aligned} \prod_{k=1}^{n-1} 1_{\{L_k(T_k) > b_k\}} &= \left(\prod_{k=1}^{n-1} \frac{1_{\{L_k(T_k) > b_k\}}}{Q_{T_{k-1}}^B(L_k(T_k) > b_k)} \right) \left(\prod_{k=1}^{n-1} Q_{T_{k-1}}^B(L_k(T_k) > b_k) \right) \\ &= \left(\prod_{k=1}^{n-1} \Lambda_{k-1}(T_{n-1}) \right) \left(\prod_{k=1}^{n-1} Q_{T_{k-1}}^B(L_k(T_k) > b_k) \right) \\ &= \left(\prod_{k=1}^{N-1} \Lambda_{k-1}(T_{n-1}) \right) \left(\prod_{k=1}^{n-1} Q_{T_{k-1}}^B(L_k(T_k) > b_k) \right), \end{aligned}$$

we have

$$V_{\text{cpn},n}(0) = E \left(\Lambda(T_{n-1}) B(T_{n+1})^{-1} X_n(T_n) \prod_{k=1}^{n-1} Q_{T_{k-1}}^B(L_k(T_k) > b_k) \right). \quad (25.16)$$

Next, taking the $\mathcal{F}_{T_{n-1}}$ -conditional expected value inside the expected value in (25.16) and using the measure \tilde{Q}^B defined by its Radon-Nikodym derivative with respect to Q^B in (25.15), we obtain

$$\begin{aligned} V_{\text{cpn},n}(0) &= \tilde{E} \left(\psi_n E_{T_{n-1}} \left(B(T_{n+1})^{-1} X_n(T_n) \right) \right), \\ \psi_n &= \prod_{k=1}^{n-1} Q_{T_{k-1}}^B(L_k(T_k) > b_k), \end{aligned}$$

and the proposition follows. \square

Remark 25.2.2. Quantities ψ_n in (25.14) can be calculated or approximated analytically since each term of the form $Q_{T_{k-1}}^B(L_k(T_k) > b_k)$ involves an expected value over a relatively short period $[T_{k-1}, T_k]$, so that short-time approximations (e.g., based on a Gaussian distribution) to the distribution of L_k over $[T_{k-1}, T_k]$ may be applied effectively. Note that the fact that L_k is a Q^B -martingale over time period $[T_{k-1}, T_k]$ would help here. In some simulation schemes, such as those considered in Sections 25.2.5 and 25.2.6 for example, the ψ_n come for free, without any extra work.

Remark 25.2.3. The measure \tilde{Q}^B is not equivalent to Q^B because $\Lambda(t)$ can be zero. However, since the value of the TARN is zero for those paths for which $\Lambda(t)$ is zero, \tilde{Q}^B and Q^B are equivalent on the relevant subspace of the sample space.

The formula (25.13) specifies that the value of a TARN can be computed by Monte Carlo simulation under the measure \tilde{Q}^B by adding values of net coupons X_n scaled by weights ψ_n . This should be contrasted to the original expression (25.9) where the weights on coupons are instead indicator functions $1_{\{Q_n < R\}}$. Obviously, the ψ_n 's are much smoother functions of a simulated path than are indicator functions, since in the former the digital discontinuities have been integrated away by computing the probabilities $Q_{T_{k-1}}^B(L_k(T_k) > b_k)$ in (25.14) (quasi)-analytically.

Another feature of note of formula (25.13) is the presence of nested expected values where the inner one is calculated under the original spot measure Q^B , while the outer uses the survival measure \tilde{Q}^B . While \tilde{Q}^B is the main simulation measure, when computing the value of the n -th coupon we must return to measure Q^B , when stepping from time T_{n-1} to time T_{n+1} .

In measure \tilde{Q}^B , the TARN never knocks out, so a simulation based on the result in Proposition 25.2.1 can be interpreted as a version of the importance sampling method in Section 25.2.1, where the measure is changed from Q^B to \tilde{Q}^B and the likelihood ratio is partially pre-integrated. Of course, to use the method in practice we need to establish how precisely to simulate model dynamics in measure \tilde{Q}^B ; we study this topic in the next three sections.

25.2.5 Simulating Under the Survival Measure Using Conditional Gaussian Draws

We first consider a special case of the Libor market model (see Chapter 14), where we use a single-factor² volatility specification with separable deterministic local volatility, see (14.13)–(14.14). Continuing with the inverse floater TARN example and assuming that the tenor structure coincides with the schedule of the TARN, Libor forwards satisfy

$$dL_i(t) = \lambda_i(t)\varphi(L_i(t))(\mu_i(t)dt + dW(t)), \quad (25.17)$$

$i = 1, \dots, N - 1$, with $W(t)$ being a one-dimensional Brownian motion under the spot measure Q^B .

TARN valuation with formula (25.13) requires us to simulate Libor rates under measure \tilde{Q}^B , i.e. in such a way that $L_n(T_n) > b_n$ for each n . To see how this would work, let us consider a simulation time step from T_{n-1} to T_n for a fixed n for all Libor rates. We note that for each n , b_n in (25.11) is $\mathcal{F}_{T_{n-1}}$ -measurable, i.e. is known at time T_{n-1} . Employing a simple Euler scheme, we can approximate the Q^B -dynamics as

²We comment on the multi-factor case later.

$$L_n(T_n) = L_n(T_{n-1}) + \lambda_n(T_{n-1}) \varphi(L_n(T_{n-1})) (\mu_n(T_{n-1})\tau_{n-1} + \sqrt{\tau_{n-1}}Z),$$

where Z is a standard Gaussian random variable. Given that we want to simulate in such a way that $L_n(T_n) > b_n$, we need to make sure that Z satisfies

$$L_n(T_{n-1}) + \lambda_n(T_{n-1}) \varphi(L_n(T_{n-1})) (\mu_n(T_{n-1})\tau_{n-1} + \sqrt{\tau_{n-1}}Z) > b_n, \quad (25.18)$$

which can be solved to yield

$$Z > Z_{\min}, \quad Z_{\min} \triangleq (b_n - m_n)/v_n, \quad (25.19)$$

where we have denoted

$$v_n \triangleq \lambda_n(T_{n-1}) \varphi(L_n(T_{n-1})) \sqrt{\tau_{n-1}}, \quad m_n \triangleq L_n(T_{n-1}) + v_n \mu_n(T_{n-1}) \sqrt{\tau_{n-1}}, \quad (25.20)$$

so that

$$L_n(T_n) = m_n + v_n Z. \quad (25.21)$$

The lower bound Z_{\min} in (25.19) is known at time T_{n-1} , and the measure change is expressed by the requirement that the random variable Z should satisfy (25.19). In Section 25.2.1 we have already discussed how to simulate a Gaussian random variable conditioned on it being below (or above) a certain level; all we need to do is to apply the idea behind the scheme (25.7) (with b set to Z_{\min} from (25.19)). In particular we can just set

$$\tilde{Z} = \Phi^{-1}(\Phi(Z_{\min}) + (1 - \Phi(Z_{\min}))U), \quad (25.22)$$

where U is a uniform draw from $[0, 1]$.

While in this new measure we, by construction, have that $L_n(T_n) > b_n$, it may not be entirely obvious that this is the measure \tilde{Q}^B as defined by (25.15), since we can satisfy the constraint (25.18) in many different ways. To check, let us denote the measure implicit in the simulation scheme above by \hat{Q}^B for a moment. We obviously have that for any $l < b_n$,

$$\hat{Q}_{T_{n-1}}^B(L_n(T_n) > l) = 1, \quad l < b_n. \quad (25.23)$$

For l such that $l > b_n$, we have

$$\hat{Q}_{T_{n-1}}^B(L_n(T_n) > l) = Q_{T_{n-1}}^B(\tilde{Z} > (l - m_n)/v_n)$$

and then, from (25.22),

$$Q_{T_{n-1}}^B(\tilde{Z} > (l - m_n)/v_n) = \frac{1 - \Phi((l - m_n)/v_n)}{1 - \Phi(Z_{\min})}.$$

Now, from (25.20)–(25.21),

$$1 - \Phi((l - m_n)/v_n) = Q_{T_{n-1}}^B(L_n(T_n) > l),$$

$$1 - \Phi(Z_{\min}) = Q_{T_{n-1}}^B(L_n(T_n) > b_n),$$

and we finally obtain

$$\widehat{Q}_{T_{n-1}}^B(L_n(T_n) > l) = \frac{Q_{T_{n-1}}^B(L_n(T_n) > l)}{Q_{T_{n-1}}^B(L_n(T_n) > b_n)}, \quad l \geq b_n$$

which, together with (25.23), demonstrates that \widehat{Q}^B is the same measure as \widetilde{Q}^B defined by (25.15).

To finish the description of the simulation scheme we note that once \tilde{Z} has been drawn, all Libor rates $L_i(t)$, $i = n, \dots, N - 1$, can be evaluated using

$$L_i(T_n) = L_i(T_{n-1}) + \lambda_i(T_{n-1}) \varphi(L_i(T_{n-1})) \left(\mu_i(T_{n-1}) \tau_{n-1} + \sqrt{\tau_{n-1}} \tilde{Z} \right),$$

for $i = n, \dots, N - 1$.

Notice that the algorithm above not only shows how to easily propagate the Libor curve forward in time, it also gives us the weights ψ_n in (25.14) without extra work. Specifically, from (25.19) we see that the one-step survival probability $Q_{T_{n-1}}^B(L_n(T_n) > b_n)$ is simply equal to $1 - \Phi(Z_{\min})$. It should be noted that the simplicity of the algorithm is partly based on the fact that we consider only a single-factor LM model in (25.17), and also by the fact that the payout is such that we can express the survival condition as a simple condition on one of the primary Libor rates over each time period. Both restrictions can, however, be lifted fairly easily. For instance, Pietersz [2005] suggests using a suitable rotation of the local volatility matrix to make sure that the survival (over a given time step) is determined by a single Gaussian draw. A different, and more general, twist is offered in Fries and Joshi [2008b]; we briefly review this approach in the next section.

25.2.6 Generalized Trigger Products in Multi-Factor LM Models

Following Fries and Joshi [2008b], we define a *generalized trigger product* to be a contract that pays a (net) coupon X_n until a knockout event, defined as the first time index n where an \mathcal{F}_{T_n} -measurable *trigger variable* G_n exceeds some *trigger level* h_n , i.e. when $G_n > h_n$, with h_n being $\mathcal{F}_{T_{n-1}}$ -measurable. In the spot measure, the so-defined security has present value

$$V_{\text{gtp}}(0) = E \left(\sum_{n=1}^{\eta-1} B(T_{n+1})^{-1} X_n \right), \quad \eta = \min \{n \geq 1 : G_n \geq h_n\} \wedge N.$$

A generalized trigger product is closely linked to barrier options we considered in Sections 23.4.2 and 24.1.2. TARNs are a special case; for a TARN the

trigger variable is in fact the n -th structured coupon C_n and the trigger barrier is given by $h_n = (R - Q_{n-1})/\tau_n$, see (20.2). Note that we do not assume any particular form for G_n or h_n at this point.

Leaning on the results in Proposition 25.2.1, the value $V_{\text{gtp}}(0)$ can be rewritten as an expectation in the survival measure \tilde{Q}^B ,

$$V_{\text{gtp}}(0) = \sum_{n=1}^{N-1} \tilde{\mathbb{E}} (\psi_n \mathbb{E}_{T_{n-1}} (B(T_{n+1})^{-1} X_n)), \quad (25.24)$$

$$\psi_n = \prod_{k=1}^{n-1} Q_{T_{k-1}}^B (G_k \leq h_k). \quad (25.25)$$

Let us now generalize (25.17). We assume that all Libor rates are driven by a d -dimensional Brownian motion,

$$dL_i(t) = \sigma_i(t)^\top (\mu_i(t) dt + dW(t)), \quad i = 1, \dots, N-1,$$

where $\sigma_i(t)$ is a general process that may depend on (potentially all) Libor rates at time t . Denoting by $\mathbf{L}(t)$ the vector of all forward Libor rates observed at t , we rewrite the dynamics in a vector format

$$d\mathbf{L}(t) = M(t) dt + \Sigma(t)^\top dW(t), \quad (25.26)$$

for a suitably defined vector function $M(t)$ and a matrix function $\Sigma(t)$ (both of which are functions of $\mathbf{L}(t)$).

Let us consider a single time step from T_{n-1} to T_n , and assume that the n -th trigger variable G_n is in fact a function $G_n(\mathbf{L}(T_n))$ of the vector of Libor rates observed at T_n . An Euler scheme for (25.26) is given by

$$\mathbf{L}(T_n) = \mathbf{L}(T_{n-1}) + M(T_{n-1}) \tau_{n-1} + \hat{\Sigma}^\top Z, \quad (25.27)$$

where Z is a d -dimensional standard Gaussian vector and $\hat{\Sigma} = \Sigma(T_{n-1})^\top \sqrt{\tau_{n-1}}$. Let us define by $\gamma(z)$ the value of the trigger variable G_n as a function of the realized Gaussian increment in the Euler scheme (25.27),

$$\gamma(z) \triangleq G_n (\mathbf{L}(T_{n-1}) + M(T_{n-1}) \tau_{n-1} + z),$$

so that

$$G_n(\mathbf{L}(T_n)) = \gamma(\hat{\Sigma}^\top Z). \quad (25.28)$$

Next, we define the normalized gradient (a row vector) of the function γ by

$$v = \nabla \gamma(0) / \|\nabla \gamma(0)\|.$$

The survival boundary is given in terms of the function γ as $\gamma(z) = h_n$. Let Y_{\max} be the solution of the linearization of this equation in direction v , i.e. let us set

$$Y_{\max} = (h_n - \gamma(0)) / \|\nabla \gamma(0)\|. \quad (25.29)$$

Then, to first order, as follows from (25.28), the survival condition $G_n(\mathbf{L}(T_n)) < h_n$ is equivalent to the following condition on the Gaussian draw Z ,

$$Y < Y_{\max}, \quad Y \triangleq v \widehat{\Sigma}^{\top} Z. \quad (25.30)$$

In (25.30) the random variable Y is Gaussian, making it straightforward to design a sampling scheme where (25.30) is always satisfied. Drawing on the same idea that lead to (25.22) and (25.7), we define

$$U = \Phi(Y/\sigma_Y), \quad \sigma_Y^2 = v \widehat{\Sigma}^{\top} \widehat{\Sigma} v^{\top},$$

and also set

$$\tilde{Y} = \sigma_Y \Phi^{-1}(U \Phi(Y_{\max})). \quad (25.31)$$

Clearly $\tilde{Y} \leq Y_{\max}$ always, and

$$Q^B(Y < K) = \Phi(Y_{\max}) Q^B(\tilde{Y} < K) \quad (25.32)$$

for any $K \in (-\infty, \Phi(Y_{\max})]$. In particular, to first order,

$$\begin{aligned} \gamma\left(\widehat{\Sigma}^{\top} Z + v^{\top}(\tilde{Y} - Y)\right) &\approx \gamma(0) + \|\nabla \gamma(0)\| v\left(\widehat{\Sigma}^{\top} Z + v^{\top}(\tilde{Y} - Y)\right) \\ &= \gamma(0) + \|\nabla \gamma(0)\| (Y + \tilde{Y} - Y) \\ &= \gamma(0) + \|\nabla \gamma(0)\| \tilde{Y} \end{aligned}$$

and, since $\tilde{Y} \leq Y_{\max}$, we have that (again to first order)

$$\gamma\left(\widehat{\Sigma}^{\top} Z + v^{\top}(\tilde{Y} - Y)\right) \leq h_n.$$

Therefore, if we replace the stepping scheme (25.27) with

$$\mathbf{L}(T_n) = \mathbf{L}(T_{n-1}) + M(T_{n-1}) \tau_{n-1} + \widehat{\Sigma}^{\top} Z + v^{\top}(\tilde{Y} - Y), \quad (25.33)$$

then $\mathbf{L}(T_n)$ will always, to first order, be in the survival region. In particular, to make a time step in the survival measure \tilde{Q}^B , we simply make an adjustment to the Gaussian draw to stay in the survival region and instead of $\widehat{\Sigma}^{\top} Z$ use $\widehat{\Sigma}^{\top} Z + v^{\top}(\tilde{Y} - Y)$. Moreover, from (25.32) we immediately obtain the weight that we need to apply to a Monte Carlo path with a particular draw Z (for time step $T_{n-1} \rightarrow T_n$) — it is simply equal to $\Phi(Y_{\max})$ from (25.29). Putting it all together, we obtain the following result (compare to Proposition 25.2.1).

Proposition 25.2.4. *A generalized trigger product can be valued by (25.24)–(25.25) where an Euler simulation in the survival measure is given by (25.33). The weights ψ_n in (25.25) are given by*

$$\psi_n = \prod_{k=1}^{n-1} Q_{T_{k-1}}^B (G_k(\mathbf{L}(T_k)) < h_k) = \prod_{k=1}^{n-1} \Phi(Y_{\max,k}),$$

where by $Y_{\max,k}$ we denote the value of Y_{\max} in (25.29) for time step $T_{k-1} \rightarrow T_k$.

25.3 Model-Based Control Variates

The method of control variates was first introduced in Section 3.4.3. As we recall from (3.83), the method boils down to replacing the standard Monte Carlo estimate

$$\frac{1}{K} \sum_{j=1}^K Y(\omega_j)$$

of $E(Y)$ with

$$\frac{1}{K} \sum_{j=1}^K (Y(\omega_j) - \beta^\top (Y^c(\omega_j) - E(Y^c))), \quad (25.34)$$

where $Y^c(\omega_j)$ are random samples of the potentially multi-dimensional control variate Y^c , chosen such that $E(Y^c)$ is available in closed form. As shown in (3.85), the variance reduction achieved by the method is directly proportional to the correlation between the primary variable Y and its control Y^c . There are multiple ways to select the control variate Y^c . For instance, if Y is the value of a security under a given model, then Y^c may represent the value of the same security under a different, but closely related model; or the value of a different (but related) security under the same model; or the value of an approximate hedging strategy of Y . Of course, we may also select a control variate that is a weighted combination of many different individual control variates, each chosen by a different strategy.

In the next few sections, we shall study several methods to design control variates. We start with the model-based control variate method, which uses the value of a security in a simplified proxy model as a control for the security value in the actual pricing model. To fix ideas, let \hat{V}_{orig} be the Monte Carlo estimate for the true security value in the original pricing model, and let \hat{V}_{proxy} be the Monte Carlo estimate for the same security in a proxy model. In the proxy model, assume that a highly accurate price estimate V_{PDE} is available, most likely computed by the PDE methods of Chapter 2. As in (25.34), let us introduce a corrected value estimate as

$$\hat{V}_{\text{corrected}} = \hat{V}_{\text{orig}} - \beta (\hat{V}_{\text{proxy}} - V_{\text{PDE}}),$$

where β is the appropriate regression coefficient. Assuming that $E(\hat{V}_{\text{proxy}}) = V_{\text{PDE}}$ to high precision, the new estimate will be practically unbiased. If the

path values used to compute \hat{V}_{orig} are positively correlated with the ones used to obtain \hat{V}_{proxy} , then the variance of the estimate is reduced.

The computational effort to estimate $\hat{V}_{\text{corrected}}$ is noticeably higher than that needed to compute \hat{V}_{orig} , since two additional valuations (one by Monte Carlo and one by PDE methods) are now needed. For the method to lead to an efficiency improvement³, the achieved variance reduction needs to be high, in turn requiring very high correlation between the original and proxy model path values of the security. This can typically be only achieved if the two models are closely related, and use random numbers in near-identical fashion to generate security path values. For instance, it is unlikely that one could successfully use a short rate model to compute a control variate when the original model is a Libor market (LM) model.

While on the topic of the LM model, we note that this particular model is particularly in need of variance reduction: not only is the LM model always implemented via Monte Carlo methods, it is also more computationally demanding than many other Monte Carlo based models⁴ and typically is used for complex, compute-intensive payoffs. To find a suitably faithful proxy model for a full-blown LM model, we note that while PDE methods are generally not available for LM models, PDE *approximations* are possible. While some believe that such approximations may serve as outright substitutes for LM model, in our opinion it is very difficult to make the approximations sufficiently accurate and robust to safely use them for actual security pricing. On the other hand, these approximations are often perfectly adequate for the model-based control variate method, since the requirements on proxy model precision and internal consistency are actually quite low — all that is needed is that the estimator \hat{V}_{proxy} is highly correlated to the true model price and has a limit that can be computed accurately by some other scheme. In fact, the proxy model does *not* have to be arbitrage free, which for LM models opens up the possibility of replacing complicated path-dependent drift terms with simpler ones that admit a PDE representation of security values.

25.3.1 Low-Dimensional Markov Approximation for LM models

We recall that an LM model is Markovian only in the full set of all forward Libor rates on the yield curve, plus any additional variables required to model unspanned stochastic volatility. As numerical methods for PDEs start becoming impractical when there are more than 3 or 4 state variables, a fair bit of simplification of the LM model is required to come up with a PDE-friendly model proxy. To show one way of proceeding, let us start with

³See Section 3.4.1 for a discussion of efficiency measures for variance reduction schemes.

⁴Such as the quasi-Gaussian (qG) model, which normally involves simulation of many fewer state variables than the LM model, see Section 13.1.9.3.

a one-factor model equipped with a deterministic local volatility (14.13)–(14.14), the spot measure dynamics of which we represent as

$$dL_n(t) = \varphi(L_n(t)) \lambda_n(t) (\mu_n(t, \mathbf{L}(t)) dt + dW(t)), \quad (25.35)$$

$$\mu_n(t, \mathbf{L}(t)) = \sum_{i=q(t)}^n \frac{\tau_i \varphi(L_i(t)) \lambda_i(t)}{1 + \tau_i L_i(t)},$$

for $n = 1, \dots, N - 1$, with $\mathbf{L}(t)$ being the vector of all forward Libor rates. $W(t)$ is here a one-dimensional Brownian motion under the spot measure Q^B ; an extension to two dimensions is studied in Section 25.3.2 below.

In a first step towards a low-dimensional Markovian approximation of (25.35), we look to get rid of the local volatility $\varphi(L_n(t))$. To that end, let us introduce the following transform,

$$f(x) = \int_{x_0}^x \frac{d\xi}{\varphi(\xi)}, \quad (25.36)$$

where x_0 is an arbitrary but fixed number (see also (2.81)–(2.82)). Defining new variables

$$l_n(t) \triangleq f(L_n(t)), \quad n = 1, \dots, N - 1, \quad (25.37)$$

we eliminate φ from the diffusion part of the SDE and get, for $n = 1, \dots, N - 1$,

$$dl_n(t) = \lambda_n(t) \left(\left(\mu_n(t, \mathbf{L}(t)) - \frac{1}{2} \lambda_n(t) \varphi'(L_n(t)) \right) dt + dW(t) \right). \quad (25.38)$$

For our purposes here, the main issue with (25.38) is the fact that the drift $\mu_n(t, \mathbf{L}(t))$ at each point in time depends on the whole vector of forward Libor rates. An easy way to deal with this is to simply replace in μ_n all Libor forwards with their values at $t = 0$, i.e.,

$$\mu_n(t, \mathbf{L}(t)) \approx \mu_n(t, \mathbf{L}(0)). \quad (25.39)$$

As the LM model drift terms are generally small, the usage of the first-order approximation (25.39) is certainly justifiable for control variate purposes⁵.

With approximation (25.39) we are a long way towards a low-dimensional Markov representation, but still need to simplify the term $\varphi'(L_n(t))$ in (25.38). For local volatility functions $\varphi(x)$ that are close to linear, we can use $\varphi'(L_n(t)) \approx \varphi'(L_n(0))$, an approximation that is exact for the important cases of log-normal and displaced log-normal model specifications. With this, we arrive at the following approximate SDE,

⁵Approximations to the drift could be improved by using the Brownian bridge techniques described in Section 14.6.2.5, but the impact of these improvements on the intended control variate applications is negligible.

$$dl_n(t) = \lambda_n(t) \left(\left(\mu_n(t, \mathbf{L}(0)) - \frac{1}{2} \lambda_n(t) \varphi'(L_n(0)) \right) dt + dW(t) \right). \quad (25.40)$$

In (25.40), each $l_n(t)$ is an integral of $\lambda_n(t)$ against a Brownian motion with deterministic drift. To make all the variables functions of the same state variable, we approximate the volatility structure with the following *separable* one,

$$\lambda_n(t) \approx \hat{\lambda}_n(t), \quad \hat{\lambda}_n(t) = \sigma_n \alpha(t), \quad n = 1, \dots, N-1. \quad (25.41)$$

In a separable volatility structure, each forward Libor volatility function equals a Libor-specific scalar multiplied by a function of time common to all Libor rates. This special structure allows us to define a one-dimensional Markovian state variable by

$$dX(t) = \alpha(t) dW(t), \quad (25.42)$$

and all variables $l_n(t)$ are then deterministic functions of $X(t)$:

$$l_n(t) = l_n(0) + d_n(t) + \sigma_n X(t), \quad (25.43)$$

$$d_n(t) = \int_0^t \lambda_n(s) \left(\mu_n(s, \mathbf{L}(0)) - \frac{1}{2} \lambda_n(s) \varphi'(L_n(0)) \right) ds.$$

Translated back to Libor forwards, we arrive at the reconstitution formula

$$L_n(t) = f^{-1}(f(L_n(0)) + d_n(t) + \sigma_n X(t)), \quad n = 1, \dots, N-1, \quad (25.44)$$

where we have made an implicit assumption that $f^{-1}(\cdot)$ exists (which is the case if, for example, $\varphi(\cdot)$ is positive in (25.36)).

With the representation above, at each point in time t , the value of any path-independent derivative V can be expressed as a function of t and $X(t)$,

$$V = V(t, X(t)),$$

where the function $V(t, x)$ satisfies the following PDE,

$$\frac{\partial V(t, x)}{\partial t} + \frac{\alpha(t)^2}{2} \frac{\partial^2 V(t, x)}{\partial x^2} = r(t, x) V(t, x), \quad (25.45)$$

subject to appropriate boundary and jump conditions. In (25.45), $r(t, X(t))$ is the discounting rate applied to any payoff over an instantaneous period of time $[t, t+dt]$, whose specific expression in terms of Libor rates (and, ultimately, $X(t)$) depends on the interpolation method used. For instance, under the assumption that instantaneous forward rates $f(t, u)$ are constant⁶ for $u \in [t, T_{q(t)}]$ (see Section 15.1.6 and equation (15.20)), we obtain

$$r(t, X(t)) = \frac{1}{\tau_{q(t)}} \ln(1 + \tau_{q(t)} L_{q(t)}(t)),$$

where the right-hand side is understood to be a function of $X(t)$ via (25.44).

⁶We generally do not recommend this interpolation scheme, but it makes for a good example.

25.3.2 Two-Dimensional Extension

Before turning to the question of how to pick the volatility term structure for the LM proxy model, let us consider the extension to LM models driven by a two-dimensional Brownian motion. To build a two-factor Markovian proxy model, assume that

$$W(t) = (W^1(t), W^2(t)),$$

and that forward Libor volatilities are two-dimensional processes,

$$\lambda_n(t) = (\lambda_n^1(t), \lambda_n^2(t)).$$

In the spot measure, the deterministic local volatility LM model then follows

$$\begin{aligned} dL_n(t) &= \sum_{k=1}^2 \lambda_n^k(t) \varphi(L_n(t)) (\mu_n^k(t, \mathbf{L}(t)) dt + dW^k(t)), \\ \mu_n^k(t, \mathbf{L}(t)) &= \sum_{i=q(t)}^n \frac{\tau_i \varphi(L_i(t)) \lambda_i^k(t)}{1 + \tau_i L_i(t)}, \end{aligned}$$

for $n = 1, \dots, N - 1$. Following the steps in Section 25.3.1, we eventually come to the point where we need to approximate the volatility structure with a separable one similar to (25.41). A naive generalization would specify

$$\begin{aligned} \lambda_n^1(t) &\approx \widehat{\lambda}_n^1(t), \quad \widehat{\lambda}_n^1(t) = \sigma_n^1 \alpha^1(t), \\ \lambda_n^2(t) &\approx \widehat{\lambda}_n^2(t), \quad \widehat{\lambda}_n^2(t) = \sigma_n^2 \alpha^2(t), \end{aligned} \tag{25.46}$$

for $n = 1, \dots, N - 1$. However, an extension is possible (and desirable, as will be clear later), as we can use the more general expression

$$\begin{aligned} \widehat{\lambda}_n^1(t) &= \sigma_n^1 \alpha^{11}(t), \\ \widehat{\lambda}_n^2(t) &= \sigma_n^1 \alpha^{21}(t) + \sigma_n^2 \alpha^{22}(t), \end{aligned} \tag{25.47}$$

while keeping the approximation Markovian. In particular we can then define two (correlated) state variables by⁷

$$\begin{aligned} dX_1(t) &= \alpha^{11}(t) dW^1(t) + \alpha^{21}(t) dW^2(t), \\ dX_2(t) &= \alpha^{22}(t) dW^2(t). \end{aligned}$$

The Libor rates can then be computed by (compare to (25.44))

⁷We can use the triangular form here because the square root of a variance-covariance matrix can always be written this way by application of the Cholesky decomposition, see Section 3.1.2.1. A more general form, however, could be beneficial for fitting as discussed later, see footnote 9.

$$L_n(t) = f^{-1} \left(f(L_n(0)) + \tilde{d}_n(t) + \sigma_n^1 X_1(t) + \sigma_n^2 X_2(t) \right), \quad (25.48)$$

for $n = 1, \dots, N - 1$, where the deterministic part $\tilde{d}_n(t)$ is suitably defined. The (two-dimensional) valuation PDE for $V = V(t, x, y)$ is now given by

$$\begin{aligned} \frac{\partial V}{\partial t} + \frac{\alpha^{11}(t)^2 + \alpha^{21}(t)^2}{2} \frac{\partial^2 V}{\partial x^2} + \alpha^{21}(t)\alpha^{22}(t) \frac{\partial^2 V}{\partial x \partial y} \\ + \frac{\alpha^{22}(t)^2}{2} \frac{\partial^2 V}{\partial y^2} = r(t, x, y)V, \end{aligned} \quad (25.49)$$

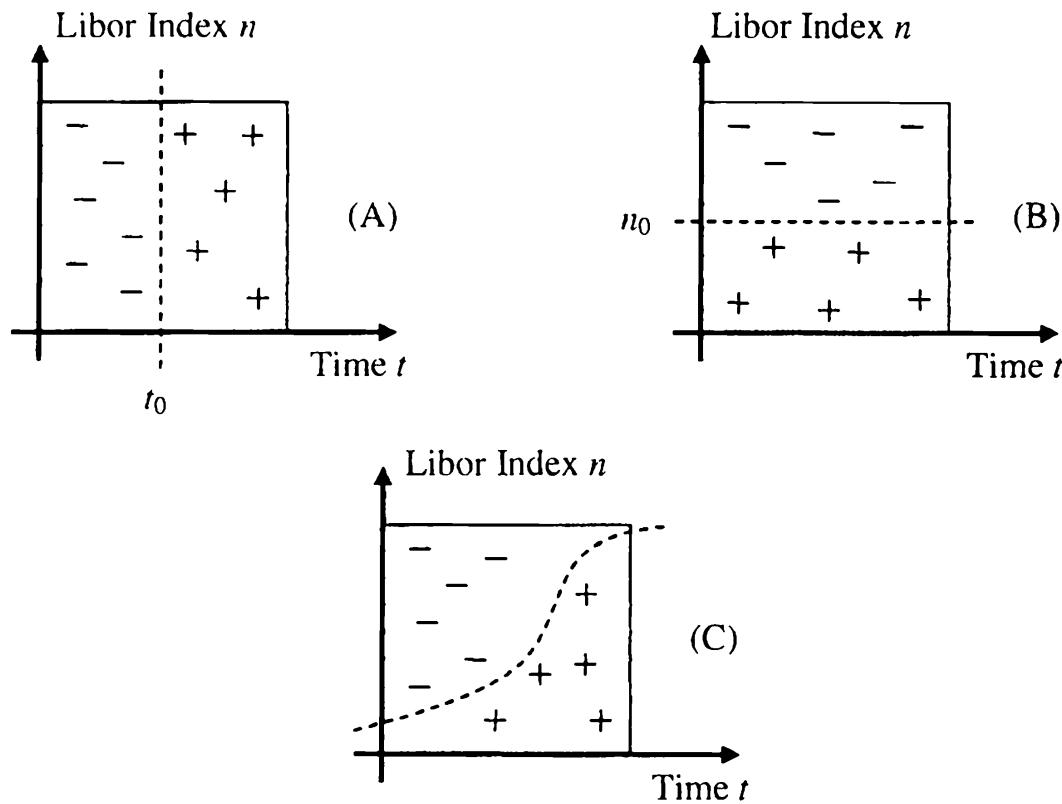
and can be solved with, for example, the Craig-Sneyd scheme from Section 2.11.2. The terminal condition $V(T, x, y)$ here is determined from the payoff of the derivative, after expressing the yield curve at time T in terms of the state variables $X_1(T), X_2(T)$ through the reconstitution formula (25.48).

The specification (25.47) is normally substantially more accurate than (25.46) in approximating the volatility structure of the original model. To understand why, recall from the principal components analysis in Section 14.3.1 that the first volatility component of the original model $\lambda_n^1(t)$ normally represents a near-parallel yield curve shift and is positive for all t and n — a shape that can be represented in the form (25.46) quite well. However, the second component $\lambda_n^2(t)$ models a yield curve twist (see Figure 14.1) which, for a fixed value of t , will require that $\{\lambda_n^2(t)\}$ crosses zero for some value of n . With this in mind, consider the approximations in (25.46) and (25.47). In the former, a function of the form $\sigma_n^2 \alpha^2(t)$ can cross zero either “vertically” ($\sigma_n^2 \alpha^2(t) = 0$ for some $t = t_0$ for all n) or “horizontally” ($\sigma_n^2 \alpha^2(t) = 0$ for some $n = n_0$ for all t). In reality, due to an imposed (or desired) time-homogeneity, $\lambda_n^2(t)$ usually crosses zero “diagonally”, in the loose sense that the function $n_0(t) \triangleq \{n : \lambda_n^2(t) = 0\}$ grows with t . Figure 25.1 demonstrates the point. On each of the three diagrams of the figure, we plot signs of the second volatility component for all points $(t, n) \in [0, T_{N-1}] \times \{1, \dots, N-1\}$. The plus symbol indicates a positive value and the minus symbol represents a negative one. Diagrams (A) and (B) represent the only two possibilities for the second volatility component of the form (25.46). The diagram (C) shows how a typical second volatility component really looks like, a behavior that can be replicated by (25.47) but *not* by (25.46).

25.3.3 Approximating Volatility Structure

So far we have glossed over the actual mechanics of approximating the original model volatility $\lambda(t)$ with a separable proxy version, as in (25.41) or (25.47). One approach would be to perform an outright calibration (see Section 14.5) of the Markov proxy model to the same swaption quotes used to calibrate the original model. We generally do not recommend this when we use the Markov LM model to generate a control variate. Instead, recall

Fig. 25.1. Sign of the Second Volatility Component



Notes: Signs of the second volatility component for separable (diagrams (A), (B)) parameterization of the form (25.46), and a typical non-separable one (diagram (C)).

from Section 3.4.3 that the variance reduction achieved by a control variate method is strongly correlation dependent, which suggests that we attempt to approximate the *factor* volatilities of the original model directly, without much consideration given to the precision with which the proxy model can price swaptions. In this spirit, calibration of, say, the two-factor separable volatility can be stated as a least-squares problem⁸

$$\begin{aligned} & \sum_{i=0}^{N-2} \sum_{n=1}^{N-1} (\lambda_n^1(T_i) - \sigma_n^1 \alpha^{11}(T_i))^2 \\ & + \sum_{i=0}^{N-2} \sum_{n=1}^{N-1} (\lambda_n^2(T_i) - (\sigma_n^1 \alpha^{21}(T_i) + \sigma_n^2 \alpha^{22}(T_i)))^2 \rightarrow \min. \quad (25.50) \end{aligned}$$

Here we optimize over all σ 's and α 's, for a total of $5 \times (N - 1)$ variables. Of course, we may extend the norm as we see fit, to include smoothing penalty terms or to use different weights on different terms or factors.

⁸Note that, in line with (14.42), we assume that λ 's and α 's are piecewise constant over $[T_i, T_{i+1})$, $i = 0, \dots, N - 2$.

Denson and Joshi [2009] propose a slightly different fitting algorithm. First, we would set

$$\sigma_n^1 = \lambda_n^1(T_0), \quad \sigma_n^2 = \lambda_n^2(T_0), \quad n = 1, \dots, N - 1$$

and

$$\alpha^{11}(T_0) = \alpha^{22}(T_0) = 1, \quad \alpha^{21}(T_0) = 0.$$

This ensures that $\lambda_n^i(t) = \widehat{\lambda}_n^i(t)$ for $t \in [T_0, T_1]$ for all $n = 1, \dots, N - 1$, and $i = 1, 2$. Then we would solve the minimization problem (notice index i starting at $i = 1$, rather than at $i = 0$)

$$\begin{aligned} & \sum_{i=1}^{N-2} \sum_{n=1}^{N-1} (\lambda_n^1(T_i) - \lambda_n^1(T_0)\alpha^{11}(T_i))^2 \\ & + \sum_{i=1}^{N-2} \sum_{n=1}^{N-1} (\lambda_n^2(T_i) - (\lambda_n^1(T_0)\alpha^{21}(T_i) + \lambda_n^2(T_0)\alpha^{22}(T_i)))^2 \rightarrow \min \end{aligned} \quad (25.51)$$

for α 's only, a quadratic optimization problem that is solved analytically⁹, e.g.

$$\alpha^{11}(T_i) = \frac{\sum_{n=1}^{N-1} \lambda_n^1(T_i)\lambda_n^1(T_0)}{\sum_{n=1}^{N-1} (\lambda_n^1(T_0))^2}, \quad i = 1, \dots, N - 2,$$

and so on. The direct and analytic linkage of $\widehat{\lambda}$'s to λ 's in this approach could lead to better performance of the control variate method when calculating risk sensitivities, especially vegas.

25.3.4 Markov Approximation as a Control Variate

To use the Markov approximation as the control variate, we define $\widehat{V}_{\text{proxy}}$ (see the beginning of Section 25.3 for notations) to be the value of the security in the Markov LM model computed by Monte Carlo, and by V_{PDE} the value computed by the PDE method in the same model. For V_{PDE} to be consistent with $\widehat{V}_{\text{proxy}}$, both the Monte Carlo and the PDE methods should be applied to the same derivative. While this may seem like a trivial point, some cases are fairly subtle and require care. For example, for callable Libor exotics, the $\widehat{V}_{\text{proxy}}$ value would often be a lower bound (Section 18.3) on the actual value of the callable derivative and, effectively, would represent the value of an exotic swap that knocks out at the estimated exercise boundary. In this case the PDE method should be applied not to a callable Libor exotic, but to a knockout swap with a knockout boundary lifted directly from the Monte

⁹Denson and Joshi [2009] also extend the specification (25.47) by allowing an extra term for $\widehat{\lambda}_n^1(t)$ with $\widehat{\lambda}_n^1(t) = \sigma_n^1\alpha^{11}(t) + \sigma_n^2\alpha^{12}(t)$, which may improve the fit somewhat.

Carlo valuation. For this to work in practice, the LS regression method for the estimation of the exercise boundary in the Markov LM model should use the Markov state variables $X(t)$ as explanatory variables¹⁰ so that relevant regression functions can be easily transferred into the PDE setup.

Achieving a high correlation between the path values of a derivative in the original and proxy models is crucial to the performance of the model-based control variate method, so care must be taken in ensuring that the simulation schemes for the original and Markov models are as similar as possible. This ranges from the obvious requirement of using the same simulation seed for random number generation in the two models, to the more subtle issue of discretization scheme compatibility. For example, suppose we use the Euler discretization scheme on the original LM model SDE (25.35). Then, for the Markov LM model we must also use the Euler scheme on the SDE

$$\begin{aligned} dL_n(t) &= \widehat{\lambda}_n(t)\varphi(L_n(t)) \\ &\times \left(\mu_n(t, \mathbf{L}(0)) - \frac{1}{2}\widehat{\lambda}_n(t)(\varphi'(L_n(0)) - \varphi'(L_n(t))) dt + dW(t), \right) \end{aligned}$$

obtained from (25.40) and (25.37). Here $\widehat{\lambda}_n(t) = \sigma_n\alpha(t)$ of course, or the equivalent for the two-dimensional model. In particular, notice that to keep the simulation of the two models in lock-step, we must resist the temptation to be clever, and avoid using special-purpose discretization schemes that take advantage of the simple form of the Markov proxy model dynamics.

The performance of the model-based control variate method can often be quite impressive, with Piterbarg [2003] reporting a reduction in sample standard deviation by a factor of 3 to 10, corresponding to a speed improvement of 10 to 100 times (see Section 3.4.1). Of course, there is extra work involved that includes an extra Monte Carlo simulation for the Markov model, and a (relatively speedy) PDE valuation. The potential downside of the method is the fact that its scope is somewhat limited by the need to perform a PDE valuation. With three dimensions probably being the practical maximum for a reasonably quick PDE scheme, the model-based control variate method is limited to either i) a two-factor Libor market model (as we developed above), ii) a three-factor model (a straightforward extension), iii) a two-factor model with stochastic volatility, or iv) a two-factor model for a path-dependent trade that could be treated in PDE by introducing an extra state variable (e.g., a TARN, see Section 20.1.5). For products and models that do not fit these categories, a proxy model may still be useful for defining dynamic control variates, as demonstrated in Section 25.5. In addition, we always have the option of using instrument-based control variates, as we describe next.

¹⁰This is a good idea for any Markovian model, see Section 18.3.9.1.

25.4 Instrument-Based Control Variates

In the model-based control variate method, we created a control variate by introducing a new model and applying it to the (unchanged) payoff of the security we look to price. In the instrument-based control variate method, in a sense we do the opposite: we keep the model fixed but change the payoff. In fixed income applications, the idea of using proxy securities as control variates is most closely associated with Bermudan swaption pricing, but the basic ideas often extend fairly naturally to more complicated callable Libor exotics.

For concreteness, let us start by considering a Bermudan swaption with $N - 1$ exercise opportunities. We follow the notation of Chapter 19 and, in particular, denote the $N - 1$ exercise values by $U_n(t)$, $n = 1, \dots, N - 1$ (see (19.1)). The K -path Monte Carlo estimate value of a Bermudan swaption, or indeed a general CLE, is given by

$$H_0(0) \approx \frac{1}{K} \sum_{k=1}^K Y(\omega_k), \quad Y(\omega) = \sum_{i=\eta(\omega)}^{N-1} B(T_{i+1}, \omega)^{-1} X_i(\omega), \quad (25.52)$$

where $\{\omega_k\}_{k=1}^K$ are the simulated paths and η the (estimate of the) optimal exercise time index.

Naively, we could try to introduce control variates based on the $N - 1$ (deflated) exercise values, as observed at the final expiry time T_N . That is, we could define controls $Y^c = (Y_1^c, \dots, Y_{N-1}^c)^\top$, where

$$Y_n^c(\omega) \triangleq \sum_{i=n}^{N-1} B(T_{i+1}, \omega)^{-1} X_i(\omega), \quad n = 1, \dots, N - 1. \quad (25.53)$$

Each control is then a sum of path values of net coupons. Alternatively, for Bermudan swaptions in particular, we can use

$$Y_n^c(\omega) = U_n(T_n, \omega), \quad n = 1, \dots, N - 1, \quad (25.54)$$

where, of course, each $U_n(t)$ is the value at time t of all net coupons from the n -th one onward (see (18.2)). For Bermudan swaptions the $U_n(t)$ are just swap values and are available bias-free in a closed form expression. Note that (25.54) constitutes a *different* set of control variates than (25.53).

Both of the control variate schemes outlined above are quite simplistic and typically fail to yield good variance reduction. One reason is that both these control variates are, in the case of a standard Bermudan swaption, essentially linear functions of rates, whereas the payoff of a Bermudan swaption is option-like and clearly not well-approximated by a linear function. We can attempt to rectify this issue by using control variates that are non-linear in rates. European swaptions are natural choices for this, but often must be ruled out due to lack of exact valuation formulas¹¹ or approximations that are

¹¹The swap market models of Section 15.4 are, however, not affected by this issue.

sufficiently accurate across a wide range of moneyness and expiries. For Libor market models, however, caps (or floors) often have exact pricing formulas, so these instruments may be a good option; we explore this idea in more detail below. First, however, let us note that a more subtle reason for the failure of the schemes (25.53) and (25.54) is due to the fact that the control variates effectively are observed “at the wrong time”. To elaborate, notice that the value of a Bermudan swaption (or a CLE) along a path ω in (25.52) involves cash flows fixed at times $T_{\eta(\omega)}, \dots, T_N$, whereas all control variates in (25.53) always include a deterministic number of cash flows. Similarly, in (25.54) each control variate is sampled at a single time only. So, compared to controls, a path value of a Bermudan swaption will often have an incorrect number of net coupons included, and will likely have low correlation with the controls.

The fact that the timing mismatch contributes significantly to de-correlation between a Bermudan swaption and naive controls was noted by Rasmussen [2005], who also proposed to rectify the issue by sampling the controls *at the Bermudan swaption exercise time*. Here is the technical result that justifies this choice.

Proposition 25.4.1. *With T_n denoting the n -th exercise date, set $U_n = U_n(T_n)$, and let Z_n , $n = 0, \dots, N$, be a martingale process with respect to $\{\mathcal{F}_n \triangleq \mathcal{F}_{T_n}\}_{n=0}^N$. Let stopping times $\eta, \sigma \in \{1, \dots, N-1\}$ be given such that $\eta \leq \sigma$. Then*

$$(\text{Corr}(U_\eta, Z_\eta))^2 \geq (\text{Corr}(U_\eta, Z_\sigma))^2.$$

Proof. The proof follows by the repeated applications of the optional sampling theorem. For the covariance term, we have

$$\begin{aligned} \text{Cov}(U_\eta, Z_\sigma) &= E(U_\eta Z_\sigma) - E(U_\eta) E(Z_\sigma) \\ &= E(U_\eta E(Z_\sigma | \mathcal{F}_\eta)) - E(U_\eta) E(E(Z_\sigma | \mathcal{F}_\eta)) \\ &= E(U_\eta Z_\eta) - E(U_\eta) E(Z_\eta) \\ &= \text{Cov}(U_\eta, Z_\eta). \end{aligned}$$

For the variance term,

$$\begin{aligned} \text{Var}(Z_\sigma) &= E(Z_\sigma^2) - (E(Z_\sigma))^2 \\ &= E(E(Z_\sigma^2 - Z_\eta^2 + Z_\eta^2 | \mathcal{F}_\eta)) - E(E(Z_\sigma | \mathcal{F}_\eta))^2 \\ &= E(E(Z_\sigma^2 - Z_\eta^2 | \mathcal{F}_\eta)) + E(Z_\eta^2 - E(Z_\eta)^2) \\ &= E(\text{Var}(Z_\sigma | \mathcal{F}_\eta)) + \text{Var}(Z_\eta) \\ &\geq \text{Var}(Z_\eta), \end{aligned}$$

and the result follows. \square

To understand the implications of Proposition 25.4.1, consider using a European option maturing at time T_n as a control variate. If for a particular

Monte Carlo path we have that $\eta < n$, then the proposition essentially suggests that we should use the value of the option at time T_η to generate a control variate, rather than wait until the maturity date T_n . Of course, since the result in the proposition deals with martingale controls, a little care is required in creation of the control variates. Specifically, most interest rate derivatives (including the prospective controls in (25.53)) pay coupons, and hence need to be adjusted to become martingales. Fortunately, this is fairly easy to do: all coupons paid to time t *should not* be dropped from the value of the control at time t , but should be rolled up (using the numeraire) to time t . In particular, since our definition of $U_n(t)$'s in (18.2) makes sense even for $t > T_n$, we define new controls $Y^c = (Y_1^c, \dots, Y_{N-1}^c)^\top$ by

$$\begin{aligned} Y_n^c \triangleq U_n(T_\eta) &= B(T_\eta) \sum_{i=n}^{\eta-1} B(T_{i+1})^{-1} X_i \\ &\quad + B(T_\eta) E_{T_\eta} \left(\sum_{i=\max(\eta, n)}^{N-1} B(T_{i+1})^{-1} X_i \right), \end{aligned} \quad (25.55)$$

for $n = 1, \dots, N - 1$ (where we use the convention that $\sum_a^b = 0$ if $b < a$). While these controls still do not exhibit the non-linear characteristic of options (we turn to this shortly), they do resolve the timing problem. By the optional sampling theorem the exact values of the controls, as required by the control variate method, are known,

$$E(U_n(T_\eta)) = U_n(0), \quad n = 1, \dots, N - 1.$$

In theory, adding more controls never increases the variance of the estimate. However, as explained in Section 3.4.1, the efficiency of the scheme can very well decrease as the additional computational effort of new controls may not be rewarded by sufficiently high decreases in variance. As a consequence, one often is best served by selecting just a few carefully chosen controls, rather than indiscriminately throwing a large set of suboptimal controls at the problem. For the case above, the set of exercise values U_n , $n = 1, \dots, N - 1$, is composed of various subsets of (net) coupons already contained in the “longest” underlying U_1 . The efficiency gains from using sums of subsets of coupons as a vector control compared to using just the sum of all coupons as a single control can rightly be questioned, suggesting that the following one-dimensional control may be useful:

$$Y^c = Y_1^c, \quad (25.56)$$

$$Y_1^c \triangleq U_1(T_\eta) = B(T_\eta) \sum_{i=1}^{\eta-1} B(T_{i+1})^{-1} X_i + B(T_\eta) E_{T_\eta} \left(\sum_{i=\eta}^{N-1} B(T_{i+1})^{-1} X_i \right).$$

To introduce non-linearity into the controls for a Bermudan swaption, we can consider using caplets and caps, as these can be valued exactly in the majority of LM models. For concreteness, let us focus on a payer Bermudan swaption with strike k , the exercise values of which are given by

$$U_n(t) = B(t) \sum_{i=n}^{N-1} E_t \left(B(T_{i+1})^{-1} (L_i(T_i) - k) \tau_i \right),$$

where $L_i(T_i)$ is a Libor rate observed at time T_i for the period $[T_i, T_{i+1}]$. To construct a suitably non-linear control for the Bermudan swaption, consider using the set of caps¹²,

$$V_{\text{cap},n}(t) \triangleq B(t) \sum_{i=n}^{N-1} E_t \left(B(T_{i+1})^{-1} (L_i(T_i) - k)^+ \tau_i \right), \quad n = 1, \dots, N-1.$$

With these caps we can construct a few possible controls. For example, in direct analogy to (25.55), we can use a collection of all caps (observed at the Bermudan exercise time) as an $(N-1)$ -dimensional control,

$$Y^c = (Y_1^c, \dots, Y_{N-1}^c)^\top, \quad (25.57)$$

$$\begin{aligned} Y_n^c &\triangleq V_{\text{cap},n}(T_\eta) = B(T_\eta) \sum_{i=n}^{\eta-1} B(T_{i+1})^{-1} (L_i(T_i) - k)^+ \tau_i \\ &+ B(T_\eta) E_{T_\eta} \left(\sum_{i=\max(\eta,n)}^{N-1} B(T_{i+1})^{-1} (L_i(T_i) - k)^+ \tau_i \right), \end{aligned}$$

for $n = 1, \dots, N-1$. Alternatively, since each cap is just a sum of different caplets $(L_i(T_i) - k)^+$, a closely related, but simpler, control can be constructed by using all caplets (instead of all caps), again sampled at the exercise time,

$$\begin{aligned} Y^c &= (Y_1^c, \dots, Y_{N-1}^c)^\top, \\ Y_n^c &\triangleq B(T_\eta) E_{T_\eta} \left(B(T_{n+1})^{-1} (L_n(T_n) - k)^+ \tau_n \right), \end{aligned}$$

for $n = 1, \dots, N-1$. Furthermore, in direct analogy to (25.56), we can use only one control, the longest cap, stopped at the exercise time,

¹²While we use the fixed rate of the swap as a strike for a cap, a potentially better control variate could be constructed by using a strike at or near the exercise boundary of the Bermudan swaption.

$$Y^c = Y_1^c, \quad (25.58)$$

$$\begin{aligned} Y_1^c &\triangleq V_{\text{cap},1}(T_\eta) \equiv B(T_\eta) \sum_{i=1}^{\eta-1} B(T_{i+1})^{-1} (L_i(T_i) - k)^+ \tau_i \\ &+ B(T_\eta) \sum_{i=\eta}^{N-1} E_{T_i} \left(B(T_{i+1})^{-1} (L_i(T_i) - k)^+ \tau_i \right). \end{aligned}$$

We note that using just the longest cap as a control variate can be interpreted as using all caplets as controls, but enforcing the same regression coefficient β for all of them. Finally, linear (e.g. (25.55)) and non-linear (e.g. (25.57)) controls can be combined together to improve variance reduction further over a wide range of strikes and maturities.

Various strategies for constructing Bermudan swaption control are tested in a LM model setup by Jensen and Svenstrup [2003]; their main conclusion is that the combination of caplets (25.57) and linear controls (25.55) performs well for a diverse set of Bermudan swaptions. Using the longest cap in combination with the longest swap resulted in only a slightly worse control. The achieved reduction in sample standard deviation is of the order 3 to 5 (i.e. speed ups of up to order 9 to 25 times).

For securities more complex than standard Bermudan swaptions, such as general callable Libor exotics, the idea of sampling the controls at the exercise time still applies. The choice of controls, however, becomes non-trivial. If the underlying can be valued in closed form, such as for callable inverse floaters or callable capped floaters on Libor rates, it should be used. The extra non-linearity in the payoff can be handled with caps of different strikes. However, when pricing CLEs we will often find that the underlying exotic swap does not permit closed-form valuation, preventing us from using the underlying as a control. For these securities, a more general, dynamic type of control variate may be an alternative, as we describe next.

25.5 Dynamic Control Variates

As demonstrated in the previous section, even for relatively simple securities such as Bermudan swaptions, finding good control variates is often challenging. For more complicated CLEs, the search becomes increasingly involved and, what is probably worse, usually has to be done on a case-by-case basis: what works for callable range accruals will probably not work for callable CMS spread options.

In contrast, *dynamic*, or delta-based, control variates are always available — at least in theory. As discussed in Section 3.4.3.2, the main idea behind this method is to select as a control variate the value of a self-financing hedging strategy for the security to be priced. Constructing the exact hedging strategy requires knowledge of the deltas of the security, at each point in

time and for each realization of the Monte Carlo paths. These are, of course, rarely available, but often we can use deltas constructed using *approximate* risk sensitivities instead.

Approximate deltas can be constructed in a number of ways. One idea, suggested by Clewlow and Carverhill [1994] and mentioned already in Section 3.4.3.2, uses deltas from a tractable proxy model. For the Libor market model, these deltas could originate from, say, an approximate Markov proxy model, as described in Section 25.3. Without resorting to proxy models, a general technique suitable for CLEs is suggested in Moni [2005], who proposes to extract approximate deltas from regressed values of CLE prices, as computed by the LS method (see Section 18.3). The method capitalizes on the fact that the regressed values of the CLE are designed to be good approximations to the actual CLE values under various market scenarios. Let us describe this method in a bit more detail.

We use the basic scheme of Section 18.3.1 as an example, with regression variables defined to be polynomials of explanatory variables $x(t) = (x_1(t), \dots, x_d(t))$ as described in Section 18.3.9.2. The regression approximation (18.11) to the hold value can be written as (using index n instead of $n - 1$ to simplify notations)

$$\tilde{H}_n(T_n) = p_n(x(T_n)), \quad n = 0, \dots, N - 1,$$

where $p_n(x)$'s are polynomials in d variables, obtained “for free” as part of the LS algorithm. With this representation, we can compute approximate sensitivities with respect to explanatory variables

$$\frac{\Delta H_n(T_n)}{\Delta x_m(T_n)} \approx \frac{\Delta \tilde{H}_n(T_n)}{\Delta x_m(T_n)} = \frac{\partial p_n}{\partial x_m}(x(T_n)), \quad m = 1, \dots, d,$$

to which corresponds the approximate hedging strategy

$$V_{hs}(T_n) = H_0(0) + \sum_{j=0}^{n-1} \left(\sum_{m=0}^d \frac{\partial p_n}{\partial x_m}(x(T_j)) (x_m(T_{j+1}) - x_m(T_j)) \right)$$

for $n = 1, \dots, N - 1$. The expected value of the hedging strategy at time 0 under some pricing measure (such as the often-used spot measure) is given by

$$E(V_{hs}(T_n)) = \sum_{j=0}^{n-1} E \left(\sum_{m=0}^d \frac{\partial p_n}{\partial x_m}(x(T_j)) (E(x_m(T_{j+1}) | \mathcal{F}_{T_j}) - x_m(T_j)) \right), \quad (25.59)$$

which needs to be known bias-free for the control variate method to work. The easiest way to calculate (25.59) is to select all explanatory variables to be martingales, so that

$$E(x(T_{j+1}) | \mathcal{F}_{T_j}) = x(T_j) \quad (25.60)$$

for each j , and

$$\mathbb{E}(V_{\text{hs}}(T_n)) = H_0(0).$$

Moreover, in this case the hedging strategy is itself a martingale,

$$\mathbb{E}(V_{\text{hs}}(T_n) | \mathcal{F}_{T_k}) = V_{\text{hs}}(T_k)$$

for $k \leq n$.

The martingale requirement (25.60) on the explanatory variables is a restriction on the set of all possible explanatory variables used in the LS method, but not a very severe one. Recall (Section 18.3.9.2) that we typically advocate using financially meaningful quantities as explanatory variables. Numeraire-deflated values of traded securities are both financially-meaningful *and* are martingales, hence they can be used for construct a dynamic control variate. For explanatory variables that themselves cannot be represented as prices of traded securities prices, slight modifications in variable selection can often be used to make them so. For example, while a swap rate is not a martingale, the closely related (deflated) swap *value* is.

In some cases, linking a required explanatory variable to a particular security price may be difficult, such as for the stochastic variance factor in a stochastic volatility Libor market model (see Section 18.3.9.3). In such cases, we can always construct a martingale from the explanatory variable by simply subtracting out its mean over each simulation time step, as already done in Section 3.5.5 for martingale construction (see e.g. (3.118)).

Once we have constructed an approximate hedging strategy, we can define a control variate as the hedging strategy stopped at the exercise time (employing a key insight from Section 25.4),

$$Y^c = Y_1^c, \quad Y_1^c \triangleq V_{\text{hs}}(T_\eta).$$

Tests in Moni [2005] show that this approach typically yield reductions in the standard error by a factor of two to three.

The quality of the hedging strategy produced by regressions will depend on the quality of the estimated future hold values implied by the regression functions. While the bias in the lower bound value of the CLE itself also depends on the quality of the regression approximation, there is an important difference: for the basic lower bound price to be of good quality, the approximations to the hold (and exercise) values need only be accurate around the exercise boundary. When using the regression to construct a dynamic control variate, however, the approximations need to be accurate over the whole range of possible values of explanatory variables. The former is obviously much easier to achieve than the latter. When using regression to produce a control variate, Moni [2005] recommends to use polynomials of a higher degree than for the basic lower-bound valuation routine (which, incidentally, means that the regression results no longer come for free from the basic valuation, although the extra cost is modest).

Dynamic hedging ideas are studied in Jensen and Svenstrup [2003], in the context of the approach described in Section 25.4. Here, the authors use hedging strategies to represent values of core European swaptions in a Bermudan swaption, as European option values (and not caps) appear to be the more natural option-based controls. As closed-form values of these controls in Libor market models are unavailable, they resort to approximate hedging strategies based on deltas generated by the swaption approximation formulas that are available for LM models. The authors report good results for this method as well.

Finally, let us note that a number of additional twists on the theme of dynamic control variates have emerged in the literature, some of which are based on information extracted from the upper bound methods of Sections 3.5.5 and 18.3.8. Representative papers on the application of these techniques to CLEs include Beveridge and Joshi [2009] and Bender et al. [2006]. In Juneja and Kalra [2009], the authors additionally suggest to use a measure change arising from multiplicative duality¹³ for importance sampling.

25.6 Control Variates and Risk Stability

We finish this chapter with a caveat. The various control variate methods discussed in this chapter often show impressive reductions in simulation error on the basic security price, but are not always equally effective in reducing simulation errors of *risk sensitivities*. This observation stems from the fact that the sources of simulation errors when calculating risk sensitivities often bear little relationship to the sources of simulation error of the value itself. The following simple example should clarify this point. Consider the problem of pricing a digital option on the underlying X with the payoff

$$1_{\{X>b\}}.$$

As X is positively correlated with $1_{\{X>b\}}$, we can use X itself as a control variate to reduce the variance. Then the value of the security using such control is effectively equal to the value of a new security with the payoff

$$1_{\{X>b\}} - \beta(X - E(X)),$$

with β the regression coefficient. Clearly, however, the risk sensitivities of this new security would exhibit the same level of simulation error as the original one, as both payoffs have the same jump discontinuity at $X = b$ which, as

¹³Multiplicative duality is developed in Jamshidian [1995] and is closely related to the (additive) duality of Section 1.10.2. A comparison of multiplicative and additive duality for upper bound simulations can be found in Chen and Glasserman [2007a].

discussed in Chapter 23, is the dominant factor affecting the simulation error and the stability of risk sensitivities here. This problem is fairly typical of variance reduction techniques in general, and control variate methods in particular. For irregular payoffs, we find that we typically get more “bang for the buck” out of techniques that focus specifically on improving risk stability, rather than on general variance reduction. Sample techniques include the smoothing methods of Chapter 23, the non-perturbation methods of Chapter 24 or, perhaps, payoff smoothing through importance sampling (Section 25.2). Of course, nothing prevents one from combining general variance reduction techniques with payoff smoothing methods, a strategy that often works very well.

Vegas in Libor Market Models

We recall that vega measures the sensitivity of a security price to moves in volatility. In interest rate models with rich volatility structures, calculating (and even defining, for that matter) vega can be surprisingly difficult, especially in a Monte Carlo setting where vega computations bring about a new layer of complexity beyond the standard challenges discussed in recent chapters. Since, as explained in Chapter 22, vega is of fundamental importance in risk management, the ability to robustly and accurately compute vega is a key requirement for any actual model implementation. This final chapter of the book is dedicated to the challenging topic of vega computations, mostly using the Libor market (LM) model as a convenient, and highly relevant, example.

26.1 Basic Problem of Vega Computations

As discussed in Section 4.4, any diffusive (HJM) model of interest rates is uniquely defined by its volatility structure $\sigma_f(t, T)$ (see Lemma 4.4.1). At the most fundamental level, vega calculations involve the computation of interest rate derivative price sensitivities to changes to this fundamental volatility structure. For a general model, $\sigma_f(t, T)$ is two-dimensional¹, depending on both calendar time t and time to maturity T . For a given interest rate derivative we, in principle, are faced with the problem of quantifying the impact on the derivative security value of *all* possible two-dimensional shocks to this volatility structure. While the space of all possible shocks to a two-dimensional surface is quite rich, in theory we can decompose each shock into a linear combination of “Dirac delta” shocks to individual points (t, T) , $0 \leq t \leq T < T_{\max}$, and measure vegas to those shocks only. This is sufficient,

¹In multi-factor models, for each t and T , $\sigma_f(t, T)$ is obviously a multi-dimensional vector, but it is not this dimensionality that we are interested in here.

as vega is a first-order sensitivity and must therefore be linear with respect to linear combination of shocks.

While interest rate vegas are fundamentally two-dimensional, simpler types of interest rate models often reduce this dimensionality for tractability. For instance, one-factor Gaussian and quasi-Gaussian models reduce the two-dimensional structure of a generic HJM volatility structure to a separable factor form

$$\sigma_f(t, T) = g(t)h(T) \quad (26.1)$$

for some $g(t), h(T)$ (see (4.44) and (13.2)). A shock to the volatility structure that preserves the form (26.1) obviously cannot be two-dimensional, and a two-dimensional “Dirac delta” shock will not preserve the factor form. Instead, if we wish to measure volatility sensitivities in models satisfying (26.1), we would have to either bump the volatility structure for a fixed t and all T (a shock to function $g(\cdot)$) or to bump it for all t but a fixed T (a shock to function $h(\cdot)$). In other words, the set of volatility shocks that preserve the volatility structure factor form (26.1) is significantly reduced relative to the general case, and we are consequently prevented from measuring the impact of many types of potentially relevant volatility shocks. Of course, if the factor decomposition is refined relative to (26.1) by using additional state variables (see Section 12.1.5, for instance, or our discussion of multi-factor quasi-Gaussian models in Section 13.3.2), then more complicated shock shapes may be approximated to arbitrarily high precision.

While the discussion above concern model vegas, i.e. sensitivities with respect to perturbations of the model volatility structure, it is often the market vegas, i.e. vegas with respect to volatilities of market-observed vanilla options (see Section 22.1.4), that are ultimately of most practical interest. The dimensionality issue touched upon earlier is equally present in market vegas, since the set of European swaptions is two-dimensional, indexed by option expiry and swap tenor². The dimensionality reduction that is implicit in models such as a one-factor Gaussian model means that we are sometimes unable to quantify the sensitivities of a given derivative to *all* market instruments. Instead, we are forced to choose, often somewhat arbitrarily, which (much reduced) set of European swaptions we wish to use for model calibration and, ultimately, for vega reporting. While we can make a reasonably informed choice for some derivatives (e.g., vanilla Bermudan swaptions where we would use coterminal European swaptions and, possibly, caplets), many securities will not allow us to easily locate the dominant vega exposure locations, should these even exist in the first place. In contrast, models that either use a high number of Markov state variables (e.g., a multi-factor quasi-Gaussian model) or do not rely on volatility factorization at all (e.g., an LM model) will better preserve the full dimensionality of the volatility structure and hence, at least in theory, could be used to *tell* us, in

²We are ignoring the strike dimension for now.

an unambiguous way, which points of the volatility structure have impact on the value of a given security.

The discussion above is clearly intimately related to our earlier analysis of the debate surrounding local versus global calibration, see Section 14.5.5. With models that require product-specific volatility calibration, the choice of the calibration option set effectively decides in which buckets the vega will be reported, sensibly or not. On the other hand, with globally calibrated models with a fully flexible volatility structure, the model itself will ultimately determine the vega bucketing, in a manner that relies little on (possibly flawed) user intuition. The distinction between model types is fundamental, and irrespective of whether we ultimately choose to use product-specific or global calibration, the ability to discover, in a largely automated way, the set of European swaption volatilities that drive the value of any given derivative can often be essential to robust risk management of interest rate product portfolios. Of course, information uncovered this way could also be used to guide more robust and accurate product-specific calibrations for the local projection method (see Sections 18.4, 20.1.3, 20.2.1).

Our focus in this chapter is squarely on globally calibrated models, with the LM model being our primary example. The same techniques could, however, be applied to any globally-calibrated model underpinned with either a genuinely two-dimensional volatility structure, or one that approximates it closely, such as a multi-factor quasi-Gaussian model.

26.2 Review of Calibration

Let us start the technical discussion by recalling some notations from Section 14.5, and also introducing some new ones. We start with G (see Section 14.5.2 and in particular (14.41)), a subset of discretized instantaneous Libor volatilities which we regard as primary model parameters to be calibrated to market data. The $(N_t \times N_x)$ -dimensional matrix G is defined by a rectangular grid of times and tenors $\{t_i\} \times \{x_j\}$, $i = 1, \dots, N_t$, $j = 1, \dots, N_x$. For the purposes of this chapter, we denote by G^{full} the full grid of instantaneous Libor volatilities $\|\lambda_{n,k}\|$ for all n, k (see Section 14.5.3 and (14.42)). The matrix G is obtained from matrix G^{full} by selecting rows and columns that correspond to times $\{t_i\}$ and tenors $\{x_j\}$. As in Section 14.5.8, we assume that the calibration, or *benchmark*, set consists of all swaptions with expiries t_i and tenors x_j , $i = 1, \dots, N_t$, $j = 1, \dots, N_x$. On the other hand, we call the set of *all* at-the-money swaptions (i.e. swaptions with expiries T_i and tenors $T_j - T_i$ for all $i = 1, \dots, N-1$, $j = i+1, \dots, N-1$) the *full swaption set*.

Recall now the sample calibration algorithm of Section 14.5.7. Given a guess of G , we interpolate it to obtain G^{full} , which is then used to calculate model volatilities of swaptions in the benchmark set that we arrange in a matrix $\Lambda(G)$ with entries

$$(\Lambda(G))_{i,j}, \quad i = 1, \dots, N_t, \quad j = 1, \dots, N_x. \quad (26.2)$$

Given $\Lambda(G)$, an objective function $\mathcal{I}(G; \widehat{\Lambda})$ may then be constructed, typically involving a sum of precision and smoothness terms (see e.g. (14.51) and (14.54)), where the precision targets measure the distance between the model and market volatilities of swaptions in the benchmark set. Here we explicitly highlight the dependence of the objective function on market volatilities of swaptions in the benchmark set; these market volatilities are here assumed arranged in an $N_t \times N_x$ matrix $\widehat{\Lambda}_{i,j}$, $i = 1, \dots, N_t$, $j = 1, \dots, N_x$. The model calibration minimizes the objective function, resulting in a calibrated grid G^* of Libor volatilities, given by

$$G^*(\widehat{\Lambda}) = \operatorname{argmin}_G \mathcal{I}(G; \widehat{\Lambda}). \quad (26.3)$$

Once the model is calibrated, the value of a given derivative security, $V = V(G^*(\widehat{\Lambda}))$, may be calculated.

26.3 Vega Calculation Methods

Having formalized in Section 26.2 above the basic dependency structure of derivative security values on market volatilities, the key question is now how to establish sensitivities with respect to these volatilities. The next few sections outline several potential methods.

26.3.1 Direct Vega Calculations

26.3.1.1 Definition and Analysis

In the *direct method* for vega calculations, we simply apply a shock to the matrix of market swaption volatilities $\widehat{\Lambda}$, redo the model calibration, and reprice our security position. Let δ be an $N_t \times N_x$ matrix characterizing the shape of the chosen shock; then, the set of shocked market volatilities is given by the matrix

$$\widehat{\Lambda} + \epsilon \delta,$$

for some small $\epsilon > 0$. We proceed to calibrate the new grid of model volatilities $G^*(\widehat{\Lambda} + \epsilon \delta)$ by solving (26.3), i.e.

$$G^*(\widehat{\Lambda} + \epsilon \delta) = \operatorname{argmin}_G \mathcal{I}(G; \widehat{\Lambda} + \epsilon \delta),$$

and then estimate the (market) vega $\nu_{\text{mkt}}(\delta)$ in direction δ by the finite difference³

³See footnote 16 in Chapter 6 for a similar definition of sensitivity to a shock to a yield curve.

$$\nu_{\text{mkt}}(\delta) = \epsilon^{-1} \left(V(G^*(\widehat{\Lambda} + \epsilon\delta)) - V(G^*(\widehat{\Lambda})) \right) \approx \frac{d}{du} V(G^*(\widehat{\Lambda} + u\delta)) \Big|_{u=0}. \quad (26.4)$$

We note that here, and throughout, we think of vegas as pure derivatives, while in reality for reporting purposes the vega is often normalized to represent a change in value of a derivative that corresponds to, say, 1% change in the quoted market volatility.

The shock δ could take many forms, starting with the most basic *flat shock* (or “parallel shift”), where $\delta = \delta_{\text{flat}}$, $(\delta_{\text{flat}})_{i,j} = 1$ for all i, j . For more granular sensitivities, e.g. to measure sensitivities to individual market volatilities, we could use *bucketed shocks* $\delta = \delta_{n,m}$,

$$(\delta_{n,m})_{i,j} = 1_{\{i=n\}} 1_{\{j=m\}}. \quad (26.5)$$

A full collection of bucketed shocks — one for each swaption in the benchmark set — gives rise to a total of $N_t \cdot N_x$ so-called *bucketed vegas*⁴.

While it is often the goal to calculate vegas to all swaptions in the benchmark set — i.e. calculate sensitivities in directions $\delta_{n,m}$ in (26.5) for all n, m — it is not necessary to use the directions $\delta_{n,m}$ directly. As we have already seen in the context of interest rate deltas in Section 6.4, as long as we have the same number of directions ($N_t N_x$) that span the set $\{\delta_{n,m}\}$, then we can always express vegas in one basis from vegas in another basis by simple linear algebra. To give an example, consider the set of *running cumulative shocks* $\delta'_{n,m}$ given by

$$(\delta'_{n,m})_{i,j} = \begin{cases} 1, & i < n \text{ or } i = n, j \leq m, \\ 0, & \text{otherwise.} \end{cases} \quad (26.6)$$

Then, since

$$\delta'_{n,m} = \delta'_{n,m-1} + \delta_{n,m}$$

(with obvious modifications for $m = 1$), we have

$$\nu_{\text{mkt}}(\delta_{n,m}) = \nu_{\text{mkt}}(\delta'_{n,m}) - \nu_{\text{mkt}}(\delta'_{n,m-1}). \quad (26.7)$$

Hence, we can calculate $\nu_{\text{mkt}}(\delta'_{n,m})$ for all n, m using the algorithm described above, and then calculate all $\nu_{\text{mkt}}(\delta_{n,m})$ using (26.7).

Another sometimes used choice for the vega calculation basis is the set of all *cumulative shocks* $\delta''_{n,m}$ defined by

$$(\delta''_{n,m})_{i,j} = 1_{\{i \leq n\}} 1_{\{j \leq m\}} \quad (26.8)$$

for all n, m . Again, we can easily express $\delta_{n,m}$ in terms of the $\delta''_{n,m}$. The motivation for introducing alternative bases is similar to that for introducing

⁴Between the two extremes of the flat and bucketed shocks lie *row shocks* δ_n of the form $(\delta_n)_{i,j} = 1_{\{i=n\}}$, $n = 1, \dots, N_t$.

alternative ways of bumping the yield curve in Section 6.4.4: cumulative shocks as a rule lead to less distortion in the internal, model-specific volatility representation. We elaborate on this point later in the chapter.

While many variations of the direct vega method are possible, ultimately the accuracy and stability of vegas obtained by direct perturbation are rarely entirely satisfactory. The main reason is the fact that the calibration (26.3) is not exact, in the sense that the model does not exactly replicate all market volatilities of the swaptions in the benchmark set,

$$\Lambda(G^*(\widehat{\Lambda})) \neq \widehat{\Lambda}.$$

This imprecision is typically caused by the presence of regularization (smoothness) terms in the objective function, by usage of low-dimensional parametric forms for the volatility structure, or by other smoothness measures introduced to prevent overfitting of the model. For a well-designed calibration procedure, the resulting calibration errors are typically within bid-ask tolerances of market data⁵ and, consequently, are of little concern in securities pricing. However, the accuracy is often insufficient for vega calculations, since the typical size of the shock applied to market data (i.e., the magnitude of ϵ in (26.4)) is usually of the same order as the calibration errors. As a result, when calculating the vega calibration errors (the “noise”) might easily be of the same order of magnitude as the sensitivities themselves (the “signal”), making vegas too noisy to be useful. To improve on this, one can try increasing the size of ϵ in (26.4), but as described in Section 3.3.1 this leads to a bias relative to the true infinitesimal volatility sensitivity. More worryingly, applying shocks of a large magnitude to small subsets of the swaption volatility surface may result in an unrealistically choppy market data scenario to which the model can no longer calibrate properly.

The noise problems described above are less severe for global shocks than for local ones. For example, calculating the flat shift vega with $\delta = \delta_{\text{flat}}$ by direct methods is often possible, and in fact can serve as a benchmark and a reality check for the more advanced methods that we introduce later. The relatively good performance for global shocks is easy to understand, as they tend to preserve the distribution of calibration error among swaptions in the base and bumped scenario; i.e. we roughly have

$$\Lambda(G^*(\widehat{\Lambda})) - \widehat{\Lambda} \approx \Lambda(G^*(\widehat{\Lambda} + \delta)) - (\widehat{\Lambda} + \delta). \quad (26.9)$$

In other words, such shocks do not affect (too much) the calibration error for individual market volatilities; when calculating vegas by (26.4), the calibration errors therefore tend to cancel out. In fact, the introduction of cumulative shocks such as (26.6) and (26.8) can, in part, be motivated by the notion of keeping calibration errors relatively constant to ensure that

⁵These are typically in the order of 0.1% in implied volatility terms (with typical market swaption volatilities being in the 10–50% range).

(26.9) holds. While usage of a cumulative shock basis can, in fact, improve the vega noise somewhat, it still rarely produces satisfactory results.

Below, we elaborate a bit more on the noise issues plaguing the direct vega method. We should note that even if one were to find a remedy for the noise problem, the direct vega method may still be unattractive due to the need to repeatedly run a computationally intensive calibration algorithm for each shocked scenario. While the calibration algorithm of Section 14.5 is often relatively fast, if multiple scenarios are required, the total computation time per security can easily become impractically large.

26.3.1.2 Numerical Example

To demonstrate how the (basic) direct vega method performs on a simple example, we set up a 20 year one-factor LM model with 6 month Libor tenors, using a relatively coarse calibration grid: $\{t_i\} = \{x_j\} = \{1y, 5y, 10y, 15y\}$. For simplicity, the LM model is log-normal with flat Libor volatilities at 20%, i.e. $\lambda_{n,m} = 20\%$ for all n, m . The yield curve is also assumed flat, at a level of 5% continuously compounded.

In our model calibration, the swaption benchmark set consists of swaptions with expiry/tenor matching Libor volatilities in the matrix G , i.e. on the grid $\{t_i\} \times \{x_j\}$. We use a global calibration as outlined in Section 26.2 with smoothing weights (in expiry and tenor direction) set high enough to remove unwanted oscillations in the model volatility surface. Additionally, the vega shocks $\delta_{n,m}$, $n, m = 1, \dots, 4$, are assumed to be the bucketed shocks (26.5) applied to the 4×4 swaption matrix $\hat{\Lambda}$. In this and subsequent examples we consider three instruments:

1. 5y5y European swaption: a European payer swaption with strike 5%, expiry 5y and tenor 5y. Note that this swaption belongs to the benchmark set.
2. 3y7y European swaption: a European payer swaption with strike 5%, expiry 3y and tenor 7y. Note that this swaption does *not* belong to the benchmark set.
3. 10nc1 Bermudan swaption: a 10 no-call 1 (see Section 5.12) Bermudan payer swaption with annual exercise rights and a 5% strike.

All three derivatives have a notional of 1, and in all examples their values are calculated by Monte Carlo with 16,384 paths.

Table 26.1 shows the vegas obtained by the direct method for the 5y5y European swaption. As the 5y5y swaption is part of the benchmark set, we would expect a non-zero bucket vega number in only the 5y5y expiry/tenor bucket. While the largest vega exposure indeed does show up in this bucket, the table results are noisy and there are non-zero vegas in most other buckets as well.

Consider now our second test instrument, the 3y7y European swaption which is not in the benchmark set; its vegas are given in Table 26.2. As this

	1y	5y	10y	15y
1y	-0.1	0.2	0.0	0.0
5y	-0.8	16.5	0.3	0.1
10y	-0.1	0.8	-1.0	
15y	-0.6	2.5		

Table 26.1. Vegas by the direct method for the 5y5y European swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

swaption is not in the benchmark set, we here do not expect only a single non-zero bucket vega; instead, a well-behaved algorithm should produce non-zero numbers only in the four buckets that immediately surround the 3y7y point (the 1y5y, 1y10y, 5y5y, and 5y10y swaptions). As is evident from the table, however, the direct vega method assigns non-zero vegas to many other buckets as well, with some of the vegas being substantially negative.

	1y	5y	10y	15y
1y	-3.0	6.6	7.9	-0.7
5y	-0.2	7.0	2.3	-0.2
10y	0.0	0.3	-0.4	
15y	-0.2	0.9		

Table 26.2. Vegas by the direct method for the 3y7y European swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

Finally, let us look at vegas of the 10-nocall-1 Bermudan swaption. While results in Tables 26.1 and 26.2 primarily served to show the deficiencies of the direct vega method, the Bermudan swaptions vegas will be used as a useful benchmark for better methods we shall develop later in the chapter. Table 26.3 lists the relevant results; we notice that there are non-zero vega numbers in buckets corresponding to swaptions with total maturity (expiry + tenor) exceeding the 10 year life of the Bermudan swaption, so again we must conclude that the vega report in the table is affected by a significant amount of noise.

26.3.2 What is a Good Vega?

In the previous section, we pointed out some obvious deficiencies of vegas computed by direct perturbation methods. Before developing other methods

	1y	5y	10y	15y
1y	1.9	3.8	1.1	0.4
5y	6.0	6.8	-0.2	-0.1
10y	3.0	-0.6	-0.2	
15y	-0.2	1.1		

Table 26.3. Vegas by the direct method for the 10nc1 Bermudan swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

for calculating vegas, it is useful to first define what characteristics a method for computing vegas should ideally have. For starters, we obviously require that all computed risk sensitivities, including vegas, to be both stable and accurate. While stability can (and should, on a regular basis) be tested empirically by observing calculated risk measure over long periods of time, accuracy is often mode difficult to measure. One relevant metric for accuracy could be the performance of the P&L predict from Section 22.2.1, since accurately calculated risk measures typically imply low unpredicted P&L. While a P&L predict analysis is always useful, tests of vega accuracy in this manner may be inconclusive, as the P&L predict measures aggregate quality of all risk sensitivities and could be thrown off by inaccuracy of greeks other than the vega. As a consequence, it is often helpful to have more tailored measures of vega “goodness”; the list below contains several relevant measures. All the equalities below should be understood as “equality within tolerance”, where tolerances are typically determined by the requirements of the trading desk that uses the LM model.

1. Additivity. We normally would expect that

$$\nu_{\text{mkt}}(\delta_1 + \delta_2) = \nu_{\text{mkt}}(\delta_1) + \nu_{\text{mkt}}(\delta_2),$$

i.e. that applying two shocks together gives a vega that is a sum of vegas that correspond to the two individual shocks. As a particularly important case, the flat-shift vega should be reproduced as a sum of bucketed vegas:

$$\sum_{n,m} \nu_{\text{mkt}}(\delta_{n,m}) = \nu_{\text{mkt}}(\delta_{\text{flat}}).$$

2. Scaling. Scaling the size of a shock should scale the vega accordingly,

$$\nu_{\text{mkt}}(c\delta) = c\nu_{\text{mkt}}(\delta)$$

for a reasonable range of values of c , e.g. $c \in [0.5, 2]$. It is often also natural to require that the vega is invariant with respect to the sign of the bump, i.e. the equality holds with $c = -1$,

$$\nu_{\text{mkt}}(-\delta) = -\nu_{\text{mkt}}(\delta).$$

3. Locality. Our notion of vega locality is similar to the one used for yield curve perturbations in Chapter 6, and loosely requires that vega exposure “lives” where we expect it to. In this requirement, we can distinguish between a few variations:

- a) Benchmark set locality. The bucketed vegas calculated for a European swaption in the benchmark set are equal to zero everywhere except in the bucket that corresponds to the swaption itself. In other words, for a swaption with expiry t_i and tenor x_j ,

$$\nu_{\text{mkt}}(\delta_{n,m}) = 0 \text{ for } n \neq i, m \neq j$$

and $\nu_{\text{mkt}}(\delta_{i,j}) = \partial V_{\text{swaption},i,j} / \partial \hat{\lambda}_{i,j}$, where the right-hand side, the vega of the European swaption, is calculated in a vanilla model compatible with the LM model used. As we saw from Table 26.1, the direct method does not have benchmark set locality.

- b) Full swaption set locality. This is a stronger version of the previous point. For standard European swaptions that are *not* part of the benchmark set, we expect the vega to be non-zero only in the four buckets that surround the swaption in question. In particular, for a European swaption with expiry T_l and final swap maturity of T_k , we expect that

$$\nu_{\text{mkt}}(\delta_{n,m}) = 0 \text{ for } n \notin \{i-1, i\} \text{ and } m \notin \{j-1, j\},$$

where

$$i = \min \{a : t_a \geq T_l\}, \quad j = \min \{b : x_b \geq T_k - T_l\}.$$

Moreover we require that

$$\nu_{\text{mkt}}(\delta_{i-1,j-1}) + \nu_{\text{mkt}}(\delta_{i-1,j}) + \nu_{\text{mkt}}(\delta_{i,j-1}) + \nu_{\text{mkt}}(\delta_{i,j})$$

equals the vega of the European swaption in the compatible vanilla model. Again, the direct method does not satisfy this property as clear from Table 26.2.

- c) Exotic locality. For many exotics derivatives such as, for example, Bermudan swaptions, we know *a-priori* in which buckets the vega is supposed to reside (and often what sign it is supposed to have). For example, for a Bermudan swaption with final maturity T_k , we would expect no vega below the coterminal diagonal⁶:

$$\nu_{\text{mkt}}(\delta_{n,m}) = 0 \text{ if } t_n + x_m > T_k.$$

Given that the direct vega method fails simpler tests of locality, it is unlikely that it will respect the theoretical location (or sign) of Bermudan swaption vegas, an observation confirmed by Table 26.3.

⁶If the benchmark set does not include all coterminal swaptions for a given Bermudan swaption, non-zero vegas are still possible immediately below the coterminal diagonal due to interpolation effects.

4. Convergence. As with all quantities calculated by Monte Carlo methods, we expect vegas to converge to some value as we increase the number of paths. In particular, for the number of Monte Carlo paths N_{MC} used, the vegas calculated with N_{MC} paths should be within required tolerances compared to vegas calculated with $2N_{MC}$ paths, and vegas calculated with two different Monte Carlo seeds should be identical to within given tolerance.
5. Stability. Again, as a general requirement on values calculated by numerical methods, we expect vegas to vary smoothly with changing market inputs.

26.3.3 Indirect Vega Calculations

26.3.3.1 Definition and Analysis

While the mapping (26.3) of market volatilities to model volatilities involves non-linear optimization that adds noise, the reverse mapping of model volatilities to market volatilities (26.2) is typically done by direct application of swaption volatility approximation formulas and is, consequently, noiseless. Hence, it is natural to think that Jacobian techniques — which we have already encountered in Sections 6.4.3 and 22.1.4 — could be fruitfully applied here, with the exact mapping (26.2) used to define the transformation from model vegas to market vegas. To motivate the method we write, informally,

$$\frac{\partial V}{\partial G} = \frac{\partial V}{\partial \Lambda} \frac{\partial \Lambda}{\partial G}, \quad (26.10)$$

where on the left hand side we have (a vector of) model vegas, i.e. sensitivities with respect to changes in model volatilities, and on the right a product of (a vector of) market vegas and (a matrix of) sensitivities of swaption volatilities with respect to model parameters. As it is the market vegas we are interested in, we solve this linear system to obtain

$$\frac{\partial V}{\partial \Lambda} = \left(\frac{\partial \Lambda}{\partial G} \right)^{-1} \frac{\partial V}{\partial G}. \quad (26.11)$$

In this equation, $\partial \Lambda / \partial G$ can be computed analytically, whereas the term $\partial V / \partial G$ (the model vegas) normally must be computed by Monte Carlo methods.

Let us develop the ideas above a bit more carefully. For a given $N_t \times N_x$ perturbation matrix δ , we define the model vega $\nu_{mdl}(\delta)$ in direction δ by

$$\nu_{mdl}(\delta) = \epsilon^{-1} (V(G^* + \epsilon\delta) - V(G^*)) \approx \frac{d}{du} V(G^* + u\delta) \Big|_{u=0},$$

where $\epsilon > 0$ is a small number, and $G^* = G^*(\widehat{\Lambda})$ as before is the matrix of model volatilities calibrated to market.

Let us consider applying, to the model volatilities, a set of $N_t \cdot N_x$ shocks denoted by $\delta_{n,m}$ for $n = 1, \dots, N_t$, $m = 1, \dots, N_x$. These could be the unit shocks of (26.5), or any of the other families we introduced in Section 26.3.1. It often helps to think of these shocks as market data scenarios, with vega hedging being the exercise of finding weights for the hedging instruments (swaptions in the benchmark set) to neutralize as much as possible the sensitivity of a given security to the chosen scenarios. The sensitivity of the volatility of the (i,j) -th swaption in the benchmark set to “scenario” $\delta_{n,m}$ is given by

$$\frac{\partial \Lambda_{i,j}}{\partial \delta_{n,m}} \triangleq \epsilon^{-1} (\Lambda_{i,j}(G^* + \epsilon \delta_{n,m}) - \Lambda_{i,j}(G^*)) \approx \left. \frac{d}{du} \Lambda_{i,j}(G^* + u \delta_{n,m}) \right|_{u=0},$$

a quantity that can easily be calculated by differentiating the formula for approximation swaption volatility in the LM model with respect to model volatilities⁷. Hence, in its most basic form, the market vega matrix $(\nu_{\text{mkt}})_{i,j}$ can be introduced as the solution to the following least-squares minimization problem

$$\sum_{m=1}^{N_x} \sum_{n=1}^{N_t} \left(\nu_{\text{mdl}}(\delta_{n,m}) - \sum_{j=1}^{N_x} \sum_{i=1}^{N_t} (\nu_{\text{mkt}})_{i,j} \frac{\partial \Lambda_{i,j}}{\partial \delta_{n,m}} \right)^2 \rightarrow \min. \quad (26.12)$$

The definition (26.12) can be extended in a number of ways, along the same lines as in Section 6.4.3. We could, for instance, use a different number of scenarios than hedging instruments, either by supplying more scenarios or by utilizing only a subset of the benchmark set for hedging purposes. We could also use different weights for different scenarios, with higher weights applied to the scenarios we care more about. In addition, we could introduce regularization weights to, say, penalize hedging positions of excessive size. A reasonably general definition of market vegas in the *indirect method* for vega calculations is then given by the solution to the following least-squares minimization problem (compare to (6.29))

$$\begin{aligned} \sum_{n,m} W_{n,m}^2 & \left(\nu_{\text{mdl}}(\delta_{n,m}) - \sum_{i,j} (\nu_{\text{mkt}})_{i,j} \frac{\partial \Lambda_{i,j}}{\partial \delta_{n,m}} \right)^2 \\ & + \sum_{i,j} U_{i,j}^2 ((\nu_{\text{mkt}})_{i,j})^2 \rightarrow \min, \end{aligned} \quad (26.13)$$

where $W_{n,m}$ are weights applied to different scenarios and $U_{i,j}$ are penalty weights for different hedges. This problem can be formulated in matrix form

⁷See Section 14.4.2 for examples of such formulas. Notice that the Jacobian $\{\partial \Lambda_{i,j} / \partial \delta_{n,m}\}$ is often available for free as part of the initial calibration of the LM model, especially if calibration relies on a gradient-based optimization method.

(see e.g. (6.30)) and could be solved by standard methods of linear algebra, as in (6.31).

The indirect method for computing vegas avoids noisy (and costly) model recalibration, and often results in a marked improvement over the direct vega method of Section 26.3.1. Still, the results are not perfect, as the vegas calculated by the indirect method will often violate several of the criteria for good vegas listed in Section 26.3.2. In particular, while the indirect vega method tends to satisfy the additivity and scalability properties, it is often quite noisy and exhibits unsatisfactory convergence and stability.

Stability and convergence issues could in principle be addressed by modifying (26.12) (or (26.13)). Specifically, we can add penalty terms that would promote smoothness of market vegas in expiry and tenor dimensions, in the same spirit as we smooth model volatilities during calibration, see Section 14.5.6 and in particular equation (14.51). For example, to promote first-order smoothness we can change (26.12) to

$$\begin{aligned} \sum_{n,m} & \left(\nu_{\text{mdl}}(\delta_{n,m}) - \sum_{i,j} (\nu_{\text{mkt}})_{i,j} \frac{\partial \Lambda_{i,j}}{\partial \delta_{n,m}} \right)^2 \\ & + w_{\partial t} \sum_{i,j} \left((\nu_{\text{mkt}})_{i,j} - (\nu_{\text{mkt}})_{i-1,j} \right)^2 \\ & + w_{\partial x} \sum_{i,j} \left((\nu_{\text{mkt}})_{i,j} - (\nu_{\text{mkt}})_{i,j-1} \right)^2 \rightarrow \min. \end{aligned}$$

We can also add second-order smoothing terms along the same lines as in (26.12). These modifications do not make the minimization problem any harder to solve, as it remains quadratic.

As one would expect, the addition of smoothing terms often significantly improves the stability and convergence characteristics of the indirect method. Unfortunately, however, extra smoothing destroys the locality of vegas: if we apply the indirect method with smoothing to a European swaption in the calibration set, then its vega will “leak out” from its native bucket to other nearby buckets. Despite this issue, we believe that the indirect method with smoothing (and its variants) is widely used in industry for vega calculations in LM models. The locality problems of the method are either ignored on pragmatic grounds (in effect choosing the lesser of two evils, non-locality over instability), or justified by the fact that vegas for actual trading books tend to be spread out over all buckets anyway. Such arguments are obviously not entirely convincing, and assigning vega to buckets where there should be none has strong negative implications for hedging and P&L explain.

26.3.3.2 Numerical Example and Performance Analysis

In order to later improve on the indirect vega method, let us first gain some understanding of the actual performance of the method. For concreteness,

we continue with the LM model example from Section 26.3.1.2 using Monte Carlo with 16,384 paths (the same as in the examples of Section 26.3.1.2), and apply shocks $\delta_{n,m}$, $n, m = 1, \dots, 4$, (assumed to be the bucketed shocks (26.5)) to the 4×4 Libor volatility matrix G . In all tests, we do not use smoothing, i.e. we compute vegas by applying the basic equation (26.12).

Considering first the 5y5y European swaption defined in Section 26.3.1.2, vegas computed by the indirect method are given in Table 26.4. Comparison with Table 26.1 shows that there is a marked improvement over the direct method, but a fair amount of noise is still apparent. For example, the vega of -0.9bp in the 5y1y bucket is clearly incorrect.

	1y	5y	10y	15y
1y	0.0	0.0	0.1	0.0
5y	-0.9	17.0	0.2	0.0
10y	0.0	0.0	0.0	
15y	0.0	0.0		

Table 26.4. Vegas by the indirect method for the 5y5y European swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

Table 26.5 lists vegas for a 3y7y European swaption. Again, we see an improvement over the vegas calculated by the direct method in Table 26.2, but non-zero values in the 1 year tenor column again indicate that the method is not completely satisfactory.

	1y	5y	10y	15y
1y	-2.9	6.4	8.2	-0.9
5y	-0.2	7.3	1.9	0.0
10y	0.0	0.0	0.0	
15y	0.0	0.0		

Table 26.5. Vegas by the indirect method for the 3y7y European swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

Finally, we consider the 10nc1 Bermudan swaption, the indirect method vegas of which are shown in Table 26.6. While overall somewhat cleaner than the vegas in Table 26.3, negative values in the 15 year expiry row indicate the presence of noise.

	1y	5y	10y	15y
1y	1.9	4.3	1.9	0.2
5y	6.2	5.1	1.4	0.5
10y	2.4	0.2	-0.3	
15y	-0.4	-0.2		

Table 26.6. Vegas by the indirect method for the 10nc1 Bermudan swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

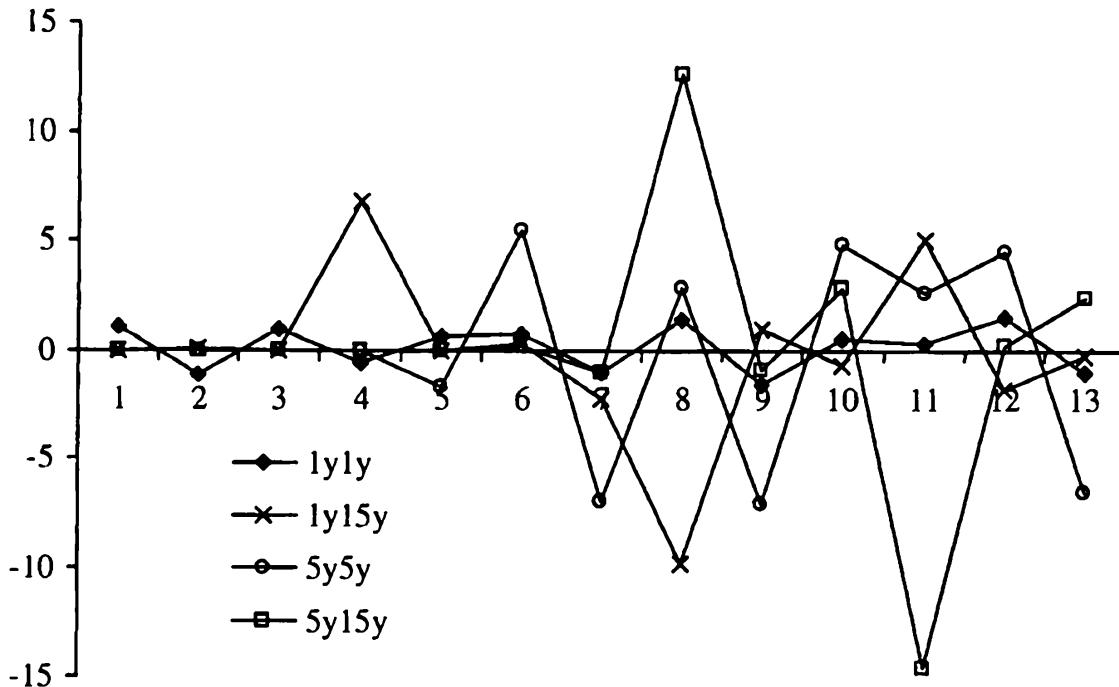
To understand why the performance of the indirect vega method is unimpressive, let us first note that the solution to (26.12) is given by (26.11), where $(\nu_{\text{mdl}})_{i,j}$ is arranged into a vector⁸ $\partial V/\partial G$, $(\nu_{\text{mkt}})_{i,j}$ is arranged into a vector $\partial V/\partial \Lambda$, and the matrix $\partial \Lambda/\partial G$ is an appropriately arranged matrix of sensitivities of swaption volatilities with respect to Libor volatilities. Some of the buckets — namely 10y15y, 15y10y and 15y15y — are outside of the 20 year model horizon and can be discarded, so the dimension of the matrix $\partial \Lambda/\partial G$ is 13×13 .

The vector of model vegas $\partial V/\partial G$ is calculated numerically in Monte Carlo, by perturbing individual entries in the matrix G and repricing the derivative⁹. This procedure induces noise in the model vegas which will be transmitted into market vegas through (26.11), by multiplication with the inverse of the matrix $\partial \Lambda/\partial G$. Let us look at this matrix in more detail, as clearly its properties will influence the propagation of noise. Using our test setup above, Figure 26.1 represents the matrix $(\partial \Lambda/\partial G)^{-1}$ graphically, with a few selected columns plotted as separate lines, showing how a shock to a particular swaption volatility in the benchmark set affects all Libor volatilities in G . Each market vega is obtained by adding up all model vegas weighted by the values of a corresponding column (line in the figure).

Two things are apparent in Figure 26.1. First, each column is rather “wiggly”, with positive and negative values alternating in a ringing pattern. This behavior — which is not due to numerical noise, since the calculation of the matrix $(\partial \Lambda/\partial G)^{-1}$ is exact — is likely to exacerbate any noise in model vegas. Second, the (absolute values of) values in the matrix are quite high, reaching values of 10 to 15. This is significant, as any noise in the market vegas is then essentially multiplied by a factor of 10 to 15. This noise-amplifying effect is confirmed by looking at the eigenvalues of the matrix $(\partial \Lambda/\partial G)^{-1}$: the lowest and highest eigenvalues equal 1 and 18, respectively. If we make the reasonable assumption that the level of

⁸This could be done arbitrarily, but for concreteness we do it in row-major order, i.e. rows of the matrix are stacked end-to-end to come up with a vector.

⁹Of course, the pathwise or likelihood differentiation methods of Chapter 24 could have been used here as well.

Fig. 26.1. Inverse Jacobian for Indirect Vega Method

Notes: Each line represents how a shock to a given swaption volatility affects forward Libor (i.e. model) volatilities in matrix G arranged in a row-major order. The lines are graphic representations of (a few selected) columns from the inverse Jacobian $(\partial \Lambda / \partial G)^{-1}$, see text for details.

accuracy in calculating model vegas is roughly the same as for deltas, then the accuracy in *market* vegas could be 10 or even 20 times worse. And this is only for our simplistic example — in real applications with models of longer tenors and larger benchmark sets, the noise amplification factor could easily be in the hundreds, making the indirect vega method perform poorly.

Incidentally, looking at the matrix $(\partial \Lambda / \partial G)^{-1}$ in Figure 26.1 also sheds more light on the poor performance of the *direct* method of calculating vegas, as described earlier in Section 26.3.1. When a particular swaption volatility is shocked and the model is recalibrated, Figure 26.1 shows that the resulting model will effectively have its Libor rate volatilities severely distorted by a large shock with irregular shape. For example, a perturbation of 1% to a 5y5y swaption volatility would move some of the Libor volatilities by almost 15% (and others by -15%). Clearly, with shocks of this size we cannot hope to accurately capture only first-order sensitivity, as second- and higher-order effects will pollute the vega we are trying to calculate.

26.3.4 Hybrid Vega Calculations

26.3.4.1 Definition and Analysis

In Section 26.3.3 we identified poor numerical invertibility (also known as *stiffness*) of the matrix $\partial A / \partial G$ as the main reason for poor performance of the basic indirect method for vega calculations. This stiffness primarily arises from the usage of shocks to model volatilities that do not adequately take into consideration the dependence of swaption volatilities on Libor volatilities. To improve the indirect vega method, it is therefore natural to change our set of simple bucketed shocks in Libor volatilities to a set of shaped shocks that will result in a better Jacobian, with less ringing and smaller noise amplification factors than in Figure 26.1.

One good choice for the Jacobian would be a unit matrix, which is both perfectly smooth and involves no amplification of noise. A unit Jacobian matrix will arise only if we use Libor volatility shocks that correspond to shocks of individual swaption volatilities in the benchmark set. It may appear that this line of reasoning simply leads us back to the direct method of vega calculations, but we here make a subtle but critical distinction: instead of outright shocking swaption volatilities and recalibrating the model, we instead construct Libor volatility shocks that are approximately equivalent to bucketed swaption volatility shocks, and then apply these shocks through the Jacobian technique outlined earlier. Avoiding recalibration and carefully controlling the shape of shocks to the Libor rate volatility surface not only leads to better computational performance, it also ultimately will lead to better vega quality¹⁰, in the sense defined in Section 26.3.2.

In light of the discussion above, the key problem we have to deal with is how to construct, in a noise-free manner, a shaped shock to Libor volatilities that approximates a shock to a particular European swaption. Here, the bootstrap LM model calibration presented in Section 14.5.8 turns out to be useful. Recall that the idea of bootstrap (or cascade) calibration is to find the instantaneous volatility of each Libor rate over each time period one at a time by solving a quadratic equation, a procedure that is enabled by doing the calculations in a certain (row-major) order. As pointed out in Section 14.5.8, bootstrap calibration is normally not suitable for a full calibration to market data, as market volatilities of swaptions typically come with some amount of noise in them or, at any rate, are not guaranteed to change smoothly across expiries and tenors. As exemplified by Figure 26.1, this leads to rapid accumulation of noise in Libor volatilities during the bootstrap and almost unavoidable calibration failure (where quadratic equations fail to have real roots). On the other hand, if the input volatilities happened to be smooth and “compatible” with a Libor market model, there would be no reason why the bootstrap calibration would not work. This suggests

¹⁰An observation also made by Pietersz and Pelsser [2004], although in a somewhat different context.

that we should apply swaption volatility shocks not to market values of swaption volatilities, but to the implied swaption volatilities returned by the calibrated model. The latter are fundamentally compatible with an LM model and, assuming that a reasonable amount of smoothing was enforced in the calibration norm, smooth enough for the bootstrap method to work.

These ideas lead us to the following *hybrid method* for calculating vegas, combining features of both the direct and indirect methods.

1. Calibrate the LM model to market data, i.e. obtain G^* from $\widehat{\Lambda}$, using our global calibration method.
2. Calculate Λ , the model-implied swaption volatilities.
3. Fix expiry t_n and tenor x_m .
4. Apply a unit shock $\alpha_{n,m}\delta_{n,m}$ with $\delta_{n,m}$ shaped as in (26.5), for sufficiently small¹¹ $\alpha_{n,m} > 0$ to Λ .
5. Bootstrap calibrate an LM model to swaption volatilities $\Lambda + \alpha_{n,m}\delta_{n,m}$, to obtain a matrix of shocked Libor volatilities $G^{*,n,m}$.
6. Calculate a Libor shock $\delta_{\text{Libor},n,m}$ by $\delta_{\text{Libor},n,m} = \beta_{n,m}(G^{*,n,m} - G^*)$, where the scaling constant $\beta_{n,m} \neq 0$ is chosen so that $\max_{i,j} |(\delta_{\text{Libor},n,m})_{i,j}| \leq \epsilon$ for a small $\epsilon > 0$.
7. Repeat Steps 3–6 for all expiries and tenors, and save all shocks $\{\delta_{\text{Libor},n,m}\}$.
8. Apply the indirect method of Section 26.3.3 with the collection of Libor volatility shocks $\{\delta_{\text{Libor},n,m}\}$.

The Jacobian matrix $\partial\Lambda_{i,j}/\partial\delta_{\text{Libor},n,m}$ will here be exactly diagonal, with the element $\alpha_{n,m}\beta_{n,m}$ on the diagonal in the position determined by the ordering of swaptions in the benchmark set; the inverse transformation (26.11) will amount to an appropriate scaling of each model vega. Of course, the choice of $\delta_{n,m}$'s in Step 4 is not unique, and instead of (26.5) we could have used other families such as (26.6) or (26.8). In this case the Jacobian would no longer be diagonal, but otherwise the method would still work the same way.

Let us discuss now the choice of various constants that appear in the algorithm. In Step 4 we need a choice for the positive constant $\alpha_{n,m} > 0$. The idea here is to apply a constant small enough that the bootstrap calibration of Step 5 works. Clearly for $\alpha_{n,m} = 0$ this is the case, so there exists a small enough $\alpha_{n,m} > 0$ that satisfies this criteria. On the other hand, we should not choose $\alpha_{n,m}$ too small as it may adversely affect the Monte Carlo simulation error when computing relevant finite differences (see Section 23.2). In practice, we may start with some reasonable large value of $\alpha_{n,m}$, say 1%, and attempt the bootstrap. If this fails, we reduce $\alpha_{n,m}$ by half and try again — and so on until we find the value of $\alpha_{n,m}$ that allows the bootstrap to succeed.

¹¹We comment below on the choice of $\alpha_{n,m}$ as well as other required constants.

As for the constant $\epsilon > 0$ required in Step 6, we should choose it in a way that ensures that shocks to Libor volatilities are small enough to prevent significant second-order effects to show up in vegas, yet big enough to control the level of Monte Carlo error in the numerically calculated sensitivity. A reasonable choice here is to set ϵ somewhere between 0.1% to 1%.

26.3.4.2 Numerical Example

To present test results for the hybrid vega method, we continue the numerical example of Sections 26.3.1.2 and 26.3.3.2. Looking first at the 5y5y European swaption, the hybrid method vegas are listed in Table 26.7. As we can see, results are much improved compared to the direct and indirect methods (see Tables 26.1 and 26.4, respectively) with very little noise and the only significant vega correctly showing up in the 5y5y bucket (as expected).

	1y	5y	10y	15y
1y	0.0	0.0	0.1	0.0
5y	-0.1	16.3	0.0	0.0
10y	0.0	0.0	0.0	
15y	0.0	0.0		

Table 26.7. Vegas by the hybrid method for the 5y5y European swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

Similar good results are obtained for the 3y7y European swaption, as shown in Table 26.8. Unlike the results in Tables 26.2 and 26.5, there is here hardly any noise visible outside of the four neighboring buckets where we expect the vega to be located.

	1y	5y	10y	15y
1y	0.0	5.3	4.9	0.0
5y	0.0	5.2	4.5	0.0
10y	0.0	0.0	0.0	
15y	0.0	0.0		

Table 26.8. Vegas by the hybrid method for the 3y7y European swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

Finally, Table 26.9 lists hybrid method vegas for a Bermudan swaption. Once again, we only see vega where it is expected to be, in contrast to Tables 26.3 and 26.6.

	1y	5y	10y	15y
1y	3.3	3.4	2.0	0.0
5y	6.5	4.6	0.6	0.0
10y	2.8	0.5	0.0	
15y	0.0	0.0		

Table 26.9. Vegas by the hybrid method for the 10nc1 Bermudan swaption as defined in the text, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

26.4 Skew and Smile Vegas

So far, we have focused our discussion of vega on the computation of sensitivities to at-the-money swaption volatilities. Together with the correlation sensitivities that we touch upon in Section 26.5, this comprises the full set of volatility sensitivities in models that parametrize each swaption volatility smile with a single number representing the overall level of volatilities across all strikes (such as a log-normal LM model). In models with richer volatility smile parameterizations, such vegas also have a meaningful interpretation as sensitivities to parallel shifts of volatility smiles of each swaption in the benchmark set. For such models, however, there are other volatility sensitivities that so far have been left out of the discussion, namely the sensitivities to changes in *shapes* of swaption volatility smiles. Sensitivities to changes in the volatility smile slope and curvature are often denoted *skew vegas* and *smile vegas*, respectively.

As with ATM vegas, in principle there are skew and smile vegas for each swaption in the benchmark set. However, it is rarely a requirement that one be able to calculate them all individually, as ATM vegas capture the majority of volatility sensitivity. More often, what is required are aggregated measures of skew and smile risk, such as a single number that corresponds to a change in slope or curvature of *all* volatility smiles together. For such aggregated measures of risk, brute-force recalibration and recomputation along the lines of the direct vega method is often sufficient. Moreover, it is typically more useful to use a scenario-based approach with large slope or curvature shocks, rather than true first-order differentiation. For example, in a displaced log-normal LM model (see Table 14.1 in Section 14.2.4) one

can switch the skew parameter from 1 (log-normal) to 0 (Gaussian) to get a good idea of the impact of the slope of volatility smile.

In the off-chance that bucketed skew/smile exposure is required, the indirect (Section 26.3.3) method or the hybrid (Section 26.3.4) approach that we developed for the ATM vegas could often be reused. In some cases skew and smile sensitivities are even easier to calculate than ATM vegas, due to a simpler connection between the model and market parameters. For example, the term swaption skew in a displaced log-normal LM model is a linear function of instantaneous Libor skews, see Section 15.2, making the Jacobian-type methods particularly easy to apply. We do not go into further detail here, as the mechanics of these calculations should be clear to the reader by now.

26.5 Vegas and Correlations

Earlier in the chapter we used a one-factor version of the LM model in our numerical examples, but in practice we are often more interested in calculating vegas in multi-factor LM models. For a q -factor LM model, yield curve dynamics are characterized by factor volatilities, i.e. q -dimensional vectors $\lambda_k(T_n)$ associated with each Libor rate $L_k(t)$ and each time period $(T_{n-1}, T_n]$. As we recall from Section 14.5.4, these are constructed from the volatility norm $\|\lambda_{n,k}\|$ (which we denoted by G^{full} in Section 26.2) and instantaneous correlations of Libor rates. In the indirect and hybrid methods of Sections 26.3.3 and 26.3.4 we shocked the elements of G (and, ultimately, G^{full}), with the understanding that instantaneous correlations of Libor rates remained fixed in all scenarios. In the direct method of Section 26.3.1, we referenced the sample calibration algorithm of Section 14.5.7, which implicitly assumed that the instantaneous correlations of Libor rates were untouched while perturbing swaption volatilities. So, in all three methods we so far have calculated interest rate vegas under the assumption that instantaneous Libor correlations are kept constant when forming the derivative with respect to volatility. While not unreasonable, this choice is not unique and several viable alternatives exist. We discuss some of these in this section.

26.5.1 Term Correlation Effects

While the correlation structure in the LM model is typically captured through a parameterization of instantaneous Libor correlations, the prices of traded correlation-sensitive instruments — CMS spread options, in particular — depend more directly on *term* correlations of *swap* rates (see Section 14.4.3.1). Importantly, when Libor volatilities are changed with instantaneous Libor correlations kept constant, term correlations of swap rates will generally change quite significantly. This effect should be intuitively clear and is a

consequence of the dependence of the formula for term correlation in Section 14.4.3.1 on Libor volatilities.

To demonstrate the magnitude of the vega effect on term correlations, we continue the numerical example of Sections 26.3.1.2, 26.3.3.2 and 26.3.4.2, but now extend our setup to a 10-factor LM model with the instantaneous Libor correlations parameterized by a function of the form (14.19) (with $\rho_\infty = 0.5$, $a_0 = 0.42$, $a_\infty = 0$, $\kappa = 0.08$) and instantaneous Libor volatilities fixed at 20%. For concreteness, let us study the sensitivity of the term correlation between the 10 and 1 year swap rates over a 10 year horizon ($\rho_{term}(0, 10)$ in the notation of Section 14.4.3.1). The base value of this correlation in our setup is about 83%. As demonstrated in Table 26.10, $\rho_{term}(0, 10)$ is quite sensitive to shocks to some of the volatilities. For example, a shock of 1% to the volatility of a 10y10y swaption would change this term correlation by -0.85%, which is highly significant.

	1y	5y	10y	15y
1y	0.1	-6.2	-3.1	2.0
5y	-5.0	-5.6	7.0	4.7
10y	14.1	45.7	-85.3	
15y	0.0	0.0		

Table 26.10. Sensitivity for the 10y1y term (10 years) swap rate correlation, in basis points (1bp = 10^{-4}) per 1% shift in volatility of each swaption in the benchmark set. All numbers are computed using the hybrid method in Section 26.3.4.2. Rows are expiries and columns are tenors of swaptions in the benchmark set.

26.5.2 What Correlations should be Kept Constant?

Since term swap rate correlations change under volatility shocks when instantaneous Libor correlations are fixed, we could instead decide to keep term swap rate correlations constant (while allowing instantaneous Libor correlations to move) under volatility shocks; this choice would lead to different vegas, of course. As we discussed in Chapter 22, this ambivalence is not unique to the problem of calculating vegas, and we often need to decide which quantities to keep constant and which to let float when calculating risk sensitivities. Ultimately, such decisions are often driven by traders' preferences for risk representation, or by the types of hedging strategies that they want to pursue. In making these decisions, traders generally (and reasonably) tend to emphasize the issue of *consistency* across different products.

A typical interest rate exotics trading desk will trade correlation-sensitive exotics (e.g., CMS spread TARNs, see Section 5.13.3), as well as vanilla

spread options. The exotic derivatives will often be risk managed in an LM model, while for spread options the desk may use a simpler vanilla model, as discussed in Chapter 17. These two models will typically have different (internal) correlation parameters: the LM model will use instantaneous Libor correlations, while a vanilla model (based, say, on a Gaussian copula) will use a term correlation between swap rates as an input. While it is natural for each model to keep its internal correlation parameters constant when calculating vegas, doing so would lead to inconsistency in the definition of vegas between the exotic derivatives and their vanilla hedges. Such inconsistency is typically quite dangerous as it could lead to a position that is deemed hedged, but in fact has an outright exposure.

In the example above, as well as many other similar situations, arguably the easiest way to maintain consistency is to use the more general model (here, the LM model) as the risk “engine” for all products in the book, exotic or not. In an LM model setting, we would then need to compute the volatility sensitivity of both exotics and vanilla securities assuming fixed instantaneous Libor correlations. For the vanilla securities, computation of this sensitivity could be done by either outright valuation of vanilla spread options in a LM model or, perhaps more pragmatically, by calculating the volatility shock impacts on the relevant term correlations of swap rates and combining the results (Jacobian-style) with known correlation sensitivities of the vanilla model. To complement the resulting vega report, it would be natural to also report correlation sensitivities, by calculating sensitivities of the portfolio to instantaneous Libor correlations¹².

While not entirely without merit, the approach outlined above has its limitations. Of course, if the portfolio is fully hedged to both vegas (under the assumption of constant Libor correlations) and to instantaneous Libor correlations, then it has no volatility risk, irrespective of how we define vega. However, a fully hedged position is rarely, if ever, achieved, in which case reported volatility and correlation sensitivities are used as a monitoring tool. As we have commented before, it is often much easier for traders to understand sensitivities expressed in terms of traded quantities, rather than in terms of non-traded quantities such as instantaneous Libor correlations. Traders therefore typically have a strong preference for seeing their correlation risk expressed in terms of market-implied swap rate correlations, which, for consistency reasons, dictates that vegas should be calculated under the assumption that market (and not model) correlations are kept constant. In addition, we should note that the LM model typically uses a fairly parsimonious correlation parameterization, often comprised of just a handful of numbers (see Section 14.3.2). Hence, correlation risk produced by the LM model would tend to be insufficiently granular for risk-managing vanilla spread options which often are quite liquid for a range of expiries

¹²But *not* sensitivities to term correlations of swap rates, which would lead to another inconsistency, with double-counting of risk.

and a reasonably large number of swap rate pairs. This issue also favors using term swap rate correlations for risk management purposes.

26.5.3 Vegas with Fixed Term Correlations

In the last section we made the case for holding term correlations of swap rates fixed when computing vegas. Let us discuss how to turn this idea into practice, by suitably modifying the various computational methods discussed earlier in this chapter. The direct method of Section 26.3.1 is the easiest to modify: all we need to do is to add the relevant¹³ term swap rate correlations as targets in the basic model calibration, with the calibration algorithm extended along the lines of Section 14.5.9. Note that it is spread option *correlations* that should be the calibration targets and not spread option *values* — we certainly expect the values of spread options to change under different volatility scenarios, even as we keep correlations constant. While it is easy to extend the direct method, its limitations with regards to the quality of vegas produced remain (or are amplified, most likely), and consequently this approach is not recommended.

Extending the indirect (or hybrid, as the procedure is more or less the same) method to control spread option correlations is somewhat more difficult. One naive choice would involve appending a correlation calibration to every shock of Libor volatilities, to ensure that term swap correlations stay fixed after each perturbation of volatilities. In other words, after applying a shock to G , we would then proceed to solve for new parameters to the instantaneous Libor correlation function in order to remain in calibration with term swap rate correlations. In this approach, we would need to run a separate optimization problem for each model vega shock, which most likely would make the method prohibitively slow and introduce extra noise due to non-exact nature of the solution of the optimization problem.

Our preferred method for extending the indirect vega computation is, once again, based on Jacobian methods. In this approach, we would i) apply shocks to model volatilities *and* to model correlations, ii) calculate the value of a derivative as well as changes to swaption volatilities and swap rate correlations, and iii) manipulate these quantities to obtain the vegas. Let us present the blueprint of the scheme using the stylized notations of (26.10) — we trust that the reader can expand our presentation into a workable algorithm.

Let ξ be the vector of instantaneous Libor correlations, and ρ the vector of term swap rate correlations (see footnote 13). We recognize security value and market data dependence on model data through the notations

¹³Note that it is impractical to include *all* swap rate correlations in the calibration set. Instead, one would typically choose a set (or perhaps a few sets) of correlations of two specific swap rate tenors, such as 10 year and 2 year, over a collection of time periods.

$$V = V(G, \xi), \quad \Lambda = \Lambda(G, \xi), \quad \rho = \rho(G, \xi).$$

Our goal is to compute the vector

$$\left. \frac{\partial V}{\partial \Lambda} \right|_{\rho=\text{const}}, \quad (26.14)$$

i.e. sensitivities to market swaption volatility shocks, keeping market correlations constant. Implicitly, G and ξ are functions of Λ and ρ ,

$$G = G(\Lambda, \rho), \quad \xi = \xi(\Lambda, \rho),$$

and so is V ,

$$V(\Lambda, \rho) = V(G(\Lambda, \rho), \xi(\Lambda, \rho)). \quad (26.15)$$

By an application of the chain rule to (26.15), we get

$$\left. \frac{\partial V}{\partial \Lambda} \right|_{\rho=\text{const}} = \frac{\partial V}{\partial G} \frac{\partial G}{\partial \Lambda} + \frac{\partial V}{\partial \xi} \frac{\partial \xi}{\partial \Lambda}. \quad (26.16)$$

Here, the sensitivities $\partial V / \partial G$ and $\partial V / \partial \xi$ may be obtained by application of model parameter shocks to the valuation of the derivative. $\partial G / \partial \Lambda$ and $\partial \xi / \partial \Lambda$ can be found by inverting the (full) Jacobian (inversion should be understood in the generalized least-squares sense as in (26.12) vs. (26.11) as the matrices involved may not even be square),

$$\begin{pmatrix} \partial G / \partial \Lambda & \partial G / \partial \rho \\ \partial \xi / \partial \Lambda & \partial \xi / \partial \rho \end{pmatrix} = \begin{pmatrix} \partial \Lambda / \partial G & \partial \Lambda / \partial \xi \\ \partial \rho / \partial G & \partial \rho / \partial \xi \end{pmatrix}^{-1}. \quad (26.17)$$

The matrix on the right-hand side of (26.17) is obtained by applying shocks to model volatilities and correlations, following the same approach as outlined for the indirect and hybrid vega methods. As a by-product of this calculation we also conveniently obtain risk sensitivities with respect to market *correlations*, since

$$\left. \frac{\partial V}{\partial \rho} \right|_{\Lambda=\text{const}} = \frac{\partial V}{\partial G} \frac{\partial G}{\partial \rho} + \frac{\partial V}{\partial \xi} \frac{\partial \xi}{\partial \rho}, \quad (26.18)$$

where $\partial G / \partial \rho$ and $\partial \xi / \partial \rho$ are obtained in (26.17).

26.5.4 Numerical Example

To demonstrate the difference between various definitions of vegas, we look at a simple, single 10 year option on the spread between 10 year and 1 year swap rates. We calculate its vegas in the 10-factor LM model used in Section 26.5.1, under the assumption of constant instantaneous Libor correlations (Table 26.11) and constant term swap rate correlations (Table

26.12). The second method puts all vega (apart from some minor noise) into the 10y1y and 10y10y buckets, unlike the first method which assigns significant vega to, for example, the 10y5y bucket. Arguably, most traders would consider the vega in Table 26.12 more intuitive, as it is exactly the shape of the vega profile that one would obtain for this spread option from a typical vanilla model.

	1y	5y	10y	15y
1y	0.00	0.10	0.05	-0.03
5y	0.08	0.10	-0.12	-0.08
10y	2.10	-0.77	2.00	
15y	0.00	0.00		

Table 26.11. Vegas by the hybrid method for the 10y option on the spread between 10y and 1y swap rates while keeping instantaneous Libor correlations constant, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

	1y	5y	10y	15y
1y	0.00	-0.01	0.00	0.00
5y	-0.01	-0.01	0.01	0.01
10y	2.36	0.05	0.47	
15y	0.00	0.00		

Table 26.12. Vegas by the hybrid method for the 10y option on the spread between 10y and 1y swap rate while keeping 10y term swap rate correlation between 10y and 1y swap rates constant, in basis points ($1\text{bp} = 10^{-4}$) per 1% shift in volatility of each swaption in the benchmark set. Rows are expiries and columns are tenors of swaptions in the benchmark set.

26.6 Deltas with Backbone

As we saw in Section 26.5, the need for consistency between exotic and vanilla models often drives the definitions of risk sensitivities. Such consistency requirements, it turns out, also affect calculations of deltas (and, of course, gammas) in the LM and other term structure models. To describe this effect in more detail, we first recall the discussion of Section 16.1.2 and, in particular, the fact that vanilla models are sometimes set up to attribute some user-specified amount of the vega to delta. If such a procedure is used,

it would be useful to ensure that the deltas computed in models for more exotic derivatives have the same meaning as in the vanilla model. In essence, this would require a link in the exotic model between the volatility smile and the level of rates.

Sometimes vanilla-exotic delta consistency is ensured automatically, as a consequence of the choice of the models in use. For example, the vanilla SV model (16.8)–(16.9) is naturally consistent with the SV LM model (14.15)–(14.16). On the other hand, if we start adjusting the backbone of the vanilla model as in Section 16.1.2, the consistency would often be lost. For example, were we to use a vanilla model of the type (16.5), we would need to modify the volatility terms for SDEs for Libor rates under the LM model to have the same form, e.g.

$$dL_n(t) = O(dt) + \left(b_n(t)L_n(t) + (m - b(t)) L_n(0) + (1 - m)L \right) \\ \times \lambda_n(t)^T dW(t), \quad n = 1, \dots, N - 1,$$

for some mixing m and level L .

A more complicated situation would arise were we to use a vanilla delta convention without a natural exotic counterpart, such as the SABR model of Section 8.6 or the SVI interpolation rule of Section 16.1.5. In these cases, it would be difficult to “internalize” the same smile move logic in the LM model dynamics. Fortunately, we can use an external brute-force approach that, in principle, works for any combination of the vanilla and exotic models.

The method we have in mind is quite straightforward, and we describe it with a log-normal LM model representing the exotic model. With f denoting the yield curve, let $\widehat{\Lambda}(f)$ be the ATM Black volatilities of the swaptions in the benchmark set, given the yield curve f . Suppose the delta is calculated by shifting the yield curve from f to f' , which causes a move in swaptions to $\widehat{\Lambda}(f')$, as dictated by our vanilla rule for smile moves. How we proceed depends on the vega calculation method in use. In the direct method, we simply recalibrate the LM model to the new set of swaption volatilities $\widehat{\Lambda}(f')$, and then proceed to use the resulting LM model parameterization together with the shifted yield curve f' to calculate the shocked value of the security in question.

In the indirect vega method, we would do conceptually the same calculation as for the direct method, except we would obtain the LM model parameterization for the shocked yield curve scenario by applying the inverse Jacobian $(\partial\Lambda/\partial G)^{-1}$ to the shifted swaption volatilities $\widehat{\Lambda}(f')$. The inverse Jacobian would automatically be available as part of the basic vega calculation. In the hybrid method, we would first apply the shift in market swaption volatilities arising from the shift in the yield curve, i.e. $\widehat{\Lambda}(f') - \widehat{\Lambda}(f)$, to the base model swaption volatilities, by setting

$$\Lambda' = \Lambda + (\widehat{\Lambda}(f') - \widehat{\Lambda}(f)).$$

Subsequently, we would bootstrap-calibrate the LM model to Λ' and, once again, use this model to calculate the shocked value of the security in question.

26.7 Vega Projections

After our short detour into delta computations, we return to LM vegas. Clearly, the reporting of vega depends on the benchmark set of swaptions used in the vega calculation method: simply put, the vega is reported only to those swaptions that are in the benchmark set. Note that, while we so far assumed that this benchmark set is the same as used for calibrating the LM model in the first place, it actually need not be. Indeed, it is often a good idea to use different sets for calibration and vega calculation. For calibration we often seek to include as many European swaptions as possible to capture the maximum amount of market information in the model calibration, but for vegas it may be preferable to choose a smaller set of benchmarks. There are several reasons for this, starting with the fact that liquidity in different European swaptions is not the same, and the desk may want to express vegas in only the most liquid swaptions¹⁴. In addition, both the computation time (which is linear in the number of shocks applied) and the numerical properties of all vega calculation methods (properties such as the stiffness of the matrix $\partial\Lambda/\partial G$ in Section 26.3.3 or the shape of the Libor volatility bumps in Section 26.3.4) tend to deteriorate as the number of benchmark swaptions grows. For crisper and quicker vegas, it is therefore often useful to cut down on the number of benchmark swaptions.

To understand the issues that a reduced benchmark set of swaptions may lead to, let us revert to the setup used for numerical results in earlier sections (see Section 26.3.1.2, for instance) and imagine that we use an LM model calibrated to European swaptions with expiries and maturities of 1y, 2y, ..., 19y (the “full swaption set” discussed in Section 26.2), yet the vega is to be calculated with the 4×4 benchmark set of swaptions used in previous numerical results (see e.g. Table 26.3). Then the vega for the 10nc1 Bermudan swaption would be reported in the 5y5y bucket but not, say, in 4y5y buckets (as the 4y5y swaption is not in the benchmark set). This, of course, is slightly misleading — it is not that the Bermudan swaption has zero vega in the 4y5y swaption bucket, but that the choice of our benchmark set effectively aggregates that sensitivity and reports it in the 5y5y and, less pronounced, 1y5y buckets. As a trading desk may want to use a fairly granular grid for keeping track of its vega exposure, we should think about how to rationally “project” our coarser vegas onto a finer grid. The idea of

¹⁴This situation would often also be reflected in the usage of different swaption weights in the calibration norm (14.51), with precision weights on illiquid swaptions set lower than on liquid ones.

just assigning our computed 5y5y LMM vega to the 5y5y bucket of the full grid is clearly suboptimal; instead we should somehow spread some of the vega around to buckets surrounding the 5y5y grid point.

There are various methods for projecting vegas from small to full grids, and they all suffer from a degree of arbitrariness as, ultimately, we are trying to create information where there is none. Perhaps the simplest method here is to interpolate (bi)linearly between the points of the small grid to get the values for all points on the full grid, and then rescale to make sure the total vega (i.e. the sum of all vegas in the grid) is the same for the full and reduced-size grids. A slightly more advanced — but nevertheless still somewhat arbitrary — method utilizes the LM model itself to come up with the interpolation scheme. To elaborate on this, let ν^{ex} be the LM vegas for some exotic derivative on the small grid $N_t \times N_x$ (which we are trying to project on a full grid). Furthermore, let the matrix $\nu^{i,j}$ be the matrix (of size $N_t \times N_x$) of vegas for (i, j) -th swaption in the full swaption set; see the start of Section 26.2 for more detail. Then we find the matrix Υ^{ex} of vegas on the full grid by solving the minimization problem

$$\left\| \nu^{\text{ex}} - \sum_{i,j} (\Upsilon^{\text{ex}})_{i,j} \nu^{i,j} \right\|^2 + \mathcal{I}_{\text{smooth}}(\Upsilon^{\text{ex}}) \rightarrow \min, \quad (26.19)$$

where the norm $\|\cdot\|$ is some suitable matrix norm (such as the Frobenius norm used in Section 3.1.3) and $\mathcal{I}_{\text{smooth}}(\Upsilon^{\text{ex}})$ is a smoothing objective along the lines of the definition (14.51); for instance, for first-order smoothness in expiry and tenor directions we would specify

$$\begin{aligned} \mathcal{I}_{\text{smooth}}(\Upsilon^{\text{ex}}) = & w_{\partial t} \sum_{i,j} \left((\Upsilon^{\text{ex}})_{i,j} - (\Upsilon^{\text{ex}})_{i-1,j} \right)^2 \\ & + w_{\partial x} \sum_{i,j} \left((\Upsilon^{\text{ex}})_{i,j} - (\Upsilon^{\text{ex}})_{i,j-1} \right)^2. \end{aligned}$$

The problem (26.19) is quadratic and easily solved with linear algebra methods.

Without a smoothing term, the problem (26.19) is under-specified as there are more free variables than constraints. The smoothing term is essential to pick a unique solution, yet it may lead to undesirable effects like affecting the locality of the vega. Another issue to keep in mind here is that it is not clear how these smoothing weights should be estimated, yet they would impact strongly the allocation of vega. Ultimately, however, one has to live with such issues since, as pointed out, we are filling information “gaps” using fairly arbitrary rules.

On the positive side, the method of projecting LM vegas on a full grid allows for consistent risk representation across a whole portfolio that a trading desk normally trades, including European swaptions, other vanilla

products, and interest rate exotics. This method also allows benchmark sets to be tailored to the features of each derivative thus, potentially, getting better risk resolution and saving computational time by minimizing the number of shocks applied to each derivative. On the other hand, it obviously also introduces a certain level of arbitrariness into the vega calculation and aggregation, and even a danger that the vega for a particular derivative will be reported in inappropriate buckets. To guard against errors, it is often advisable to also calculate LM vegas for all products on the same — and relatively large — set of benchmark swaptions. This could be done relatively infrequently, for instance as part of weekly or monthly control calculations, while leaving daily vegas to be calculated with smaller, product-specific benchmark sets.

Besides the problem of projecting benchmark vegas “up” to a large common grid, we could also contemplate the possibility of projecting vegas “down”, to a smaller grid of potentially different benchmark swaptions. While this capability may seem rather esoteric, some traders find it useful to be able to express their vegas in terms of different sets of European swaptions. It could also be useful for other functions within a bank, such as the risk management department who may use a volatility grid of different shape for calculating risk numbers such as the VaR (see Section 22.3). As the “down” projection compresses information rather than creates it, it is easy to imagine a reasonable algorithm — one just needs to decide how to aggregate “old” buckets into “new” ones. This could be done by, for example, adding up all vegas in the old buckets that are within a certain distance (in expiry/tenor space) from a given new bucket.

26.8 Some Notes on Computing Model Vegas

In all of the vega calculation methods covered in this section, at some point we still need to compute a sensitivity, most typically by Monte Carlo. The standard advice from chapters in Part V of this book for calculating these risk sensitivities apply; however, let us emphasize a few salient points.

- Just like for deltas, the main source of noise for model vegas of *callable* securities is the jumps in the exercise indicators, so the exercise boundary should be kept constant when calculating model vegas, see Section 24.1.1.2.
- More generally, pathwise differentiation of Section 24.3 could be used for model vegas. SDEs for model vegas can be derived by differentiating the SDEs for primary Libor rates with respect to volatility.
- The likelihood ratio or hybrid methods that we just touched upon in Section 24.4 actually work rather better for vegas than for deltas, and could be a viable alternative. Intuitively, a shock to the initial value of a forward Libor rate to compute a delta affects a Monte Carlo path

only up until the first event time (such as an option exercise or a barrier check), whereas the bump to a vega affects the whole path. As the time to the first event goes to zero, the likelihood weight for the delta then explodes, but the one for the vega does not.

- Smoothing of the payoffs by tube Monte Carlo (see Section 23.4) or by importance sampling (see Section 25.2) benefits vegas as well as deltas.
- The variance reduction method in Section 25.3 based on a Markovian approximation could be applied to vegas, but a direct linkage between the original volatility structure and the Markovian approximation is needed, so we should use (25.51) instead of (25.50).

Markovian Projection

Markovian projection is a powerful method for simplifying complex process dynamics to a form that enables rapid calibration of model parameters to quoted option prices. We use the method several times in this book, see for instance Chapters 13 and 15. The usefulness of Markovian projection, however, extends beyond interest rate modeling applications. In this appendix, we develop the relevant theory behind the method and present additional examples and applications.

A.1 Marginal Distributions of Ito Processes

Models used in quantitative finance generally serve to define dynamics of market observables. Some models impose such dynamics directly on the observables; this is, for instance, the case for vanilla models (see Chapters 7, 8, 16 and 17) where the evolution of swap rates is modeled explicitly. Outside of interest rate modeling, equity and FX models typically fall in this category as well. In other cases, the dynamics of market observables are specified indirectly, through modeling of abstract Markovian state variables that drive the market observables through functional relationships (i.e., reconstitution formulas). This style of modeling is common in term structure models for commodities and interest rates (see Chapter 13 for a typical model).

Regardless of type, all models ultimately need to be calibrated to liquidly traded options, most often European call/put options on market observables. To facilitate efficient model calibration, it is generally helpful if exact or approximate analytical expressions exist for European options. Often, the derivation of such results is significantly aided by an initial simplification of the underlying dynamic processes, either because these processes are outright too complex to handle analytically, or because non-linear reconstitution formulas translate simple state variable dynamics into intractable dynamics for market observables. The viability of such simplifications stems from

the simple structure of European options, which only depend on the one-dimensional marginal distributions of the market observable process. As it turns out, irrespective of how complicated a process for a particular market observable is, it is often possible to find much simpler process that preserves the marginal distributions.

A systematic way of finding process simplifications is based on the following fundamental result, see Theorem 4.6 in Gyöngy [1986].

Theorem A.1.1 (Gyöngy). *Let $X(t)$ be given by an SDE*

$$dX(t) = \lambda(t) dW(t), \quad (\text{A.1})$$

where $W(t)$ is a one-dimensional Brownian motion under some probability measure P . Assume that the process $\lambda(t)$ is adapted, bounded, and uniformly bounded away from 0, such that (A.1) admits a unique strong solution. Define $b(t, x)$ by

$$b(t, x)^2 = \mathbb{E}(\lambda(t)^2 | X(t) = x), \quad (\text{A.2})$$

where \mathbb{E} is the expected value operator for the pricing measure P . Then the SDE

$$dY(t) = b(t, Y(t)) dW(t), \quad Y(0) = X(0), \quad (\text{A.3})$$

admits a weak solution Y that has the same one-dimensional distributions as X .

Remark A.1.2. The original result by Gyöngy also includes a drift in the dynamics of X , considering

$$dX(t) = \mu(t) dt + \lambda(t) dW(t)$$

instead of (A.1). The theorem then still holds with (A.3) replaced by

$$dY(t) = a(t, Y(t)) dt + b(t, Y(t)) dW(t), \quad Y(0) = X(0),$$

where

$$a(t, x) = \mathbb{E}(\mu(t) | X(t) = x).$$

In financial applications we nearly always have $\mu(t) = 0$ as we tend to consider the dynamics of X in its own martingale measure. For this reason we do not consider drifts in what follows.

Proof. The original proof in Gyöngy [1986] is fairly involved. A rigorous proof under much weaker assumptions than we stated (see Proposition A.1.4 below) is given in Brunick [2008] and is also highly technical. We do not reproduce either of these proofs, but instead we present a somewhat informal argument¹ originally due to Dupire, see Dupire [1994], Dupire [1997], who independently discovered essentially the same results as in Gyöngy [1986].

¹A version of which we have already seen in Proposition 7.4.2.

The function $b(t, x)$ is often called the *Dupire local volatility* function for the process X .

Let us denote

$$c(t, K) \triangleq c(0, S(0); t, K) = E \left((X(t) - K)^+ \right)$$

to be the values of European call options on X for expiries t and strikes K . It follows from Proposition 7.4.2 that, if we define Y by

$$dY(t) = b(t, Y(t)) dW(t) \quad (\text{A.4})$$

with

$$b(t, K)^2 = \frac{2 \frac{\partial}{\partial t} c(t, K)}{\frac{\partial^2}{\partial K^2} c(t, K)}, \quad (\text{A.5})$$

then the values of European call options in the model (A.4) will be equal to $c(t, K)$ for all expiries t and strikes K , i.e. will be the same as in the model (A.1). To compute the right-hand side, we first write (the use of Dirac delta functions in the integrands can be justified by Tanaka's formula, see Section 1.9.2 and Karatzas and Shreve [1997])

$$d(X(t) - K)^+ = 1_{\{X(t) > K\}} dX(t) + \frac{1}{2} \delta(X(t) - K) \lambda(t)^2 dt$$

and, since $X(t)$ is a P-martingale,

$$E(X(t) - K)^+ - (X(0) - K)^+ = \frac{1}{2} \int_0^t E(\delta(X(t) - K) \lambda(t)^2) dt.$$

Clearly

$$E(\delta(X(t) - K) \lambda(t)^2) = E(\delta(X(t) - K)) \times E(\lambda(t)^2 | X(t) = K)$$

and

$$E(\delta(X(t) - K)) = \frac{\partial^2}{\partial K^2} E(X(t) - K)^+ = \frac{\partial^2}{\partial K^2} c(t, K).$$

In particular,

$$\begin{aligned} \frac{\partial}{\partial t} c(t, K) &= \frac{\partial}{\partial t} \left(E(X(t) - K)^+ - (X(0) - K)^+ \right) \\ &= \frac{1}{2} \frac{\partial^2}{\partial K^2} c(t, K) \times E(\lambda(t)^2 | X(t) = K). \end{aligned}$$

Substituting this equality into (A.5) we obtain

$$b(t, K)^2 = E(\lambda(t)^2 | X(t) = K),$$

consistent with (A.2). \square

Since X and Y have the same one-dimensional marginal distributions, the prices of European options on X and Y for all strikes K and expiries T will be identical (a result that is also implicit in our proof of the theorem). Thus, for the purposes of European option valuation, a potentially complicated process X can be replaced with a simpler Markov process Y ; we call Y the *Markovian projection* of X . Notice that the process Y conveniently is of the local volatility type considered in Chapter 7, for which we have developed many exact or approximate methods for valuation of European options.

Theorem A.1.1 can be extended in a number of ways. Possibly the simplest extension involves relaxing the assumption that the Brownian motion $W(t)$ in (A.1) is one-dimensional, in which case the following trivial corollary holds.

Corollary A.1.3. *Suppose that X follows multi-factor dynamics*

$$dX(t) = \lambda(t)^\top dW(t) \quad (\text{A.6})$$

with $W(t)$ a d -dimensional Brownian motion and $\lambda(t)$ a d -dimensional adapted process whose norm is bounded and uniformly bounded away from 0. Define the SDE

$$dY(t) = b(t, Y(t)) d\widetilde{W}(t), \quad Y(0) = X(0), \quad (\text{A.7})$$

where $\widetilde{W}(t)$ a one-dimensional Brownian motion, and

$$b(t, x)^2 = \mathbb{E}(\lambda(t)^\top \lambda(t) | X(t) = x).$$

Then (A.7) admits a weak solution $Y(t)$ that has the same one-dimensional distributions as $X(t)$.

Proof. Clearly X can be written in one-dimensional form

$$dX(t) = (\lambda(t)^\top \lambda(t))^{1/2} d\widetilde{W}(t),$$

where

$$d\widetilde{W}(t) = \frac{1}{(\lambda(t)^\top \lambda(t))^{1/2}} \lambda(t)^\top dW(t).$$

Simple quadratic variance calculations show that $\widetilde{W}(t)$ is a one-dimensional Brownian motion, and the corollary then is a direct consequence of Theorem A.1.1. \square

The original result by Gyöngy required that the variance process $\lambda(t)^\top \lambda(t)$ in (A.6) be both bounded and uniformly bounded away from 0. These are rather severe limitations that are violated in some standard models of mathematical finance, including the Heston model of Chapter 8. Brunick [2008] has proved the same result under much milder regulatory conditions and also extended it to the case where the asset process X itself is multi-dimensional.

Proposition A.1.4. Let $X(t)$ be a p -dimensional stochastic process given by the strong solution of the SDE

$$dX(t) = \lambda(t)^\top dW(t),$$

where now $\lambda(t)$ is a $(d \times p)$ -matrix-valued adapted process and $W(t)$ is a d -dimensional Brownian motion. We assume that

$$\mathbb{E} \left(\int_0^t \|\lambda(s)^\top \lambda(s)\| ds \right) < \infty$$

for all $t \geq 0$. Then there exists a $(d \times p)$ -matrix-valued function b such that

$$b(t, x)^\top b(t, x) = \mathbb{E} (\lambda(t)^\top \lambda(t) | X(t) = x)$$

for any $t \geq 0$, the SDE

$$dY(t) = b(t, Y(t))^\top dW(t), \quad Y(0) = X(0),$$

admits a weak solution Y , and the random vector $Y(t)$ has the same distribution as the random vector $X(t)$ for any $t \geq 0$.

More generally, Brunick [2008] also proves that we can construct a “mimicking” process Y such that marginal distributions of some *functional* of X , rather than of X itself, are matched by a suitable functional of Y . The definition of functionals for which this result works is fairly technical, but the allowed set includes such financially relevant cases as functions of the running average of X , functions of the running maximum (or minimum) of X , and many others. To avoid technicalities, let us consider only the case of a running maximum of a one-dimensional asset process. Specifically, let us define $M(t, \xi)$ to be the running maximum for a given process ξ , i.e.

$$M(t, \xi) = \sup_{0 \leq s \leq t} \xi(s), \quad t \geq 0.$$

Proposition A.1.5. Let X follow a one-dimensional diffusion

$$dX(t) = \lambda(t) dW(t),$$

where $W(t)$ is a one-dimensional Brownian motion, and $\lambda(t)$ is a scalar adapted process that satisfies

$$\mathbb{E} \left(\int_0^t \lambda(s)^2 ds \right) < \infty, \quad t \geq 0.$$

Then i) there exists a deterministic function $b(t, x, m)$ such that

$$b(t, x, m)^2 = \mathbb{E} (\lambda(t)^2 | X(t) = x, M(t, X) = m)$$

holds for all $t \geq 0$; ii) the SDE

$$dY(t) = b(t, Y(t), M(t, Y)) dW(t), \quad Y(0) = X(0),$$

admits a weak solution Y ; and iii) the pair $(Y(t), M(t, Y))$ has the same distribution as the pair $(X(t), M(t, X))$ for any $t \geq 0$.

Example A.1.6. A European up-and-out barrier call option (see Section 2.1) with expiry T , strike K and barrier level B is an option with the payoff

$$1_{\{M(T,X) < B\}}(X(T) - K)^+.$$

The values of such options on X for all expiries, strikes and barrier levels match those on the mimicking process Y defined by Proposition A.1.5.

While the projection defined by Proposition A.1.5 matches all (up-and-out) barrier option prices, the standard projection result in Theorem A.1.1 does not, as it does not preserve joint distributions of the process observed at multiple times. This is an important point that should be kept in mind in a calibration setting: the Markovian projection of Theorem A.1.1 should be used solely as the means to calibrate to *European* option prices and not to more complicated derivatives².

A.2 Approximations for Conditional Expected Values

According to Theorem A.1.1, the coefficients for the SDE of the Markovian projection are obtained by calculating conditional expected values as in (A.2). This, in the majority of interesting cases, is a non-trivial task. Below, we consider several possible approximations.

A.2.1 Gaussian Approximation

Of the few probability distributions that allow us to calculate conditional expected values in closed form, the most important is, of course, the Gaussian distribution (a fact we use extensively in many places of the book, see for instance Chapter 17). Not surprisingly, we can get good mileage out of the idea of approximating the original distributions of $X(t)$ and $\lambda(t)$ ² with Gaussian distributions, in order to calculate the conditional expected value in (A.6). Many variations are possible here; we present one approach in Proposition A.2.1 below. To fix our setup, we assume that X follows the SDE (A.1) with a process $\lambda(t)$ given by the SDE

$$d\lambda(t)^2 = \nu(t) dt + \varepsilon(t) dZ(t),$$

with two adapted stochastic processes $\nu(t)$ and $\varepsilon(t)$, and a Brownian motion $Z(t)$.

Proposition A.2.1. *The conditional expected value in (A.2) can be approximated by*

²Although some creative approximations for such securities can occasionally be derived from Markovian projection, see Section 13.1.9.4.

$$b(t, x)^2 \approx \overline{\lambda(t)^2} + s(t)(x - X(0)), \quad (\text{A.8})$$

$$s(t) = \frac{\int_0^t \bar{\varepsilon}(s) \bar{\lambda}(s) \bar{\rho}(s) ds}{\int_0^t \overline{\lambda(s)^2} ds},$$

where $\bar{\varepsilon}(t)$, $\overline{\lambda(t)^2}$, $\bar{\rho}(t)$, $\bar{\lambda}(t)$ are deterministic approximations to $\varepsilon(t)$, $\lambda(t)^2$, $\rho(t) \triangleq \langle dW(t), dZ(t) \rangle / dt$, and $\lambda(t)$, respectively. In particular, we can take

$$\bar{\varepsilon}(t) = \mathbb{E}(\varepsilon(t)), \quad \overline{\lambda(t)^2} = \mathbb{E}(\lambda(t)^2) = \mathbb{E}\left(\int_0^t \nu(s) ds\right),$$

$$\bar{\rho}(t) = \mathbb{E}(\langle dW(t), dB(t) \rangle / dt), \quad \bar{\lambda}(t) = \sqrt{\overline{\lambda(t)^2}}.$$

Proof. First, we approximate the dynamics of $(X(t), \lambda(t)^2)$ with the Gaussian processes

$$dX(t) \approx \bar{\lambda}(t) dW(t),$$

$$d\lambda(t)^2 \approx \bar{\nu}(t) dt + \bar{\varepsilon}(t) d\bar{Z}(t),$$

where $\bar{\nu}(t) = \mathbb{E}(\nu(t))$ and $\bar{Z}(t)$ is a Brownian motion such that $\langle dW(t), d\bar{Z}(t) \rangle = \bar{\rho}(t) dt$. The result then follows from the standard conditioning formula for Gaussian variables U, V :

$$\mathbb{E}(U|V) = \mathbb{E}(U) + \frac{\text{Cov}(U, V)}{\text{Var}(V)} (V - \mathbb{E}(V)). \quad (\text{A.9})$$

□

Rather than approximating $X(t)$ and $\lambda(t)^2$ directly by Gaussian processes, we can instead use deterministic functions of Gaussian processes. For instance, if we use exponential “mapping functions”, we would then arrive at a log-normal (rather than Gaussian) approximation. Furthermore, instead of approximating the drift of $\lambda(t)^2$ as a deterministic function, we could instead approximate it by a linear function of $\lambda(t)^2$ itself, which would retain a Gaussian distribution for $\lambda(t)^2$. Ultimately, the original form of the SDEs for the asset and variance in a given model would typically suggest the most proper usage of the Gaussian approximation principle.

As evident from (A.8), the local variance function that emerges from the Gaussian approximation method is linear, whereby the approximating model (A.3) will always generate a monotonic *implied* volatility smile. On the other hand, we may know *a-priori* that the true volatility smile of the original model (A.1) is not close to linear — it could be U-shaped, say. In such cases, a wholesale replacement of the original dynamics by the approximating linear local variance model is unlikely to be satisfactory. As it turns out, it is possible to apply the Gaussian approximation in a more sophisticated manner, leading to a better approximating model. We discuss this in Section A.3.1 below, but first we outline an alternative approach to estimating conditional expected values.

A.2.2 Least-Squares Projection

Section 16.6.2 develops the least-squares projection method for conditional expected values, using the insight that a conditional expected value can be defined as a projection onto a suitable functional space. If we project onto some subspace of the full functional space, we obtain an approximation to the conditional expected value, as stated formally in Proposition 16.6.2. For our purposes here, we focus only on the particularly tractable case of linear subspaces, utilized in Section 16.6.4 to produce a linear least-squares projection. Restating the linear projection in the notations of this appendix, we obtain the following result.

Proposition A.2.2. *The linear least-squares approximation to the conditional expected value in (A.2) is given by*

$$b(t, x)^2 \approx \overline{\lambda(t)^2} + s(t)(x - X(0)), \quad (\text{A.10})$$

where

$$s(t) = \frac{\mathbb{E}(\lambda(t)^2(X(t) - X(0)))}{\mathbb{E}((X(t) - X(0))^2)}, \quad \overline{\lambda(t)^2} = \mathbb{E}(\lambda(t)^2) \quad (\text{A.11})$$

or, more compactly,

$$s(t) = \frac{\text{Cov}(\lambda(t)^2, X(t))}{\text{Var}(X(t))}.$$

We notice a strong similarity between the expressions for the Dupire local volatility approximations in Propositions A.2.1 and A.2.2. By design (as we projected the variance only on linear functions of $X(t)$), the local variance function approximation in (A.10) is still linear in x , and the expression (A.11) reduces to the formulas of Proposition A.2.1 if we apply a Gaussian approximation. This is not surprising, as the linear least-squares projection is known to produce an exact result for conditional expectations of Gaussian variables.

As Propositions A.2.1 and A.2.2 approximate local *variance* with a linear function of spot, both suggest that the SDE for the approximating process is of the displaced square-root type:

$$dY(t) = \left(\overline{\lambda(t)^2} + s(t)(Y(t) - Y(0)) \right)^{1/2} dW(t).$$

While such processes are analytically tractable (see Sections 7.2.4 and 10.2), it is often more convenient to work with processes of the displaced log-normal type (see Proposition 7.2.12) where the local *volatility* function is linear:

$$dY(t) = \sigma(t)(1 + b(t)(Y(t) - Y(0))) dW(t). \quad (\text{A.12})$$

To obtain an approximating process of this type, we can expand the square root function to the first order around $x = Y(0)$, yielding the following result.

Proposition A.2.3. *The displaced log-normal approximation to the process $X(t)$ in (A.1) is given by (A.12), where*

$$\sigma(t) = \sqrt{\lambda(t)^2} = \sqrt{E(\lambda(t)^2)},$$

$$b(t) = \frac{s(t)}{2\sigma(t)^2} = \frac{\text{Cov}(\lambda(t)^2, X(t))}{2E(\lambda(t)^2) \text{Var}(X(t))}.$$

The same result was obtained in Antonov and Misirpashaev [2009a] by a direct application of the least-squares method, i.e. by solving the minimization problem

$$E\left(\left(\lambda(t)^2 - \sigma(t)^2(1 + b(t)(X(t) - X(0)))^2\right)^2\right) \rightarrow \min$$

in $\sigma(t)$ and $b(t)$ and keeping only the leading $O(\lambda(t)^2)$ terms.

A.3 Applications to Local Stochastic Volatility Models

A.3.1 Markovian Projection onto a Stochastic Volatility Model

In applying the Markovian projection method, we are limited by the accuracy of approximations to conditional expected values. As reviewed in Section A.2, the methods that are generally available approximate local volatility or variance functions with linear³ functions which, as discussed in Section A.2.1, is insufficient for the case where $X(t)$ has complex dynamics. Fortunately, Theorem A.1.1 provides us with means to approximate a given model by a model of essentially *any* type, not just local volatility. The following, borderline trivial, corollary to Theorem A.1.1 is the key⁴.

Corollary A.3.1. *If two processes X_1 and X_2 have the same Markovian projections (i.e., they imply identical Dupire local volatility functions), then European put/call option prices on X_1 and X_2 both are identical for all strikes and expiries.*

Let us demonstrate the usefulness of Corollary A.3.1 by applying it to a stochastic volatility model. Let $X_1(t)$ follow a stochastic volatility SDE (recall from Section 8.1 that for non-linear functions $b(t, x)$ such models are sometimes called local stochastic volatility, or LSV, models)

$$dX_1(t) = b_1(t, X_1(t)) \sqrt{z_1(t)} dW(t),$$

³Markovian projection on processes with *quadratic* local volatility (see Section 7.3) was developed in Antonov and Misirpashaev [2009b], using Wiener chaos expansion techniques. Predictably it outperforms projections on linear local volatility processes but is significantly more complicated.

⁴Dupire [1997] dubs this corollary the *universal law of volatility*.

where $z_1(t)$ is some variance process. Suppose we would like to derive approximations for European options on X_1 . One possibility is to approximate X_1 with a local volatility model, using Theorem A.1.1 directly. As we discussed, this is unlikely to work well, at least if we compute conditional expected values with the approximations developed in Section A.2. Instead, we can use Corollary A.3.1 and approximate X_1 with a *stochastic* volatility process that employs a more tractable process for stochastic variance. Let us call this variance process z_2 , and consider a model of the form

$$dX_2(t) = b_2(t, X_2(t)) \sqrt{z_2(t)} dW(t). \quad (\text{A.13})$$

Then Corollary A.3.1 and Theorem A.1.1 imply that to match European option prices in the two models for all strikes and expiries, we need to set $b_2(t, x)$ such that

$$b_2(t, x)^2 = b_1(t, x)^2 \frac{\mathbb{E}(z_1(t) | X_1(t) = x)}{\mathbb{E}(z_2(t) | X_2(t) = x)}. \quad (\text{A.14})$$

While we still need to apply formulas from Section A.2 to approximate conditional expected values in (A.14), the fact that we calculate the *ratio* of two expected values gives us some hope for error cancellation — i.e. even if each individual approximation is not particularly accurate, they are inaccurate “in the same way” and the overall error diminishes when the ratio is formed. To maximize the error cancellation effect, it is obviously beneficial to choose z_2 as similar to z_1 as possible, while still retaining analytical tractability.

Using the SV model (A.13) as the target model for Markovian projection exercise will benefit from the fact that the model (A.13) is quite rich, even for linear local volatility functions $b_2(t, x)$. If z_2 is a square-root process and $b_2(t, x)$ is linear (in x), then the resulting model reduces to the displaced Heston model (8.3)–(8.4) which, as we saw in Chapter 8, is both tractable and capable of generating a wide variety of implied volatility smiles.

In general, limitations of available conditional expected value approximations impose certain restrictions on designing approximating models. In particular, as we want the “output” local volatility $b_2(t, x)$ to be as close to linear as possible — so that the inevitable linear approximation is not far off — we should choose the stochastic variance process z_2 in such a way that the characteristics of this process, and not the shape of the local volatility, explain as much of the *curvature* of the implied volatility smile of the model for X_1 as possible (note that this also holds true for X_1 processes of the non-SV type).

Let us turn our attention to another common application of Corollary A.3.1. Suppose that X_1 follows the SDE

$$dX_1(t) = \lambda(t) \sqrt{z(t)} dW(t), \quad (\text{A.15})$$

where $z(t)$ is a stochastic variance process and $\lambda(t)$ is now a stochastic process in its own right. For instance, $\lambda(t)$ could be a complicated function

of state variables in a term structure model of interest rates, which is a relevant example when $S(t)$ represents a swap rate. We would like to replace the SDE (A.15) with a local stochastic volatility model,

$$dX_2(t) = b(t, X_2(t)) \sqrt{z(t)} dW(t),$$

where we use *the same* stochastic variance process $z(t)$ as in (A.15). Then, according to Corollary A.3.1 we need to set

$$b(t, x)^2 = \frac{\mathbb{E}(\lambda(t)^2 z(t) | X_1(t) = x)}{\mathbb{E}(z(t) | X_2(t) = x)}. \quad (\text{A.16})$$

This formula can be simplified when $\lambda(t)$ and $z(t)$ are (approximately) conditionally independent given $X_1(t)$, in which case we get

$$b(t, x)^2 \approx \mathbb{E}(\lambda(t)^2 | X_1(t) = x) \frac{\mathbb{E}(z(t) | X_1(t) = x)}{\mathbb{E}(z(t) | X_2(t) = x)}. \quad (\text{A.17})$$

In many situations, it can be safely assumed that

$$\mathbb{E}(z(t) | X_1(t) = x) \approx \mathbb{E}(z(t) | X_2(t) = x),$$

in which case the formula simplifies further,

$$b(t, x)^2 \approx \mathbb{E}(\lambda(t)^2 | X_1(t) = x). \quad (\text{A.18})$$

This formula forms the basis for European swaption approximations in term structure models with stochastic volatility; we use it for both quasi-Gaussian models (see Section 13.3.3) and Libor market models (see Proposition 15.2.1).

A.3.2 Fitting the Market with a Local Stochastic Volatility Model

While the focus of this appendix is on approximating more complicated models with simpler ones, direct calibration of local stochastic volatility (LSV) models to the market is another possible application of the techniques we consider.

Let $S(t)$ be the value of a given market variable (for example, a swap rate or an equity price), with initial value $S(0) = S_0$. Suppose market prices of European call or put options are known for all expiries T and strikes K . These can be easily converted (see (A.5)) into a market-implied Dupire local volatility $b_{\text{mkt}}(t, x)$, such that the market European option prices are reproduced by the model

$$dY(t) = b_{\text{mkt}}(t, Y(t)) dW(t), \quad Y(0) = S_0.$$

Suppose we postulate a stochastic variance process $z(t)$ (such as the square-root process of the Heston model, or whatever multi-factor variance process

is in favor in equity modeling circles at the moment), and aim to construct a stochastic volatility model

$$dS(t) = b(t, S(t)) \sqrt{z(t)} dW(t), \quad S(0) = S_0, \quad (\text{A.19})$$

consistent with market European option prices. As follows from Theorem A.1.1, the local volatility function $b(t, x)$ is then given by

$$b(t, x)^2 = \frac{b_{\text{mkt}}(t, x)^2}{\mathbb{E}(z(t)|S(t) = x)}. \quad (\text{A.20})$$

As mentioned before, the challenge of computing $\mathbb{E}(z(t)|S(t) = x)$ makes the method difficult to implement in practice. One choice is to apply finite difference methods to compute the conditional expected value numerically in a forward Kolmogorov PDE for $(S(t), z(t))$, see Ren et al. [2007]. Here, we instead look for analytic approximations.

Linear projections of the type explored in Section A.2 are possible, but are likely to be inaccurate in this case if $b_{\text{mkt}}(t, x)$ has a high degree of convexity. We instead wish to explore methods based on comparing (A.19) to another stochastic volatility model in which European option prices can be cheaply computed. Suppose we have identified such a “proxy” model, defined through a known local volatility function $\tilde{b}(t, x)$,

$$dX(t) = \tilde{b}(t, X(t)) \sqrt{z(t)} dW(t), \quad X(0) = S_0, \quad (\text{A.21})$$

where $z(t)$ is the same process as in (A.19). For tractability, we often choose $\tilde{b}(t, x)$ to be a linear function of x . We will have more to say about the parameterization of the proxy model in Section A.3.3, but for now it suffices to assume that this model allows us to quickly and efficiently calculate European call (or put) option prices. These prices can be turned into a “proxy” Dupire local volatility function $b_{\text{proxy}}(t, x)$ by means of (A.5). Rewriting (A.20), we then have

$$\mathbb{E}(z(t)|X(t) = x) = \frac{b_{\text{proxy}}(t, x)^2}{\tilde{b}(t, x)^2}. \quad (\text{A.22})$$

In other words, from a proxy stochastic volatility model which easily computed European option prices, we can efficiently compute the conditional expected values $\mathbb{E}(z(t)|X(t) = x)$. One way to take advantage of this observation is to combine (A.20) and (A.22) as follows.

Proposition A.3.2. *The local volatility function $b(t, x)$ that makes the model (A.19) consistent with the market is given by*

$$b(t, x) = \tilde{b}(t, x) \frac{b_{\text{mkt}}(t, x)}{b_{\text{proxy}}(t, x)} \frac{\mathbb{E}(z(t)|X(t) = x)}{\mathbb{E}(z(t)|S(t) = x)},$$

where $X(t)$ follows the “proxy” model (A.21) with a known local volatility function $\tilde{b}(t, x)$.

The ratio $b_{\text{mkt}}(t, x)/b_{\text{proxy}}(t, x)$ can, as discussed, usually be computed efficiently. Approximating

$$\frac{\mathbb{E}(z(t)|X(t)=x)}{\mathbb{E}(z(t)|S(t)=x)} \approx 1, \quad (\text{A.23})$$

we obtain the following useful corollary.

Corollary A.3.3. *Under the approximation (A.23) we have that*

$$b(t, x) \approx \tilde{b}(t, x) \frac{b_{\text{mkt}}(t, x)}{b_{\text{proxy}}(t, x)}.$$

To obtain a more sophisticated approximation than that of Corollary A.3.3, we can attempt to improve on (A.23) by looking for an (approximate) functional relationship between $X(t)$ and $S(t)$. Denote

$$\begin{aligned} h(t, x) &= \int_{x_0}^x \frac{dy}{b(t, y)}, \quad \tilde{h}(t, x) = \int_{x_0}^x \frac{dy}{\tilde{b}(t, y)}, \\ H(t, x) &= \tilde{h}^{-1}(t, h(t, x)), \end{aligned} \quad (\text{A.24})$$

where $\tilde{h}^{-1}(t, x)$ is the inverse of $\tilde{h}(t, x)$ in the second (i.e., x) argument. Furthermore, denote

$$\tilde{X}(t) = H(t, S(t)).$$

Then

$$\begin{aligned} d\tilde{X}(t) &= \frac{\partial}{\partial x} H(t, x) \Big|_{x=S(t)} dS(t) + O(dt) \\ &= \tilde{b}\left(t, \tilde{X}(t)\right) \sqrt{z(t)} dW(t) + O(dt). \end{aligned}$$

We see that $\tilde{X}(t)$ and $X(t)$ have identical diffusion coefficients, which suggests the approximation

$$X(t) \approx H(t, S(t)). \quad (\text{A.25})$$

This leads to the following result.

Proposition A.3.4. *The local volatility function $b(t, x)$ that makes the model (A.19) consistent with the market is approximately given by*

$$b(t, x) \approx \tilde{b}(t, H(t, x)) \frac{b_{\text{mkt}}(t, x)}{b_{\text{proxy}}(t, H(t, x))}, \quad (\text{A.26})$$

with $H(t, x)$ given by (A.24).

Proof. By (A.25),

$$\begin{aligned}\mathbb{E}(z(t)|S(t)=x) &= \mathbb{E}(z(t)|H(t,S(t))=H(t,x)) \\ &\approx \mathbb{E}(z(t)|X(t)=H(t,x)).\end{aligned}$$

From (A.22),

$$\mathbb{E}(z(t)|X(t)=H(t,x)) = \frac{b_{\text{proxy}}(t,H(t,x))^2}{\tilde{b}(t,H(t,x))^2},$$

and the result follows from (A.20). \square

We emphasize that $H(t,x)$ depends on the (unknown) function $b(t,x)$, hence (A.26) is, in fact, an equation for $b(t,x)$. This equation can be solved. Let us first denote

$$h_{\text{proxy}}(t,x) = \int_{x_0}^x \frac{dy}{b_{\text{proxy}}(t,y)}, \quad h_{\text{mkt}}(t,x) = \int_{x_0}^x \frac{dy}{h_{\text{mkt}}(t,y)}. \quad (\text{A.27})$$

The following then holds (see Henry-Labordé [2009]).

Proposition A.3.5. *The mapping function $H(t,x)$ is given by*

$$H(t,x) = h_{\text{proxy}}^{-1}(t, h_{\text{mkt}}(t,x)), \quad (\text{A.28})$$

where $h_{\text{proxy}}^{-1}(t,x)$ is the inverse of $h_{\text{proxy}}(t,x)$ defined by (A.27) in the second (x) argument. Furthermore,

$$b(t,x) \approx \tilde{b}(t, h_{\text{proxy}}^{-1}(t, h_{\text{mkt}}(t,x))) \frac{b_{\text{mkt}}(t,x)}{b_{\text{proxy}}(t, h_{\text{proxy}}^{-1}(t, h_{\text{mkt}}(t,x)))}. \quad (\text{A.29})$$

Proof. Differentiating $H(t,x)$ in (A.24) with respect to x we obtain

$$\frac{\partial H(t,x)}{\partial x} = \left. \frac{\partial h(t,x)}{\partial x} \right/ \left. \frac{\partial \tilde{h}(t,f)}{\partial f} \right|_{f=H(t,x)} = \frac{\tilde{b}(t,H(t,x))}{b(t,x)}.$$

Therefore, we can rewrite (A.26) as an (approximate) equation

$$\frac{\partial H(t,x)}{\partial x} = \frac{b_{\text{proxy}}(t,H(t,x))}{b_{\text{mkt}}(t,x)}. \quad (\text{A.30})$$

Treating this as an ODE in x for fixed t , we solve it to find

$$\int_{x_0}^{H(t,x)} \frac{dy}{b_{\text{proxy}}(t,y)} = \int_{x_0}^x \frac{dy}{b_{\text{mkt}}(t,y)}, \quad (\text{A.31})$$

resulting in (A.28). Then (A.29) follows from (A.26) and (A.28). \square

Remark A.3.6. Henry-Labordé [2009] notices that in fact (A.30) implies a condition more general than (A.31), namely that

$$\int_{H_0}^{H(t,x)} \frac{dy}{b_{\text{proxy}}(t,y)} = \int_{x_0}^x \frac{dy}{b_{\text{mkt}}(t,y)}$$

for any H_0 . He proposes to choose H_0 so that the difference in drifts of $S(t)$ and $X(t)$ is minimized (our approximation above matches diffusion terms only).

A.3.3 On Calculating Proxy Local Volatility

The techniques in Section A.3.2 above all hinge on the critical assumption that we can pick a proxy stochastic volatility model that allows for efficient computation of the Dupire local volatility function $b_{\text{proxy}}(t,x)$. Provided that $z(t)$ follows the standard mean-reverting square-root process (as in the Heston model), then an obvious choice for the proxy model is a displaced Heston model, with

$$\tilde{b}(t,x) = \tilde{b}_1 x + \tilde{b}_2. \quad (\text{A.32})$$

Special cases include the $\tilde{b}(t,x) = \tilde{b}_1 x$ (the original Heston model) and $\tilde{b}(t,x) = \tilde{b}_2$ (the “Gaussian” Heston model), but any choice of constants \tilde{b}_1 and \tilde{b}_2 will allow for quick pricing of European put and call options; see Chapter 8 for details. Using the averaging techniques of Chapter 9, we can, in fact, extend (A.32) to linear local volatility function with time-dependent coefficients,

$$\tilde{b}(t,x) = \tilde{b}_1(t)x + \tilde{b}_2(t). \quad (\text{A.33})$$

The time-dependent coefficients in (A.33) can be chosen to make the proxy model resemble as much as possible the true model for S , thereby improving the quality of various approximations made. Ideally we should use

$$\tilde{b}_2(t) = b(t, S_0), \quad \tilde{b}_1(t) = \frac{\partial}{\partial x} b(t, S_0),$$

which is the first-order approximation to the Dupire local volatility $b(t,x)$ along the forward value of $S(t)$. Of course, the value and derivative of $b(t,x)$ are unknown *a-priori*, but one can easily envision various approximations or, perhaps, an iterative procedure where an approximation for the Dupire local volatility in step n is used to define the proxy local volatility \tilde{b} for step $n+1$.

The choice of a mean-reverting square-root process for $z(t)$ for variance leads to a specification that is quite amenable to the methods of Section A.3.2 and is often sufficient. However, more advanced applications such as those considered in Section 15.7 (and some popular models in equity modeling, see e.g. Bergomi [2009]) involve multi-factor stochastic variance

dynamics and require extra effort as European options are then not always easy to calculate or approximate. As should be clear from Section A.3.2, we actually do not need to be able to calculate European option prices in the proxy model; all we really need is the Dupire local volatility for the proxy model, $b_{\text{proxy}}(t, x)$. This function is defined by (see (A.22))

$$b_{\text{proxy}}(t, x)^2 = \tilde{b}(t, x)^2 \mathbb{E}(z(t) | X(t) = x). \quad (\text{A.34})$$

It turns out that there exists a reasonably efficient algorithm for calculating the right-hand side of this question for a large selection of proxy models. To make this statement precise, let us define the proxy stochastic volatility model by

$$dX(t)/X(t) = \lambda \sqrt{z(t)} \left(\rho dW(t) + (1 - \rho^2)^{1/2} dW_X(t) \right), \quad (\text{A.35})$$

$$dz(t) = \nu(t) dt + \varepsilon_1(t) dW(t) + \varepsilon_2(t) dW_z(t), \quad (\text{A.36})$$

for a deterministic $\lambda > 0$ (in the notation of (A.21) therefore $\tilde{b}(t, x) = \lambda x$). Here $W(t)$, $W_X(t)$ and $W_z(t)$ are independent Brownian motions and $\nu(t)$, $\varepsilon_1(t)$ and $\varepsilon_2(t)$ are sufficiently regular adapted processes. The algorithm is based on the following result (see Romano and Touzi [1997] and Lee [2001]).

Proposition A.3.7. *For the model (A.35)–(A.36), we have*

$$\mathbb{E}(z(t) | X(t) = x) = \frac{\mathbb{E}(z(t)\xi(t, x))}{\mathbb{E}(\xi(t, x))}, \quad \xi(t, x) = \frac{1}{D(t)^{1/2}} \phi \left(\frac{\ln(x) - m(t)}{D(t)^{1/2}} \right), \quad (\text{A.37})$$

where

$$m(t) = \ln(X(0)) + \lambda \rho \int_0^t \sqrt{z(s)} dW(s) - \frac{\lambda^2}{2} \int_0^t z(s) ds,$$

$$D(t) = \lambda^2 (1 - \rho^2) \int_0^t z(s) ds,$$

and $\phi(z) = (2\pi)^{-1/2} e^{-z^2/2}$ is the standard Gaussian PDF.

Proof. Proceeding informally, we observe that

$$\mathbb{E}(z(t) | X(t) = x) = \frac{\mathbb{E}(z(t)\delta(X(t) - x))}{\mathbb{E}(\delta(X(t) - x))}, \quad (\text{A.38})$$

where $\delta(x)$ is the Dirac delta function. Let \mathcal{F}_t^z be the filtration generated by $W(s)$ and $W_z(s)$, $0 \leq s \leq t$. Clearly $z(t)$ is adapted to \mathcal{F}_t^z . We can write

$$X(t) = X(0) \exp \left(\lambda \rho \int_0^t \sqrt{z(s)} dW(s) - \frac{\lambda^2}{2} \int_0^t z(s) ds \right)$$

$$\times \exp \left(\lambda (1 - \rho^2)^{1/2} \int_0^t \sqrt{z(s)} dW_X(s) \right),$$

where the first exponential is adapted to \mathcal{F}_t^z and the second is driven by a Brownian motion W_X that is independent of this filtration. Conditioned on the filtration \mathcal{F}_t^z , $\ln(X(t))$ is Gaussian with known moments,

$$\ln(X(t))| \mathcal{F}_t^z \sim \mathcal{N}(m(t), D(t)). \quad (\text{A.39})$$

Notice that

$$\mathbb{E}(\delta(X(t) - x)) = \mathbb{E}(\mathbb{E}(\delta(X(t) - x)| \mathcal{F}_t^z))$$

where $\mathbb{E}(\delta(X(t) - x)| \mathcal{F}_t^z)$ is the log-normal density for the conditional distribution of $X(t)$ defined by (A.39). We therefore have

$$\mathbb{E}(\mathbb{E}(\delta(X(t) - x)| \mathcal{F}_t^z)) = \mathbb{E}\left(\frac{1}{xD(t)^{1/2}}\phi\left(\frac{\ln(x) - m(t)}{D(t)^{1/2}}\right)\right),$$

and

$$\begin{aligned} \mathbb{E}(z(t)\delta(X(t) - x)) &= \mathbb{E}(z(t)\mathbb{E}(\delta(X(t) - x)| \mathcal{F}_t^z)) \\ &= \mathbb{E}\left(\frac{z(t)}{xD(t)^{1/2}}\phi\left(\frac{\ln(x) - m(t)}{D(t)^{1/2}}\right)\right), \end{aligned}$$

The result of the proposition then follows from (A.38). \square

To compute $b_{\text{proxy}}(t, x)$ in a proxy model of the type (A.35)–(A.36), Henry-Labordére [2009] suggests simulating paths of the variance process $z(t)$ and then tabulating the values of $\mathbb{E}(z(t)|X(t) = x)$ by calculating the formula (A.37) for a selection of values of x and t . To obtain the values of $b_{\text{proxy}}(t, x)$ for all t and x , we would then use (A.34) with some sort of interpolation to fill in values of $\mathbb{E}(z(t)|X(t) = x)$ between the tabulated ones.

A.4 Basket Options in Local Volatility Models

So far, our primary application for Markovian projection has been the pricing of European options on scalar processes. However, Markovian projection can easily be extended to options on multiple underlyings, e.g. basket options. Such options appear naturally in interest rate modeling — for example, we can often think of swap rates as baskets of Libor rates — a representation that we use in Section 15.2 together with a Markovian projection to approximate the swap rate process in a Libor market model. CMS spread options could also be thought of as options on baskets of two assets with weights ± 1 . Furthermore, baskets serve as a good example of the practical usage of the Markovian projection method and the various “tricks” that go along with it.

Let us consider a collection of N assets $\mathbf{S}(t) = (S_1(t), \dots, S_N(t))^T$, each driven by its own local volatility model

$$dS_n(t) = \sigma_n(t)\varphi_n(t, S_n(t)) dW_n(t), \quad n = 1, \dots, N. \quad (\text{A.40})$$

Let us denote by $\underline{\alpha}(t)$ the approximation to $\alpha(t)$ obtained by “freezing” the $S^n(t)$ ’s to their initial values in (A.44),

$$\sigma(t)^2 = \mathbb{E}(\alpha(t)^2), \quad b(t) = \frac{2\sigma(t)^2 \text{Var}(S(t))}{\text{Cov}(\alpha(t)^2, S(t))}. \quad (\text{A.47})$$

with

$$d\underline{S}(t) = \sigma(t) \left(1 + b(t)(\underline{S}(t) - S(0)) \right) dW(t), \quad (\text{A.46})$$

As in Proposition A.2.3, let us approximate the dynamics of S with a displaced log-normal process. We then have $\underline{S} \approx S$, where $\alpha(t)$ is here a complicated function of the vector of asset prices $\mathbf{S}(t)$. Where $W(t)$ is easily seen to be a standard Brownian motion. The process

$$d\underline{S}(t) = \alpha(t) dW(t),$$

then

$$dW(t) \stackrel{!}{=} \sum_{n=1}^N w_n \sigma_n(t) \phi_n(t, S^n(t)) dW^n(t), \quad (\text{A.45})$$

$$\alpha(t)^2 \stackrel{!}{=} \sum_{n,m=1}^N w_n w_m \sigma_n(t) \phi_n(t, S^n(t)) \sigma_m(t) \phi_m(t, S^m(t)) p_{n,m}, \quad (\text{A.44})$$

If we define

$$\sum_{n=1}^N w_n \sigma_n(t) \phi_n(t, S^n(t)) dW^n(t) = (t) dS(t). \quad (\text{A.43})$$

We wish to calculate values of European options on $S(t)$, a problem first considered by Avellaneda et al. [2002] (but using quite different methods). Applying Ito’s lemma to $S(t)$, we see that

$$\sum_{n=1}^N w_n S^n(t) = (t) S. \quad (\text{A.42})$$

We define by $S(t)$ the value of the basket (also sometimes called the *index*)

$$\langle dW^i(t), dW^j(t) \rangle = \rho_{i,j} dt, \quad i, j = 1, \dots, N.$$

with correlations

The Brownian motions $(W_1(t), \dots, W_N(t))$ are assumed to be correlated

$$\phi_n(t, x) = 1 + b_n(t)(x - S^n(0)). \quad (\text{A.41})$$

We assume from now on that each local volatility can be well-approximated by a linear function, i.e. that

$$\bar{\lambda}(t)^2 \triangleq \sum_{n,m=1}^N w_n w_m \sigma_n(t) \sigma_m(t) \rho_{n,m}. \quad (\text{A.48})$$

In the same spirit, we approximate $W(t)$ in (A.45) by $\bar{W}(t)$, where

$$d\bar{W}(t) \triangleq \frac{1}{\bar{\lambda}(t)} \sum_{n=1}^N w_n \sigma_n(t) dW_n(t),$$

is also a Brownian motion. Finally, let us denote

$$\rho_n(t) \triangleq \langle dW_n(t), d\bar{W}(t) \rangle / dt = \frac{1}{\bar{\lambda}(t)} \sum_{m=1}^N w_m \sigma_m(t) \rho_{n,m}. \quad (\text{A.49})$$

With all this notation in place, we then can state the following proposition.

Proposition A.4.1. *Let the skew functions φ_n be as stated in (A.41). A displaced log-normal approximation to the dynamics of the basket $S(t)$ in (A.43) is then given by (A.46) where*

$$\sigma(t)^2 \approx \sum_{n=1}^N \sum_{m=1}^N w_n w_m \sigma_n(t) \sigma_m(t) \rho_{n,m}, \quad (\text{A.50})$$

$$b(t) \approx \frac{\sum_{n=1}^N \bar{\lambda}(t) \sigma_n(t) \rho_n(t) \left(\int_0^t \bar{\lambda}(s) \sigma_n(s) \rho_n(s) ds \right) w_n b_n(t)}{\bar{\lambda}(t)^2 \int_0^t \bar{\lambda}(s)^2 ds}, \quad (\text{A.51})$$

with the ρ_n 's defined in (A.49) and $\bar{\lambda}(t)$ defined in (A.48).

Proof. The approximation (A.50) follows from

$$\sigma(t)^2 = E(\lambda(t)^2) \approx E(\bar{\lambda}(t)^2) = \bar{\lambda}(t)^2$$

and (A.48). To find the covariance term in (A.47), we recall (A.41) and write

$$\begin{aligned} & \text{Cov}(\lambda(t)^2, S(t)) \\ &= \sum_{n,m=1}^N w_n w_m \sigma_n(t) \sigma_m(t) \rho_{n,m} \\ & \quad \times \text{Cov}((1 + b_n(t))(S_n(t) - S_n(0)), (1 + b_m(t))(S_m(t) - S_m(0))), S(t)) \\ & \approx 2\bar{\lambda}(t) \sum_{n=1}^N w_n \sigma_n(t) \rho_n(t) b_n(t) \text{Cov}(S_n(t), S(t)), \end{aligned}$$

where we disregarded the terms $\text{Cov}((S_n(t) - S_n(0))(S_m(t) - S_m(0)), S(t))$ as being of higher order in volatility. Furthermore,

$$\begin{aligned}
& \text{Cov}(S_n(t), S(t)) \\
&= \sum_{m=1}^N w_m \text{Cov}(S_n(t), S_m(t)) \\
&= \sum_{m=1}^N \int_0^t w_m \rho_{n,m} \sigma_n(s) \sigma_m(s) \mathbb{E}(\varphi_n(s, S_n(s)) \varphi_m(s, S_m(s))) ds \\
&\approx \int_0^t \sigma_n(s) \bar{\lambda}(s) \rho_n(s) ds.
\end{aligned} \tag{A.52}$$

To summarize,

$$\text{Cov}(\lambda(t)^2, S(t)) \approx 2\bar{\lambda}(t) \sum_{n=1}^N w_n \sigma_n(t) \rho_n(t) b_n(t) \int_0^t \sigma_n(s) \bar{\lambda}(s) \rho_n(s) ds.$$

In the same spirit, from (A.52),

$$\begin{aligned}
\text{Var}(S(t)) &= \sum_{n=1}^N w_n \text{Var}(S_n(t), S(t)) \\
&\approx \sum_{n=1}^N w_n \int_0^t \sigma_n(s) \bar{\lambda}(s) \rho_n(s) ds = \int_0^t \bar{\lambda}(s)^2 ds,
\end{aligned}$$

and the proposition follows from (A.47). \square

Remark A.4.2. Proposition A.4.1 depends on numerous ad-hoc approximations, the nature of which we did not characterize rigorously. A more detailed analysis can be found in Antonov and Misirpashaev [2009a] where it is shown that our expressions for $\sigma(t)$ and $b(t)$ are leading order terms in the small-volatility limit.

The parameter $b(t)$ in (A.46) represents the slope of the local volatility function for the basket S , and the expression (A.51) relates this slope to a weighted average of the slopes $b_n(t)$ of the individual volatility functions for the basket components. This approximation works best when the skews $b_n(t)$ are of the same sign and the weights w_n are all positive. For skews b_n (or weights w_n) of mixed signs, the approximation is not very accurate, however. This should be intuitively clear, since the difference of two processes with positive skews (say) can easily have a *U*-shaped smile, which is obviously not well-approximated by a projection of the difference onto a displaced log-normal process. To handle this case, one possible solution is to use a projection on a stochastic volatility process (even though the components of the basket are local volatility processes), as described in Section A.5 below. Alternatively, we can use a projection on a local volatility process with a quadratic local volatility as mentioned in footnote 3. Finally, we can always fall back on the copula-based methods of Chapter 17.

A.5 Basket Options in Stochastic Volatility Models

We continue investigating basket options, but now augment the model (A.40), (A.41) with stochastic volatility. Specifically, we replace (A.40) with

$$dS_n(t) = \sigma_n(t)\varphi_n(t, S_n(t)) \sqrt{z_n(t)} dW_n(t), \quad n = 1, \dots, N,$$

where (A.41) still holds, and where individual stochastic variance processes $z_n(t)$ are defined by

$$dz_n(t) = \theta_n(t)(1 - z_n(t)) dt + \eta_n(t)\sqrt{z_n(t)} dW_{N+n}(t), \quad z_n(0) = 1, \quad (\text{A.53})$$

$n = 1, \dots, N$. Here, $(W_1(t), \dots, W_{2N})$ is a $2N$ -dimensional Brownian motion with correlations

$$\langle dW_i, dW_j \rangle = \rho_{i,j} dt, \quad i, j = 1, \dots, 2N.$$

To simplify the already cumbersome notation, let us absorb the basket weights w_n into a redefinition of the asset processes: $S_n \leftarrow w_n S_n$. The basket value $S(t)$ is now

$$S(t) = \sum_{n=1}^N S_n(t),$$

where we now have

$$dS(t) = \sum_{n=1}^N \sigma_n(t)\varphi_n(t, S_n(t)) \sqrt{z_n(t)} dW_n(t) = \lambda(t) dW(t),$$

with

$$\begin{aligned} \lambda(t)^2 &= \sum_{n,m=1}^N \sigma_n(t)\varphi_n(t, S_n(t)) \sigma_m(t)\varphi_m(t, S_m(t)) \rho_{n,m} \sqrt{z_n(t)z_m(t)}, \\ &\quad (\text{A.54}) \end{aligned}$$

$$dW(t) = \frac{1}{\lambda(t)} \sum_{n=1}^N \sigma_n(t)\varphi_n(t, S_n(t)) \sqrt{z_n(t)} dW_n(t).$$

S above is driven by a multi-dimensional stochastic volatility process, and, as we discussed previously in Section A.3.1, projecting S on a displaced log-normal local volatility process is unlikely to lead to accurate option approximations. Following Antonov et al. [2009], we instead investigate projections on a displaced Heston process.

Let us first assume that the skew parameter $b(t)$ of the target approximation is given exogenously (we will discuss its computation later in this section). With $b(t)$ given, we rewrite the dynamics of the process $S(t)$ in a way more suitable for approximations:

$$dS(t) = \sigma(t)(1 + b(t)(S(t) - S(0)))\sqrt{z(t)}dW(t), \quad (\text{A.55})$$

where

$$z(t) \triangleq \frac{\Lambda(t)^2}{\sigma(t)^2}, \quad \Lambda(t)^2 \triangleq \frac{\lambda(t)^2}{(1 + b(t)(S(t) - S(0)))^2}, \quad (\text{A.56})$$

and

$$\sigma(t)^2 \triangleq \mathbb{E}(\Lambda(t)^2). \quad (\text{A.57})$$

For future reference, we apply Ito's lemma to $z(t)$ and write

$$dz(t) = \nu(t)dt + \varepsilon(t)\sqrt{z(t)}dZ(t), \quad (\text{A.58})$$

where $Z(t)$ is a Brownian motion such that $\langle dW(t), dZ(t) \rangle = \chi(t)dt$, where the exact form of stochastic processes $\nu(t)$, $\varepsilon(t)$, $\chi(t)$ is not important for the moment.

By the multi-dimensional extension of Gyöngy's theorem in Proposition A.1.4, we replicate the exact distribution of the pair $(S(t), z(t))$ for each $t \geq 0$ with $(\tilde{S}(t), \tilde{z}(t))$, where

$$\begin{aligned} d\tilde{S}(t) &= \sigma(t)\left(1 + b(t)(\tilde{S}(t) - \tilde{S}(0))\right)\sqrt{\tilde{z}(t)}dW(t), \\ d\tilde{z}(t) &= \mathbb{E}\left(\nu(t) \mid S(t) = \tilde{S}(t), z(t) = \tilde{z}(t)\right)dt \\ &\quad + \left(\mathbb{E}\left(\varepsilon(t)^2 \mid S(t) = \tilde{S}(t), z(t) = \tilde{z}(t)\right)\right)^{1/2}dZ(t). \end{aligned} \quad (\text{A.59})$$

We cannot calculate the conditional expectations in (A.59) exactly, so we proceed to assume a particular parametric form for the process $(\tilde{S}(t), \tilde{z}(t))$,

$$d\tilde{S}(t) \approx \sigma(t)\left(1 + b(t)(\tilde{S}(t) - \tilde{S}(0))\right)\sqrt{\tilde{z}(t)}dW(t), \quad (\text{A.60})$$

$$d\tilde{z}(t) \approx \theta(t)(1 - \tilde{z}(t))dt + \eta(t)\sqrt{\tilde{z}(t)}dZ(t), \quad (\text{A.61})$$

with $\langle dW(t), dZ(t) \rangle = \rho(t)dt$. The unknown parameters $\theta(t)$, $\eta(t)$ and $\rho(t)$ emerge as solutions to the following optimization problems,

$$\begin{aligned} \mathbb{E}\left((\varepsilon(t)^2 - \eta(t)^2z(t))^2\right) &\rightarrow \min, \\ \mathbb{E}\left((\varepsilon(t)\chi(t)z(t) - \eta(t)\rho(t)z(t))^2\right) &\rightarrow \min, \\ \mathbb{E}\left((\nu(t) - \theta(t)(1 - z(t)))^2\right) &\rightarrow \min. \end{aligned} \quad (\text{A.62})$$

From the usual first-order optimality conditions, we then find

$$\begin{aligned} \eta(t)^2 &= \frac{\mathbb{E}(\varepsilon(t)^2z(t))}{\mathbb{E}(z(t)^2)}, \\ \rho(t) &= \frac{\mathbb{E}(\varepsilon(t)\chi(t)z(t)^2)}{\eta(t)\mathbb{E}(z(t)^2)}, \\ \theta(t) &= \frac{\mathbb{E}(\nu(t)(1 - z(t)))}{\mathbb{E}((1 - z(t))^2)}. \end{aligned} \quad (\text{A.63})$$

The expectations required in these equations can be approximated from the definition of $z(t)$ in (A.56) and the coefficients $\nu(t)$, $\varepsilon(t)$, $\chi(t)$ in (A.58). The (laborious) calculations are performed in Antonov et al. [2009], and we do not reproduce them here.

We have not yet specified how to set the skew function $b(t)$ in (A.60), originally appearing in (A.55) and (A.56). One idea here is to use the results from the local volatility approximation in Section A.4 and simply use (A.51). Another alternative suggested in Antonov et al. [2009] finds $b(t)$ by minimizing the *defect* $D_\theta(t)$ for the solution (A.63) of the problem (A.62),

$$D_\theta(t) = \mathbb{E}(\nu(t)^2) - \frac{(\mathbb{E}(\nu(t)(1-z(t))))^2}{\mathbb{E}((1-z(t))^2)}. \quad (\text{A.64})$$

The defect $D_\theta(t)$ measures the error in the objective function in (A.62) for the solution (A.63), and clearly is a function of $b(t)$. By minimizing the defect, an ODE for $b(t)$ arises that can be solved in closed form. Again, we refer the interested reader to Antonov et al. [2009] for details.

Our final topic in this section is the calculation of the effective volatility $\sigma(t)$ in (A.57). Among all parameters of the model (A.60)–(A.61), $\sigma(t)$ typically is the one that has the biggest impact on the quality of approximations. From (A.54) and (A.56), the standard “freezing” (where $S_n(t) = S_n(0)$ for all $n = 1, \dots, N$) approximation gives us

$$\sigma(t)^2 = \sum_{n,m=1}^N \sigma_n(t)\sigma_m(t)\rho_{n,m}\mathbb{E}\left(\sqrt{z_n(t)z_m(t)}\right). \quad (\text{A.65})$$

In principle, the same freezing idea could be applied to z_n ’s, leading us to an approximation that is identical to the one we already derived for the local volatility case, see (A.50). This is, indeed, the leading term of the expansion of $\mathbb{E}(\Lambda(t)^2)$ in small volatilities (see Antonov et al. [2009]), but for typical parameter settings encountered in interest rate modeling the quality of approximations based on this choice for $\sigma(t)$ is rather poor — volatility of variance parameters η_n are often simply too far from being “small” (larger than 1 is typical).

To improve the accuracy of the overall Markovian projection onto a displaced Heston model, we need to find a way to calculate $\sigma(t)$ in (A.65) without the assumption of small variances of z_n ’s. Clearly, our ability to do so hinges on accurate approximations to $\mathbb{E}(\sqrt{z_n(t)z_m(t)})$ in (A.65), which could be interesting for other purposes as well. We discuss such a “non-perturbative” approximation in Appendix A.A (see Proposition A.A.1) where we also consider a related problem of approximating $\mathbb{E}(\sqrt{z_n(t)})$ (Lemma A.A.2).

With these approximations in place, Antonov et al. [2009] demonstrate excellent performance of Markovian projection onto a displaced Heston for basket and spread options.

To conclude, let us quickly summarize the entire algorithm. First, we calculate a non-perturbative approximation to $\sigma(t)$ given by (A.65), with

the square roots obtained by the methods of Appendix A.A. Second, we calculate the optimal skew function $b(t)$ by minimizing the function $D_\theta(t)$ in (A.64). As it turns out, this minimization problem (or, rather, a collection of minimization problems indexed by t) leads to a first-order ODE on $b(t)$ with coefficients that only depend on $\sigma(t)$ (in the small volatility limit), as derived in Antonov et al. [2009]. Finally, having established $\sigma(t)$ and $b(t)$, we can now solve for optimal coefficients for the stochastic variance process (and its correlation with the asset process) by solving (A.62) for each $t > 0$.

A.A Appendix: Approximations for $E(\sqrt{z_n(t)z_m(t)})$ and $E(\sqrt{z_n(t)})$

As discussed in Section A.5, European option approximations in multi-dimensional SV models require calculations of certain expected values of square-root processes, a subject we consider in this section. Let us simplify the notations of (A.53) a little, and consider a two-dimensional square-root process

$$dz_i(t) = \theta_i(1 - z_i(t)) dt + \eta_i \sqrt{z_i(t)} dW_i(t), \quad z_i(0) = 1, \quad (\text{A.66})$$

where $i = 1, 2$, and $\langle dW_1(t), dW_2(t) \rangle = \rho dt$. We first consider the problem of approximating $E(\sqrt{z_1(t)z_2(t)})$.

Proposition A.A.1. *For a two-dimensional square-root process (A.66), let us define*

$$\begin{aligned} E(\theta_1, \theta_2, \eta_1, \eta_2, \rho) &= \left(1 + \rho \eta_1 \eta_2 q_1 q_2 \frac{1 - e^{-(\theta_1 + \theta_2 - \rho \eta_1 \eta_2 q_1 q_2)t}}{\theta_1 + \theta_2 - \rho \eta_1 \eta_2 q_1 q_2} \right) \\ &\quad \times E(\sqrt{z_1(t)}) E(\sqrt{z_2(t)}), \end{aligned} \quad (\text{A.67})$$

where q_i , $i = 1, 2$, are obtained by solving

$$1 = \frac{\left(E(\sqrt{z_i(t)}) \right)^2}{2\theta_i - \eta_i^2 q_i^2} \left(2\theta_i - \eta_i^2 q_i^2 e^{-(2\theta_i - \eta_i^2 q_i^2)t} \right).$$

This function gives an approximation to $E(\sqrt{z_1(t)z_2(t)})$, which is

1. Exact in the limit $\rho = 0$.
2. Exact in the limit $\rho = 1$, $\theta_1 = \theta_2$, $\eta_1 = \eta_2$.
3. Has correct leading behavior in the expansion in powers of η_1 and η_2 .

The proposition is proved in Section A.A.1. We note that enforcement of the first two non-perturbative conditions substantially improves the accuracy of the approximation to $\sigma(t)$ in (A.65) when $\eta_1, \eta_2 \geq 1$.

The approximation (A.67) relies on our ability to calculate $E(\sqrt{z_i(t)})$, the calculation that is also of independent interest sometimes. To state the result, we assume that $z(t)$ follows the square-root process (8.4):

$$dz(t) = \theta(z_0 - z(t)) dt + \eta\sqrt{z(t)} dZ(t), \quad z(0) = z_0, \quad (\text{A.68})$$

and recall the definitions (8.6):

$$d = 4\theta z_0/\eta^2, \quad n(t, T) = \frac{4\theta e^{-\theta(T-t)}}{\eta^2 (1 - e^{-\theta(T-t)})}.$$

The following result, proven in Section A.A.2, derives the required representation.

Lemma A.A.2. *In the model (A.68) we have*

$$E(\sqrt{z(t)}) = \left(\frac{2e^{-\theta t}}{n(0,t)} \right)^{1/2} e^{-n(0,t)/2} \sum_{j=0}^{\infty} \frac{(n(0,t)/2)^j}{j!} \frac{\Gamma(d/2 + j + 1/2)}{\Gamma(d/2 + j)}, \quad (\text{A.69})$$

where $\Gamma(z)$ is the Gamma function.

Remark A.A.3. The series in (A.69) converges rapidly, so only the first few terms need to be computed. In each term, the ratio of Gamma functions can be evaluated by standard algorithms present in most numerical software packages. We find that the following approximation is sufficient for our purposes:

$$\Gamma(x + 1/2)/\Gamma(x) = \begin{cases} x(0.5619x^2 - 1.3353x + 1.6651), & x \in [0, 0.9], \\ \sqrt{x}(1 - \frac{1}{8x} + \frac{1}{128x^2}), & x > 0.9, \end{cases}$$

where the first line is obtained by fitting $\Gamma(x+1/2)/\Gamma(x)$ with a third-degree polynomial over the interval $[0, 1]$, and the second line is the truncation of the series expansion of $\Gamma(x + 1/2)/\Gamma(x)$ valid for $x > 1$, see Cevher et al. [2007]. The cut-off point of 0.9 is chosen to make the function continuous and (nearly) C^1 .

A.A.1 Proof of Proposition A.A.1

A.A.1.1 Step 1. Reduction to Covariance

Taking expected values of both sides of (A.66) we see that $E(z_i(t)) = 1$, $i = 1, 2$. Furthermore, we note that

$$\begin{aligned} d(z_1(t)z_2(t)) &= z_1(t) \left(\theta_2(1 - z_2(t)) dt + \eta_2\sqrt{z_2(t)} dW_2(t) \right) \\ &\quad + z_2(t) \left(\theta_1(1 - z_1(t)) dt + \eta_1\sqrt{z_1(t)} dW_1(t) \right) \\ &\quad + \eta_1\eta_2\rho\sqrt{z_1(t)z_2(t)} dt, \end{aligned}$$

so that

$$\frac{d}{dt} \mathbb{E}(z_1(t)z_2(t)) = (\theta_1 + \theta_2)(1 - \mathbb{E}(z_1(t)z_2(t))) + \eta_1\eta_2\rho \mathbb{E}\left(\sqrt{z_1(t)z_2(t)}\right). \quad (\text{A.70})$$

Hence, calculating $\mathbb{E}(\sqrt{z_1(t)z_2(t)})$ is equivalent to calculating $\mathbb{E}(z_1(t)z_2(t))$.

A.A.1.2 Step 2. Linear Approximation to Volatility

To proceed, we consider a linear approximation to (A.66) obtained by substituting

$$\sqrt{z_i} \rightarrow p_i + q_i z_i, \quad (\text{A.71})$$

with the coefficients to be found. Define by $c(t) = c(t; p_1, q_1, p_2, q_2)$ the value of $\mathbb{E}(\tilde{z}_1(t)\tilde{z}_2(t))$ in the approximate model,

$$d\tilde{z}_i(t) = \theta_i(1 - \tilde{z}_i(t)) dt + \eta_i(p_i + q_i \tilde{z}_i(t)) dW_i(t), \quad \tilde{z}_i(0) = 1. \quad (\text{A.72})$$

Simple calculations yield an ODE on $c(t)$,

$$\begin{aligned} \frac{d}{dt}c(t) &= -(\theta_1 + \theta_2)(c(t) - 1) \\ &\quad + \rho\eta_1\eta_2(p_1 + q_1)(p_2 + q_2) + \rho\eta_1\eta_2q_1q_2(c(t) - 1), \end{aligned}$$

which can be solved,

$$c(t) = 1 + \rho\eta_1\eta_2(p_1 + q_1)(p_2 + q_2) \frac{1 - e^{-(\theta_1 + \theta_2 - \rho\eta_1\eta_2q_1q_2)t}}{\theta_1 + \theta_2 - \rho\eta_1\eta_2q_1q_2}. \quad (\text{A.73})$$

Then, from (A.70) and the equality (A.73) that we use to approximate $\mathbb{E}((z_1(t)z_2(t))$ in the original model (A.66), we get

$$\begin{aligned} \mathbb{E}\left(\sqrt{z_1(t)z_2(t)}\right) &\quad (\text{A.74}) \\ &\approx \frac{(p_1 + q_1)(p_2 + q_2)}{\theta_1 + \theta_2 - \rho\eta_1\eta_2q_1q_2} \\ &\quad \times \left(\theta_1 + \theta_2 - \rho\eta_1\eta_2q_1q_2e^{-(\theta_1 + \theta_2 - \rho\eta_1\eta_2q_1q_2)t}\right) \\ &= (p_1 + q_1)(p_2 + q_2) \left(1 + \rho\eta_1\eta_2q_1q_2 \frac{1 - e^{-(\theta_1 + \theta_2 - \rho\eta_1\eta_2q_1q_2)t}}{\theta_1 + \theta_2 - \rho\eta_1\eta_2q_1q_2}\right). \end{aligned}$$

A.A.1.3 Step 3. Coefficients of the Linear Approximation

Applying (A.74) to $z_2(t) = z_1(t)$ we obtain

$$1 = \mathbb{E}((z_1(t))) = \frac{(p_1 + q_1)^2}{2\theta_1 - \eta_1^2q_1^2} \left(2\theta_1 - \eta_1^2q_1^2e^{-(2\theta_1 - \eta_1^2q_1^2)t}\right), \quad (\text{A.75})$$

which gives one equation on p_1, q_1 . The other one is obtained by using as $z_2(t)$ an independent copy of $z_1(t)$,

$$\left(E\left(\sqrt{z_1(t)}\right)\right)^2 = (p_1 + q_1)^2. \quad (\text{A.76})$$

The coefficients p_1 and q_1 are determined as a solution to the system (A.75), (A.76), provided $E(\sqrt{z_1(t)})$ is known. Similar system holds for p_2, q_2 . This completes the derivation of the approximating function (A.67).

A.A.1.4 Step 4. Order of Approximation

The first two features of the approximation listed in the Proposition A.A.1 are valid by construction. The validity of the last feature is obvious because the usage of a linear approximation (A.71) is consistent with the first order of perturbative expansion in volatilities. This is also easy to verify directly. Indeed, a straightforward perturbative calculation gives⁵

$$E\left(\sqrt{z_i(t)}\right) = 1 - \frac{\eta_i^2 (1 - e^{-2\theta_i t})}{16\theta_i} + O(\eta_i^3),$$

$$\frac{E\left(\sqrt{z_1(t)z_2(t)}\right)}{E\left(\sqrt{z_1(t)}\right) E\left(\sqrt{z_2(t)}\right)} = 1 + \frac{\rho\eta_1\eta_2 (1 - e^{-(\theta_1+\theta_2)t})}{4(\theta_1 + \theta_2)} + O((\max(\eta_1, \eta_2))^3),$$

which is in agreement with the leading order of expansion of (A.67) in η_1 and η_2 . Note that we used the zero order expansion for coefficients $q_i = 1/2 + O(\eta_i)$ to prove the agreement.

A.A.2 Proof of Lemma A.A.2

By Proposition 8.3.2, $z(t)$ is distributed as $e^{-\theta t}/n(0, t)$ times a non-central chi-square distributed random variable ξ with $\nu = d$ degrees of freedom and non-centrality parameter $\gamma = n(0, t)$. Thus we obtain that

$$E\left(\sqrt{z(t)}\right) = \left(\frac{e^{-\theta t}}{n(0, t)}\right)^{1/2} E\left(\sqrt{\xi}\right),$$

where ξ has the density (see (8.5))

$$\frac{\partial}{\partial z} \Upsilon(z; \nu, \gamma) = e^{-\gamma/2} \sum_{j=0}^{\infty} \frac{(\gamma/2)^j}{j! 2^{\nu/2+j} \Gamma(\nu/2 + j)} z^{\nu/2+j-1} e^{-z/2}.$$

⁵Due to the reflecting symmetry of the Brownian motion, an arbitrary average of the variables $z_1(t)$ and $z_2(t)$ is an even function of η_1 and η_2 . Thus, the order of the approximations below is effectively higher.

Then

$$\begin{aligned}
 E(\sqrt{\xi}) &= e^{-\gamma/2} \sum_{j=0}^{\infty} \frac{(\gamma/2)^j}{j! 2^{\nu/2+j} \Gamma(\nu/2 + j)} \int_0^{\infty} z^{1/2+\nu/2+j-1} e^{-z/2} dz \\
 &= e^{-\gamma/2} \sum_{j=0}^{\infty} \frac{(\gamma/2)^j}{j! 2^{\nu/2+j} \Gamma(\nu/2 + j)} 2^{1/2+\nu/2+j} \Gamma(\nu/2 + j + 1/2) \\
 &= \sqrt{2} e^{-\gamma/2} \sum_{j=0}^{\infty} \frac{(\gamma/2)^j}{j!} \frac{\Gamma(\nu/2 + j + 1/2)}{\Gamma(\nu/2 + j)}.
 \end{aligned} \tag{A.77}$$

Index

- absorbing boundary, *see* diffusion,
absorbing barrier
accrual factor, *see* year fraction
ADI, *see* PDE, ADI scheme
adjusters method, *see* out-of-model
adjustment, adjusters method
affine short rate model, 429–442,
510–518
bond reconstitution formula, 431,
513–515
calibration, 439–441
 multi-pass bootstrap, 440
calibration to yield curve, 435–437
characteristic function, 432
European swaption, 437
 Fourier integration, 437
 Gram-Charlier expansion, 437
extended transform, 431
 constant parameters, 432, 434
 piecewise constant parameters, 434
Feller condition, 319, 430
importance sampling, 1065
moment-generating function, 432,
437
Monte Carlo, 442
multi-factor, 510–518
 bond dynamics, 514
 bond reconstitution formula,
 513–515
 existence and uniqueness, 513
 exponential affine, 511
 Feller condition, 513
 forward rate correlation, 514
forward rate dynamics, 514
regularity issues, 512–513
short rate state dynamics, 512
one-factor, 429–442
PDE, 442
regularity issues, 430
short rate domain, 430
short rate dynamics, 429
short rate state dynamics, 435
swap rate volatility, 438
 affine approximation, 438
 time averaging, 438
time-dependent, 431
volatility skew range, 431
volatility smile, 430
almost surely, 4
American capped straddle, 936
American swaption, 893–898
 accrued current coupon, 893
 approximating with Bermudan
 swaption, *see* Bermudan swap-
 tion, approximating American
 swaption
 discontinuity of exercise value in
 time, 893
PDE, 895–897
 extra state variable, 896–897
 proxy Libor rate method, 895–896
American/Bermudan option, 30–42
 Bellman principle, 32, 33, 69
 Black-Scholes model, 837
 capped, 936
 conditional on no exercise, 31

- continuation region, 33
- discontinuity at expiry, 39
- duality, 36
- early exercise boundary, 37
- early exercise premium, 36, 39, 42
- exercise never optimal, 36
- exercise policy, 30
- exercise region, 33
- exercise value, 30
- high contact condition, 38
- hold value, 32
- integral representation, 39, 41
- lower bound, *see* Monte Carlo, lower bound for American option
- marginal exercise value decomposition, 41
- Monte Carlo, 158–165
 - confidence interval for value, 164
 - random tree, 164
 - stochastic mesh, 165
- PDE jump condition, 34
- perfect foresight bias, 160
- short-maturity asymptotics, 39
- smooth pasting condition, 37, 38
- supermartingale, 31
- upper bound, *see* Monte Carlo, upper bound for American option
- annuity mapping function, *see* terminal swap rate model, annuity mapping function
- annuity measure, *see* measure, annuity
- arbitrage opportunity, 8
- arbitrage pricing, 11
- arithmetic put-call symmetry, 940
- Arrow-Debreu security, 21, 76, 78, 79, 456, 460, 1048
 - backward Kolmogorov equation, 456
 - forward Kolmogorov equation, 456
- art of derivatives trading, 980
- Asian option, 70
 - Black model, 920
 - Monte Carlo, *see* Monte Carlo, Asian option
 - PDE, *see* PDE, Asian option
- ATM backbone, *see* volatility smile, ATM backbone
- autocorrelation, *see* inter-temporal correlation
- averaging, *see* calibration, time averaging
- averaging cash flow, 201, 720–721
 - convexity adjustment, 720
- averaging swap, *see* averaging cash flow
- Bachelier model, *see* Normal model
- backbone, *see* volatility smile, backbone
- backward Kolmogorov equation, *see* Kolmogorov backward equation
- balance-guarantee swap, 898
- band swap, *see* flexi-swap
- “bang-bang”, 900
- barrier option, 44
 - Broadie adjustment for sampling frequency, *see* Monte Carlo, sampling extremes, adjusting barrier for sampling frequency
 - continuous barrier, 64
 - discrete barrier, 66
 - importance sampling, 1074–1077
 - Markovian projection, *see* Markovian projection, barrier option
- Monte Carlo, *see* Monte Carlo, barrier option
- on capped straddle, 937
- one-touch, 939
- pathwise differentiation method, 1041–1044
 - recursion, 1043
- payoff smoothing, *see* payoff smoothing, barrier option
- PDE jump condition, 66
- rebate, 64
- semi-static replication, 939
- step-down, 64
- step-up, 64
- tube Monte Carlo, 1025
- up-and-out, 44, 64, 66, 124, 126, 1134
- basis point, 169
- basis risk, *see* yield curve, basis risk
- basket option, 205, 1146
 - Black model, 922
 - displaced log-normal approximation, 1147
 - local volatility model, 1145

- Monte Carlo, *see* Monte Carlo,
 Asian option on basket
 slope of volatility smile, 1148
 stochastic volatility model, 1149
- BDT model, *see* Black-Derman-Toy
 model
- Bermudan cancelable swap, *see*
 Bermudan swaption; cancelable
 note
- Bermudan option, *see* Ameri-
 can/Bermudan option
- Bermudan swaption, 207, 873–918
 accreting, *see* Bermudan swaption,
 non-standard
 American, *see* American swaption
 amortizing, *see* Bermudan swaption,
 non-standard
 approximating American swaption,
 894
 bullet, *see* Bermudan swaption,
 vanilla
 carry, 906, 913
 impact on exercise decision, 913
 control variate, 1090
 exercise fee, 897
 exercise value, XXXVIII, 208, 873
 flexi-swap, *see* flexi-swap
 gamma-theta mismatch, 912
 hold value, XXXVIII, 208
 lockout, 207, 873
 mid-coupon, 895, 897–898
 no-call, *see* Bermudan swaption,
 lockout
 non-standard, 878–898
 calibration by payoff matching,
 882, 883
 calibration by PVBP matching,
 882–884
 calibration by tenor matching, 881
 calibration to basket, 885–887
 calibration to representative
 swaption, 882
 calibration to row of European
 swaptions, 886
 Gaussian short rate model, 886
 global calibration, 879, 881
 Libor market model, 885
 local projection method, 879, 881
 lower bound, 891, 907
- Markov-functional model, 879
 quadratic Gaussian model, 886
 quasi-Gaussian model, 879, 886,
 889
 representative swaption for
 accreting Bermudan, 884
 representative swaption for
 amortizing Bermudan, 883, 884
 super-replication, 888–892
 upper bound, 889, 890, 907
- non-vanilla, *see* Bermudan swaption,
 non-standard
- PDE jump condition, *see* Ameri-
 can/Bermudan option, PDE
 jump condition
- strike, 873
- survival measure, 1047
- vanilla, 878
- zero-coupon, 892–893
- Bermudan swaption calibration
 adjusters method, 955
 local projection method, 552,
 874–878
 Gaussian short rate model, 875
 non-standard Bermudan, *see*
 Bermudan swaption, non-
 standard
 quadratic Gaussian model, 875
 quasi-Gaussian model, 875
 smile calibration, 876–878
 at-the-money, 876
 exercise boundary, 877
 strike, 876
- Bermudan swaption greeks
 pathwise differentiation method,
 1044–1050
 forward induction, 1049–1050
 performance, 1050
 survival density, 1048
 survival measure, 1047
 portfolio replication for hedging, 911
 Principal Components Analysis, 911
 robust hedging, 910–913
 static hedging, 911
- Bermudan swaption valuation, 820–871
 control variate, 1086
 non-linear, 1089
 sampled at exercise time, 1087
 fast pricing, 914

- impact of forward volatilities, 874
 impact of inter-temporal correlation, 552, 875
 impact of mean reversion, 552, 874
 impact of the number of factors, 875
 Monte Carlo, 903–910
 - exercise strategy, 904
 - explanatory variables, 903
 - parametric lower bound, 904–910
 - regression lower bound, 903
 Bermudanality, 877
 Bessel function of the first kind, 282
 Bessel process, 281, 282
 best-of option, *see* MAX-option
 best-of-calls option, 780
 BGM model, *see* Libor market model
 Black model, XXXVIII, 22, 24, 202, 279, 283
 Asian option, *see* Asian option, Black model
 basket option, *see* basket option, Black model
 call option, 24
 CMS spread, 774
 delta, 350, 696
 effects of volatility mis-specification, 987
 Fourier integration, 329
 gamma-vega, 981
 log-likelihood ratio, 1060
 moment-generating function, 329
 PDE, 25
 stochastic interest rates, 28, 30
 strike-specific volatility, 696
 time-dependent parameters, 27, 983–985
 vega, 696
 - use in calibration, 702
 - with dividends, 28
 Black shadow rate model, 450
 Black-Derman-Toy model, 443–445
 - mean-fleeting, 445
 - short rate dynamics, 444
 Black-Karasinski model, 445
 Black-Scholes model, *see* Black model
 Black-Scholes-Merton model, *see* Black model
 BMA index, 192, 265
 BMA rate, 192
 Boltzman-Gibbs distribution, *see* out-of-model adjustment, path re-weighting method, Boltzman-Gibbs distribution
 Bond Market Association, *see* BMA index
 box smoothing method, *see* payoff smoothing, box smoothing
 break-even rate, *see* forward swap rate
 Broadie adjustment for sampling
 - frequency of barriers, *see* Monte Carlo, sampling extremes, adjusting barrier for sampling frequency
 Brownian bridge, 125, 645, 646
 - conditional moments, 129
 Libor market model, *see* Libor market model valuation, Monte Carlo, Brownian bridge
 path construction, *see* Brownian motion, path construction by Brownian bridge
 sampling extremes, *see* Monte Carlo, sampling extremes, with Brownian bridge
 Brownian motion, 4
 - geometric, 16
 - Haar function decomposition, *see* Brownian motion, path construction by Brownian bridge
 Ito integral, *see* Ito integral
 Karhunen-Loeve decomposition,
 - see* Brownian motion, path construction by Principal Components
 path construction, 106
 path construction by Brownian bridge, 128, 129
 path construction by Principal Components, 130
 Stratonovich integral, *see* Stratonovich integral
 BSM model, *see* Black model
 C^0 , XXXVIII
 C^1 , XXXVIII
 C^2 , XXXVIII
 C^n , XXXVIII
 calibration, 299

- calibration norm, 628–631
 - fit, 632
 - regularity, 632
- cold start, 631
- forward induction, 443, 456, 953
- Levenberg–Marquardt, 631
- local projection method, *see* local projection method
- Markovian projection method, *see* Markovian projection
- most likely path, 990
- stochastic optimization method, 953
- time averaging, 301, 307, 363, 370–381, 548, 581, 666
 - algorithm, 376–381
 - non-zero correlation, 376
 - skew, 373–374
 - volatility, 371–373
 - volatility of variance, 374–376
- callable Libor exotic, *see* CLE
- callable zero, *see* Bermudan swaption, zero-coupon
- cancelable note, 214, 827, 828
 - ATM, 858
 - carry, 856, 913
- cancelable swap, *see* cancelable note
- cap, 186, 202
 - caplet volatility from cap volatility, 704
 - interpolation, 705
 - precision norm, 705
 - relaxation, 706
 - smoothness norm, 706
 - splitting scheme, 706
 - digital, 203, 209
 - valuation formula, 202
- Capital Asset Pricing Model, 357
- capped floater, 209
- Cauchy distribution, 98, 101
 - Monte Carlo, 98
- certificate of deposit, 194
- CEV model, 280–286
 - attainability of zero, 280
 - displaced, 285
 - European call option value, 282, 283
 - explosion, 280
 - regularization, 284
 - relation to Bessel process, 281
 - strict supermartingale, 280
- time-dependent, 304
 - effective parameter, 305
 - volatility skew, 284
- characteristic function, 20
- Cheyette model, *see* quasi-Gaussian model
- chi-square distribution, 100
 - Monte Carlo, 100, 102
 - non-central, *see* non-central chi-square distribution
 - PDF, 100
- chooser cap, *see* flexi-cap
- chooser swap, *see* flexi-swap
- CIR model, *see* Cox–Ingersol–Ross model
- CLE, 213, 216, 626, 815–871, 873
 - accreting at coupon rate, 216, 868
 - carry, 857, 906, 913
 - impact on exercise decision, 847, 857
 - definition, 820
 - exercise value, XXXVIII, 215, 820
 - hold value, XXXVIII, 215, 820, 821
 - lockout, 213
 - marginal exercise value decomposition, 822
 - multi-tranche, 217
 - no-call, *see* CLE, lockout
 - optimal exercise, 822
 - single-rate, 862
 - smooth function of Monte Carlo path, 1029
 - snowball, 216, 870
- CLE calibration, 815–820
 - local projection method, 862–868
 - calibration targets, 863
 - core swap rate analog, 865
 - local models, 864–865
 - quadratic Gaussian model, 865
 - quasi-Gaussian model, 864
 - two-factor Gaussian model, 864
 - two-strike calibration, 865
 - vega, 867
 - low-dimensional models, 862–868
 - model choice, 819
 - single-rate, 862–863
 - to forward volatility, 819
- CLE greeks, 1036–1040
 - as sum of coupon greeks, 1037

- discontinuity in Monte Carlo, 1041
 freezing exercise boundary, 833, 1039, 1040
 freezing exercise time, 1038–1040
 likelihood ratio method, *see* likelihood ratio method
 pathwise differentiation method, 1035–1040, 1058–1060
 computational complexity, 1052
 forward induction, 1049–1050
 survival density, 1048
 survival measure, 1047
 perturbation method, 1040, 1059
 computational complexity, 1053
 portfolio replication for hedging, 911
 recursion, 1036
 source of noise, 1040
 tube Monte Carlo, 1029
 CLE regression, 823–862
 automatic selection of regression variables, 855
 boundary optimization, 831
 cancelable note, 827–828
 choice of regression variables, 848–854
 decision only, 828–830
 discrepancy principle, 859
 excluding suboptimal points, 856
 exercise value, 825–827
 explanatory variables, 850–854
 classification, 851
 CMS spread, 851
 core swap rate, 851
 stochastic volatility, 854
 with convexity, 852–854
 general-to-specific approach, 856
 generalized cross-validation, 859
 L-curve method, 859
 Libor market model, 849, 850
 state variables, 849
 lower bound, 831–833
 perfect foresight bias, 832
 pseudo-inverse method, 860
 quadratic Gaussian model, 849
 quasi-Gaussian model, 849
 regression operator, 824
 regression variables, 823
 rescaling, 861
 reuse exercise boundary, *see* CLE greeks, freezing exercise boundary
 ridge regression, *see* CLE regression, Tikhonov regularization
 robust implementation, 858–862
 singular value decomposition, 104
 stabilization, 859
 state variables, 848–849
 Libor market model, 849
 SVD decomposition, 860, 861
 connection to Tikhonov regularization, 861
 Tikhonov regularization, 162, 255, 859–861
 connection to SVD, 861
 truncated SVD decomposition, 162, 860, 861
 two-step, 857
 upper bound, 837–848
 alternative methods, 847
 computational cost, 841
 improvements to algorithm, 845–847
 nested simulation algorithm, 837–847
 non-analytic exercise values, 843–845
 simulation within a simulation, *see* CLE regression, upper bound, nested simulation algorithm
 CLE valuation, 215, 820–871
 as cancelable note, 827
 boundary optimization, 831
 confidence interval for value, 842
 control variate, *see* Bermudan swaption valuation, control variate
 discontinuous function of Monte Carlo path, 1041
 duality, 836, 1093
 multiplicative, 1093
 duality gap, 839, 842, 908, 909
 in stochastic volatility models, 910
 exercise policy consistency conditions, 833
 fast pricing, 916
 Hamilton-Jacobi-Bellman equation, 821

- impact of forward volatility, 818
- impact of inter-temporal correlation, 863
- impact of volatility smile dynamics, 819
- Libor market model, 824
- lower bound, 834, 841, 845, 848
 - by regression, *see* CLE regression, lower bound
 - iterative improvement, 833
 - iterative improvement by nested simulation, 835
 - quality test, 1060
- LS method, *see* CLE regression
- Monte Carlo, 823–862, 903
- optimal exercise policy, 833, 835, 1039
- PDE, 868–871
 - accreting at coupon rate, 868
 - path-dependent, 868–871
 - similarity reduction, 869
 - snowball, 870
- perfect foresight bias, 832
- policy fixing, 846
- recursion, 821
- regression method, *see* CLE regression
- tube Monte Carlo, 1029
- upper bound, 836–848
 - cancelable note, 844
 - nested simulation algorithm, 839, 908
 - non-analytic exercise values, 843–845
- weighted coupon decomposition, 916
- CMS, 206
 - annuity to forward measure change, 734–737
 - convexity adjustment, 721–744
 - annuity mapping function, *see* terminal swap rate model, annuity mapping function
 - correcting arbitrage, 732–733
 - density integration method, 736
 - impact of mean reversion, 733–734
 - impact of volatility smile, 733
 - impact on implied volatility, 774
 - Libor market model, 729–731
 - linear TSR model, 726–728
 - out-of-model adjustment, 963, 964
 - quasi-Gaussian model, 728–729
 - replication method, 722–724
 - stochastic volatility model, 738
 - swap-yield TSR model, 726
 - vega hedging, *see* terminal swap rate model, linear TSR model, vega hedging
 - hedging portfolio, 723
 - quanto, *see* quanto CMS
- CMS cap, 207, 695
 - impact of CMS convexity on volatility smile, 739
 - link to European swaptions, 739
- CMS digital spread option, 789
 - dimensionality reduction, 789
- CMS floor, 207
- CMS rate, 206
 - distribution in forward measure, 734–737
- CMS spread option, 210, 211, 619, 688, 763, 774
 - by integration, 775
 - copula method, 774–782
 - dimensionality reduction, 787
 - floating digital, 790
 - Gaussian copula, 775
 - correlation impact, 776
 - vega to swaptions, 776
 - implied copula, 779
 - implied correlation, 776
- Libor market model, 617–619, 634, 690, 806
 - closed-form approximation, 808
- Libor market model calibration, 634
- local volatility model, 1145
- Margrabe formula, 810
- Markovian projection, 1145, 1149
- multi-stochastic volatility, *see* multi-stochastic volatility model
- non-standard gearing, 775, 789
 - dimensionality reduction, 789
- Normal spread volatility, 774
- one-dimensional integration, 787
- out-of-model adjustment, 964, 966
- power Gaussian copula, 779
- quadratic Gaussian model, 808
 - closed-form approximations, 808

- risk management with one-factor model, 971
- stochastic volatility
- correlation impact, 805
- stochastic volatility de-correlation, 962
- stochastic volatility model, 1149
- correlation impact, 803
 - vega in Libor market model, 1116
- CMS swap, 206, 695
- valuation formula, 207
- CMS-linked cash flow, 721–744
- direct integration method, 734
 - replication method, 723
- coherent risk measure, *see* risk measure, coherent
- collateral, 192, 266
- complementary Gamma function, 281
- complete market, 11
- compounded rate, 200
- conditional expected value, 19
- iterated conditional expectations, *see* iterated conditional expectations
- projection approximation, *see* Markovian projection, conditional expected value by projection
- constant elasticity of variance model, *see* CEV model
- constant maturity swap, *see* CMS swap
- contingent claim, *see* derivative security
- continuity correction, *see* payoff smoothing, continuity correction
- control variate, 146–149, 330, 652, 653, 1077–1094
- adjusters method, 955
 - construction from MC upper bound, 1093
- dynamic, 148, 653, 1090–1093
- regression-based, 1091
- efficiency, 147
- impact on risk stability, 1093
- instrument-based, 1086–1090
- model-based, 675, 1077–1086
- non-linear controls, 147–149
- path re-weighting method, 961
- proxy Markov LM model, 1078
- proxy model, *see* control variate, model-based
- convexity adjustment
- averaging swap, *see* Libor-with-delay, convexity adjustment
- CMS, *see* CMS, convexity adjustment
- futures, *see* ED future, convexity adjustment
- Libor-in-arrears, *see* Libor-in-arrears, convexity adjustment
- Libor-with-delay, *see* Libor-with-delay, convexity adjustment
- moment explosion, 759–762
- second moment, 759
- copula, 768
- Archimedean, 770
 - Monte Carlo, 798
 - Clayton, 770
 - conditional CDF, 790
 - Frechet bounds, 769
 - Gaussian, 766
- CMS spread option, *see* CMS spread option, Gaussian copula integration, 787
 - joint CDF, 767
 - joint PDF, 767, 775
 - mixture, 772
 - Monte Carlo, 797
- Gumbel, 770, 771
- implied, 779
- independence, 768
- mixture, 772
- Monte Carlo, 798
- perfect anti-dependence, 769
- perfect dependence, 768
- power Gaussian, 773, 778
- parameter impact, 779
- product, 773
- Monte Carlo, 798
- reflection, 771
- Monte Carlo, 798
- Sklar's theorem, 769
- copula density, 770
- copula method, 766
- CMS spread option, *see* CMS spread option, copula method
- dimensionality reduction, 787–796
- by conditioning, 791–795

- by measure change, 795–796
- forward swaption straddle, 949
- integration, 784–796
 - inverse CDF caching, 785
 - singularities, 786
- limitations, 799–800
- mapping function, 793
- Monte Carlo, 797–799
- observation lag, 782
- quanto options, 747
- volatility swap, 934
- core correlations, *see* inter-temporal correlation
- core volatilities, 863, 874
- correlation extractor, *see* Libor market model, correlation extractor
- correlation risk sensitivity, 1119
- correlation smile, 776
- Cox-Ingersol-Ross model, 430
 - multi-factor, 518
 - two-factor, 516
- Crank-Nicolson scheme, *see* PDE, Crank-Nicolson scheme
- credit risk, 260, 975
- credit value adjustment, 266, 914
- cross-currency basis swap, *see* floating-floating cross-currency basis swap
- cross-currency basis swap spread, 262, 265
- CRX basis swap, *see* floating-floating cross-currency basis swap
- CRX spread, *see* cross-currency basis swap spread
- cumulant-generating function, 154
- curve cap, 211, 764
 - range accrual, *see* range accrual, curve cap
- CVA, *see* credit value adjustment
- date rolling convention, 224
- day count convention, 223–226
 - 30/360, 225
 - Actual/360, 224
 - Actual/365.25, 224
- day count fraction, *see* year fraction
- deflator, 9
- delta, 18, 132, 355, 980
- bucketed interest rate deltas, 251, 1045
- forward rate, 253
- Jacobian method, *see* risk sensitivities, Jacobian method
- par-point, 251, 252, 256, 257, 993
- parallel, 257
 - with backbone, 1120–1122
- delta hedge, 18
- density process, 9
- derivative security, 11
 - attainable, 11
 - pricing, 11
- diffusion, 4, 15
 - absorbing barrier, 281, 289
 - displaced, 285
 - Feller boundary classification, 280
 - Feller condition, 319
 - Fubini's theorem, 407
 - integration by parts, 120
 - Ito integral, *see* Ito integral
 - Ito process, 4
 - local time, 26, 294
 - Ornstein-Uhlenbeck process, 411
 - polynomial growth condition, 19
 - predictable process, 7
 - scale measure, 280
 - SDE, 15
 - generator, 19
 - linear, 16
 - locally deterministic, 172, 539
 - strong Markov, 15
 - strong solution, 15
 - weak solution, 15
 - speed measure, 280
- diffusion invariance principle, 14
- discount bond, XXXVIII, 23, 167
 - valuation formula, 172
- discount curve, *see* yield curve
- displaced CEV model, *see* CEV model, displaced
- displaced log-normal model, 285
 - basket option, 1147
 - canonical form, 286
 - explicit solution to SDE, 312
 - Fourier integration, 328
 - implied correlation, 809
 - moment matching, 920
 - moment-generating function, 329

- time-dependent, 304
- effective skew, 305
- explicit solution to SDE, 307
- range for process, 306
- Dupire local volatility, 1131
 - proof by Tanaka extension, 294, 1131
- duration, 246
- DVF model, *see* local volatility model
- Dybvig parameterization, *see* short rate model, Dybvig parameterization
- early exercise, 30
- ED future, 168–170, 196–197, 695, 748–759
 - convexity adjustment, 187, 197, 748–759
 - from market inputs, 751
 - Gaussian HJM model, 186
 - impact of volatility smile, 750, 756
 - Libor market model, 751, 756
 - replication method, 751, 755
- delivery arbitrage, 170
- futures rate, 169
 - definition, 196
 - instantaneous, 170, 172, 173
 - martingale in risk-neutral measure, 172, 749
 - martingale in spot Libor measure, 749
 - simple, 169
 - to forward rate, 754, 758
- mark to market, 169
- yield curve construction, 231, 992
- ED futures contract, *see* ED future
- effective volatility
 - local volatility model, *see* local volatility model, effective volatility
 - stochastic volatility model, *see* stochastic volatility model, effective volatility
- envelope theorem, 1038
- Eonia, 193, 200
- equivalent martingale measure, *see* measure, equivalent martingale
- Esscher transform, *see* exponential twisting
- Eurodollar futures contract, *see* ED future
- European call option, 24
 - at-the-money, 24
 - Fourier integration, 324
 - in-the-money, 24
 - out-of-the-money, 24
 - probability density from, *see* volatility smile, probability density from
- European digital call option, 60
- European option
 - Fourier integration, 326
- European put option, 24
 - at-the-money, 24
 - in-the-money, 24
 - out-of-the-money, 24
- European swaption, 203, 695–703
 - cash-settled, 205, 742–744
 - payoff, 743
 - put-call parity, 743
 - replication method, 742, 743
- core swaptions, 422, 817
- coterminal swaptions, *see* European swaption, core swaptions
- diagonal swaptions, *see* European swaption, core swaptions
- forward swaption straddle, *see* forward swaption straddle, 943
- midcurve, 223
- non-standard, *see* Bermudan swaption, non-standard
 - Black formula, 887
- physically-settled, 205
- SV model calibration, 701–702
- swap-settled, 205, 743
- swaption grid, 205, 701
- swaption strip, 421
- tenor, 204
- valuation formula, 204
- volatility cube, 696
- European-style option, 95
 - replication method, 337
 - valuation by volatility mixing, 339
- exchange market, 193
 - Chicago Mercantile Exchange, 196
 - London International Financial Futures and Options Exchange, 196

- Marché à Terme International de France, 196
- exotic swap, 205, 208, 209, 820, 951
 - CMS spread, 764
 - CMS-based, 210
 - digital CMS spread, 764
 - global cap, 219
 - global floor, 219
 - knock-out, 218
 - Libor-based, 209
 - multi-rate, 210, 764
 - path-dependent, 212
 - principal amount, 208
 - range accrual, *see* range accrual
 - snowball, 212
 - spread-based, 210
 - structured coupon, 208–211
 - expectations hypothesis, 173
 - expected hedging P&L, 988
 - exponential distribution, 98
 - Monte Carlo, 98
 - exponential integral, 334
 - exponential twisting, 154
 - extra state variable method, *see* PDE, path-dependent options
- “The Fed Experiment”, 450
- Federal funds future, 201
- Federal funds rate, 192, 200, 201, 266
 - effective, 192
 - target, 192
- Federal funds/Libor basis swap, 201, 266
- Feller condition, *see* diffusion, Feller condition
- Feynman-Kac solution, 21
- FFT, *see* stochastic volatility model, Fourier integration
- filtration, 3, 4
 - usual condition, 3
- flexi-cap, 71
- flexi-swap, 898–903
 - decomposition into Bermudan swaptions, 899
- local projection method, 899
- marginal exercise value decomposition, 901
- narrow band limit, 902
- PDE, 899, 901
 - purely local bounds, 899
 - “flip-flop”, 210
- floating digital, 790, 792
 - dimensionality reduction, 790
- floating digital spread option, 790
 - dimensionality reduction, 790
- floating-floating cross-currency basis swap, 262, 264, 265
- floating-floating single-currency basis swap, 201, 268
- floor, *see* cap
- Fokker-Plank equation, *see* Kolmogorov forward equation
- Fong-Vasicek model, 452–453, 515
 - bond reconstitution formula, 452
- forward CMS straddle, 941, 944, 945
 - swaption, *see* forward swaption straddle
 - volatility, *see* forward volatility
- forward contract, 195
- forward Kolmogorov equation, *see* Kolmogorov forward equation
- forward Libor model, *see* Libor market model
- forward Libor rate, XXXVIII, 168, 191, 192, 196
 - accrual end date, 224
 - accrual period, 224
 - accrual start date, 224
 - martingale in forward measure, 174
 - tenor, 168
 - variance by replication method, 756
 - year fraction, *see* year fraction
- forward par rate, *see* forward swap rate
- forward price, 24, 168
- forward rate, 167
 - continuously compounded, XXXVIII, 168
 - instantaneous, XXXVIII, 169
 - simple, 168
 - tenor, 168
 - volatility hump, 416, 492
- forward rate agreement, *see* forward contract
- forward starting option, 222
- forward swap rate, XXXVIII, 171, 199

- distribution in forward measure,
see CMS rate, distribution in forward measure
 expiry, 171
 fixing date, 171
 linking forward and annuity measure, 735
 market-implied variance, 555
 martingale in swap measure, 178
 non-standard, 879
 decomposition, 880
 tenor, 171
 weighted average of Libor rates, 171, 256
 forward swaption straddle, 223, 945–950
 copula method, 949
 relation to CMS spread option, 948
 triangulation, *see* forward volatility, triangulation
 vanilla model, 946
 vega exposure, 948
 volatility, *see* forward volatility
 forward volatility, 222
 connection to inter-temporal correlations, *see* inter-temporal correlation, connection to forward volatilities
 hedging, 912
 impact of rate correlation, 918
 impact of volatility smile, 945
 Libor rate, *see* volatility, forward volatility of Libor rate
 triangulation, 948
 forward volatility derivative, 220, 222
 forward swaption straddle, *see* forward swaption straddle
 implied Normal volatility contract, 223
 midcurve swaption, *see* European swaption, midcurve
 volatility swap, *see* volatility swap
 forward yield, *see* forward rate
 Fourier transform, 325
 inverse, 325
 FRA, *see* forward contract
 Frobenius norm, *see* matrix, Frobenius norm
 fundamental matrix, 484
 fundamental theorem of arbitrage, 10
 fundamental theorem of derivatives trading, 987
 futures contract, *see* ED future
 futures rate, *see* ED future, futures rate
 fuzzy logic, *see* payoff smoothing, fuzzy logic
 FX rate, 179, 745, 746
 dynamics in domestic risk-neutral measure, 180
 forward, 178
 martingale in domestic forward measure, 180
 Gâteaux derivative, 253
 gamma, 980
 pathwise differentiation method, *see* pathwise differentiation method, gamma
 payoff smoothing, 1019
 relationship to vega, 981
 gamma distribution, 100
 Monte Carlo, 100, 102
 PDF, 100
 Gamma function, XXXVII
 incomplete, *see* incomplete Gamma function
 quick approximation, 1153
 Gauss-Hermite quadrature, *see* quadrature, Gauss-Hermite
 Gaussian copula, *see* copula, Gaussian
 Gaussian distribution, XXXVII
 conditional distribution, 646
 cumulant-generating function, 154
 imaginary mean, 796
 inverse CDF, 99, 165
 linear transform, 103
 measure change, 795
 multi-dimensional PDF, 103
 quadratic form, 522
 moment-generating function, 522, 533
 moments, 534
 Gaussian HJM model, 184–187
 caplet, 186
 ED future convexity adjustment, *see* ED future, convexity adjustment, Gaussian HJM model

- time-stationary, 416
- zero-coupon bond option, 185
- Gaussian multi-factor short rate model,
 - see* Gaussian short rate model, multi-factor
- Gaussian one-factor short rate model,
 - see* Gaussian short rate model
- Gaussian short rate model, 406, 413–429, 478–510
 - as special case of affine model, 430
- Bermudan swaption, *see* Bermudan swaption calibration, local projection method, Gaussian short rate model
- bond dynamics, 415
- bond reconstitution formula, 414
 - efficient calculation, 415
- calibration, 421
 - bootstrap, 422
- calibration to yield curve, 414
- European swaption, 418, 421
 - Jamshidian decomposition, 418
- fast pricing of Bermudan swaptions, 914
- forward rate dynamics, 413
- forward rate volatility, 413
 - dynamics, 417
- humped volatility structure, 416
- in spot measure, 428
- in terminal measure, 428
- mean reversion, *see* mean reversion
- mean reversion calibration, *see* mean reversion calibration
- Monte Carlo, 425–429
 - approximate, 427
 - Euler scheme, 427
 - exact, 425
 - other measures, 428
- multi-factor, 478–510
 - benchmark rate parameterization, 506–508
 - benchmark rates, 506
 - benchmark tenors, 506
 - bond reconstitution formula, 478, 481, 483
 - bond volatility, 479
 - calibration, 506
 - classic development, 485–488
 - correlated Brownian motions, 489
- correlation stationarity, 488
- European swaption, 500–505
- European swaption by Jamshidian decomposition, 503
- factors and loadings, *see* Gaussian short rate model, multi-factor, statistical approach
- forward rate correlation, 488–489
- forward rate volatility, 482
- Gaussian swap rate approximation, 504–505
- loadings, 499
- mean reversion matrix diagonalization, 487–488
- Monte Carlo, 508–509
- PDE, 510
 - rotations, 484
 - separability, 478–485
 - short rate dynamics, 479
 - short rate state distribution, 485, 509
 - short rate state dynamics, 479–485
 - short rate state dynamics, integrated, 485, 509
 - single Brownian motion, 496
 - statistical approach, 495–500
 - swap rate volatility, 505
- PDE, 423–425
 - boundary conditions from PDE, 424
- short rate distribution, 426
- short rate dynamics, 413
- short rate state dynamics, 414, 425
 - integrated, 425
- swap rate dynamics in annuity measure, 420
- swap rate volatility, 420
- time-stationary, 416
- two-factor, 489–495
 - bond reconstitution formula, 490, 500
 - CLE, *see* CLE calibration, local projection method, two-factor Gaussian model
 - correlated Brownian motions, 490
 - correlation stationarity, 491
 - doubly mean-reverting form, 493
 - European swaption by Jamshidian decomposition, 500–504

- forward rate correlation, 490–491
- forward rate dynamics, 490
- forward rate volatility, 490–491, 493, 494
- short rate state conditional distribution, 502
- short rate state correlation, 490
- short rate state dynamics, 490
- single Brownian motion, 495
- volatility hump, 492–493
- Gaussian two-factor short rate model, *see* Gaussian short rate model, two-factor
- generalized trigger product, 1074
- importance sampling, 1074–1077
- pathwise differentiation method, 1041–1044
- payoff smoothing, 1074–1077
- trigger variable, 1074
- tube Monte Carlo, *see* barrier option, tube Monte Carlo
- Girsanov's theorem, 12, 13
- Gaussian distribution, 795
- Gram-Charlier expansion, 368, 437
- greeks, *see* risk sensitivities
- Green's function, 20
- grid shifting, *see* payoff smoothing, grid shifting
- GSR model, *see* Gaussian short rate model
- Gyöngy theorem, *see* Markovian projection, Gyöngy theorem

- H^2 , 5
- Hagan and Woodward parameterization, *see* short rate model, Hagan and Woodward parameterization
- hat smoothing method, *see* payoff smoothing, hat smoothing
- Heath-Jarrow-Morton model, *see* HJM model
- hedge, 251
- best hedging strategy, 355
- beta, 357
- minimum variance, 355–357
- model-independent, 716
- semi-static, *see* replication method, semi-static

- shadow delta, *see* volatility smile, shadow delta hedging
- sub-replicate, 717
- super-replicate, 717, 979
- zero-beta, 357
- Hermite matrix, 270
- Heston model, *see* stochastic volatility model
- HJM model, 181–189
- bond dynamics, 181
- forward bond dynamics, 182
- forward rate dynamics, 182
- Gaussian, *see* Gaussian HJM model
- Gaussian Markov, 187–189
- short rate dynamics, 188
- log-normal, 189
- Markovian, 405
- separable, 413
- short rate dynamics, 183
- stochastic basis, *see* HJM model, two-curve
- two-curve, 678–681
- forward rate spread dynamics, 679
- Gaussian basis spread, 681
- index bond dynamics, 680
- index forward rate dynamics, 680
- index short rate dynamics, 680
- quanto correction, 681
- Ho-Lee model, 406–410
- bond dynamics, 409
- bond reconstitution formula, 408
- calibration to yield curve, 407
- drawbacks, 410
- forward rate dynamics, 409
- short rate dynamics, 408
- hybrid differentiation method, 1061

- implied volatility, *see* volatility, implied
- importance sampling, 146, 149–158, 1063–1077
- application to payoff smoothing, 1067
- barrier option, *see* barrier option, importance sampling
- density formulation, 149
- efficiency, 151

- generalized trigger product, *see* generalized trigger product,
- importance sampling
- least-squares, 154
- likelihood ratio, 150, 153, 155
- rare events, 154
 - approximately optimal mean shift
 - in multi-variate case, 158
 - asymptotic optimality, 158
 - efficiency, 156
 - minimal variance, 155
 - multi-variate, 156
- SDE, 151–154
- short rate model, *see* short rate model, importance sampling
- survival measure, 1067
 - simulation under, 1072, 1074, 1076
- TARN, *see* TARN, importance sampling
- incomplete Gamma function, XXXVII, 281
- index, 206
 - index option, *see* basket option
- infinitesimal operator of SDE, *see* diffusion, SDE, generator
- infinitesimal perturbation analysis, 136
- information theory, 957
- instantaneous futures rate, *see* ED future, futures rate, instantaneous integration by parts for diffusion process, *see* diffusion, integration by parts
- inter-temporal correlation, 422, 552, 818, 863, 874
 - connection to forward volatilities, 818
 - hedging, 875, 912
 - impact of mean reversion, 552
 - impact of volatility smile, 945
 - impact on Bermudan swaption, *see* Bermudan swaption valuation, impact of inter-temporal correlation
 - impact on CLEs, *see* CLE valuation, impact of inter-temporal correlation
 - impact on TARNs, 929
- mean reversion calibration to, *see* mean reversion calibration, to inter-temporal correlations
- interbank money market, 192
- International Swaps and Derivatives Association, 192, 266
- intrinsic value, 27
- inverse floater, 209
- iterated conditional expectations, 176
- Ito integral, 4, 5
- Ito isometry, 5
- Ito's lemma, 6
- Ito-Taylor expansion, 118
- Jacobian, *see* risk sensitivities, Jacobian method
- Jamshidian decomposition
 - American/Bermudan option, *see* American/Bermudan option, Jamshidian decomposition
 - European swaption, *see* Gaussian short rate model, European swaption, Jamshidian decomposition
- Kolmogorov backward equation, 19, 20
- Kolmogorov forward equation, 20, 386, 457, 1048
 - correct boundary conditions, 386
 - discrete consistency with backward equation, 458
- Kullback-Leibler relative entropy, 957
- kurtosis, 375
- L^1 , XXXVIII, 4
- L^2 , XXXVIII, 4
- ladder, 985
- ladder swap, *see* ratchet swap
- Lagrange basis functions, *see* PDE, Lagrange basis; payoff smoothing, Lagrange basis
- Lagrange multiplier, 249, 958
- least squares method, *see* CLE regression
- LIA, *see* Libor-in-arrears
- Libor curve, *see* yield curve
- Libor market model, 449, 589–692, 729, 866, 910
 - annuity mapping function, 730, 731
 - asset-based adjustment, 963

- back stub, 655–660
 - arbitrage-free, 657–659
 - from Gaussian model, 659–660
 - simple, 656–657
- choosing number of factors, 612
- CLE, 819
- CMS convexity adjustment, 964
- correlation extractor, 863
- deflated bond dynamics, 649
- delta with backbone, 1120–1122
- drift approximation, 644
 - Brownian bridge, 1079
- drift freezing, 1052
- exercise boundary, 910
- exercise strategy, 907
- expected value of Libor rate in annuity measure, 669
- front stub, 660–666
 - exogenous volatility, 661–664
 - from Gaussian model, 665–666
 - simple interpolation, 664–665
 - zero volatility, 660–661
- in hybrid measure, 640
- index function, *see* tenor structure, index function
- Libor rate correlation, 601–612, 757
 - correlation PCA, 609
 - covariance PCA, 624
 - historical estimation, 604
 - majorization, 611
 - parametric form, 606, 607
 - PCA, 602–604
 - poor man’s correlation PCA, 612
 - regularization, 608
- Libor rate dynamics, 591–601
 - annuity measure, 731
 - in forward measures, 592–593
 - in hybrid measure, 595
 - in spot measure, 594
 - in terminal measure, 594, 639
- Libor rate inter-temporal correlation, 757
- Libor rate volatility
 - from volatility norm, 623–625
 - functional form, 620
 - grid-based, 620–621
 - interpolation, 622–623
- Libor rate volatility link to HJM
 - forward rate volatility, 596
- link to HJM, 595
- local volatility, 596–598
 - CEV, 597
 - displaced log-normal, 597
 - existence and uniqueness, 597, 598
 - LCEV, 597
 - log-normal, 597
- Markov, 674–675, 1078–1086
 - as control variate, 1084
 - Brownian bridge, 1079
 - calibration, 1082
 - one-factor, 1079
 - one-factor reconstitution formula, 1080
 - separable volatility, 1080
 - two-factor, 1081
 - two-factor reconstitution formula, 1081
- Markovian projection, 666, 668, 1139
- model risk, 627
- multi-stochastic volatility, 688–692, 962
 - caplet, 690
 - CMS spread option, 690
 - European swaption, 690
 - moment-generating function, 690
- Musiela parameterization, 602
- pathwise derivative
 - forward Libor rate, 1051
 - forward swap rate, 1055
 - numeraire, 1054
 - structured coupon, 1055
 - stub bond, 1054
- pathwise differentiation method, 1051–1058
 - computational complexity, 1052
- PCA, *see* Principal Components Analysis
- portfolio replication, 912
- stochastic basis, *see* Libor market model, two-curve
- stochastic variance dynamics, 688
- stochastic volatility, 599–601
 - moment-generating function, 687
 - non-zero correlation, 686
- stub volatility, 662, 666
- swap rate correlation, 618–619
- swap rate dynamics, 615, 667

- approximate, 616
- time-stationary, 621
- tool to extract forward volatility, 819
- two-curve, 682–686
 - deterministic spread, 685
 - European swaption, 684
 - Libor rate dynamics, 683
 - Monte Carlo, 684
 - swap rate dynamics, 684
- vega, *see* vega, Libor market model
- Libor market model calibration,
 - 620–635
 - algorithm, 631, 634, 674
 - bootstrap, 633
 - for vega, 1111
 - cascade, *see* Libor market model calibration, bootstrap
 - choice of instruments, 625
 - effective skew, 670
 - effective volatility, 669
 - global, 626
 - grid-based, *see* Libor market model calibration, global
 - local, 626
 - objective function, 628
 - PCA, 624
 - row-by-row, 631, 632
 - to spread options, 633, 806
 - volatility skew, 635
 - volatility smile, 672
- Libor market model valuation
 - Bermudan swaption, *see* Bermudan swaption valuation, Monte Carlo caplet, 613
 - CLE, *see* CLE valuation, Libor market model
 - CMS convexity adjustment, *see* CMS, convexity adjustment, Libor market model
 - CMS spread option, *see* CMS spread option, Libor market model
 - curve interpolation, 655–666
 - European swaption, 614, 616, 666
 - Libor-with-delay, *see* Libor-with-delay, Libor market model
 - Monte Carlo, 635
 - analysis of computational effort, 637
 - antithetic variates, 652
 - Brownian bridge, 645
 - choice of numeraire, 640
 - control variate, 652
 - discretization bias, 637
 - Euler scheme, 636
 - front stub, 662
 - high-order schemes, 648
 - importance sampling, 653
 - lagging predictor-corrector, 642
 - large time steps, 639, 644–647
 - log-Euler scheme, 636
 - martingale discretization, 648–651
 - Milstein scheme, 648
 - predictor-corrector, 641, 642, 645, 651
 - survival measure, 1072, 1075
 - two-curve, 684
 - variance reduction, 651–653
 - multi-rate vanilla derivative, 806
 - PDE, *see* Libor market model, Markov
 - TARN, *see* TARN, Libor market model
 - volatility swap, *see* volatility swap, Libor market model
- Libor rate, *see* forward Libor rate
- Libor-in-arrears, 200, 714–717
 - convexity adjustment, 715
 - replication method, 716
 - sub-replicating portfolio, 717
 - super-replicating portfolio, 717
- Libor-with-delay, 717–721
 - convexity adjustment, 718
 - Libor market model, 718, 720
 - quasi-Gaussian model, 718, 719
 - replication method, 718, 720
 - swap-yield TSR model, 718
- likelihood ratio method, 139–142, 1060–1061
 - discontinuous payoff, 138
 - exploding variance, 1061
 - for Euler scheme, 141–142
 - for Milstein scheme, 142
 - log-likelihood ratio, 140
 - score function, 140
 - vega, 1124
- linear regression, 146
- Lipschitz function, 137

- LM model, *see* Libor market model
 local projection method, 558, 862, 863, 953, 1097
 Bermudan swaption, *see* Bermudan swaption calibration, local projection method
 CLE, *see* CLE calibration, local projection method
 non-standard Bermudan swaption, *see* Bermudan swaption, non-standard, local projection method
 TARN, *see* TARN, local projection method
 volatility swap, *see* volatility swap, local projection method
 local stochastic volatility model, 316, 1137–1145
 calibration, *see* Markovian projection, LSV calibration
 Markovian projection, *see* Markovian projection, LSV calibration
 local time, *see* diffusion, local time
 local volatility model, 277–312
 approximation with displaced log-normal model, 286
 asymptotic expansion, 295–299
 basket option, *see* Markovian projection, basket option in LV model
 CEV, *see* CEV model
 displaced log-normal, *see* displaced log-normal model
 effective convexity, 307–312
 effective skew, 301–312
 effective volatility, 301
 expansion around displaced log-normal model, 296
 expansion around Gaussian model, 298
 forward equation for call options, 293
 PDE, 292–295
 simultaneous for multiple parameters, 293
 space discretization, 292
 transform to constant diffusion coefficient, 88, 292
 quadratic volatility, *see* quadratic volatility model
 range-bound, 287
 small-noise expansion, *see* volatility, small-noise expansion
 smile dynamics, 279, 350, 352
 time-dependent, 299–312
 separable, 300
 log-normal distribution, XXXVII, 16
 moment matching, *see* moment matching
 moments, 16
 Monte Carlo, 101
 Longstaff-Schwartz method, *see* CLE regression
 Longstaff-Schwartz model, 516–517
 bond reconstitution formula, 516
 lookback option, 124
 Monte Carlo, *see* Monte Carlo, lookback option
 LS method, *see* CLE regression
 LSV model, *see* local stochastic volatility model
 LVF model, *see* local volatility model
 Malliavin calculus, 142, 1042, 1060
 Margrabe formula for spread option, 810
 mark-to-model, 816
 Markov process, 15
 Feynman-Kac theorem, *see* Feynman-Kac solution
 strong, 15
 transition density, 20
 Markov-functional model, 470–476
 calibration to yield curve, 473
 criticism, 476
 Libor parameterization, 471
 log-normal, 472
 no-arbitrage condition, 471
 non-standard Bermudan swaption, 879
 numeraire, 470
 numeraire mapping, 470
 Libor parameterization, 471
 non-parametric, 474
 swap parameterization, 474
 PDE, 475
 state process, 470

- swap parameterization, 473
 transition density, 470
 Markovian projection, 803, 1129–1156
 average option, 1133
 barrier option, 1134
 basket option in LV model,
 1145–1148
 basket option in SV model,
 1149–1152
 CMS spread option, 1145
 conditional expected value by
 Gaussian approximation,
 1134–1135
 conditional expected value by
 projection, 725, 1136–1137
 displaced Heston model, 1149, 1151
 non-perturbative approximation,
 1151
 displaced log-normal model, 1136,
 1146
 Gyöngy theorem, 1130
 LSV calibration, 1139–1145
 mapping function, 1142
 proxy model, 1143–1145
 quadratic volatility model, 1137,
 1148
 quasi-Gaussian model, *see* quasi-
 Gaussian model, Markovian
 projection
 spread option, 1151
 stochastic volatility model, 1138
 martingale, 5
 Doob-Meyer decomposition, 35
 exponential, 12
 Doleans exponential, XXXVII, 12
 local, 5
 bounded, 288
 martingale representation theorem,
 6
 Novikov condition, 12
 optional sampling theorem, 35
 Snell envelope, 31, 821
 square-integrable, 5
 stopping time, *see* stopping time
 submartingale, 5
 supermartingale, 5, 360
 CEV, *see* CEV model, strict
 supermartingale
 quadratic volatility, *see* quadratic
 volatility model, strict super-
 martingale
 SV model, *see* SV model with
 general variance process, strict
 supermartingale
 matrix
 exponential, 484
 Frobenius norm, 105, 608, 609, 624,
 625, 849
 infinity norm, 53
 positive semi-definite, 103
 Cholesky decomposition, 103
 rank-deficient, 106
 spectral norm, 53
 stiffness, 1111
 tri-diagonal, 47
 MAX-option, 906
 mean reversion, 316, 411, 550, 571
 effects, 550–552
 inter-temporal correlation, 552
 swaption volatility ratio, 551
 mean reversion calibration, 550–558,
 571
 to inter-temporal correlations,
 555–557
 to row of European swaptions, 553,
 886
 to volatility ratios, 552–555
 mean-reverting square-root process,
 see square-root process
 measure, XXXVII
 absolutely continuous, 1067
 annuity, 178, 204
 change of numeraire, *see* numeraire,
 change of numeraire
 domestic, 744
 equivalent, 9, 1067
 equivalent martingale, 8, 9, 14, 171
 foreign, 744
 hybrid, 176
 local martingale, 10
 risk-neutral, XXXVII, 23, 172
 domestic and foreign, 179, 180
 spot, XXXVII, 175
 survival density, 1047
 survival for Bermudan swaption,
 see Bermudan swaption, survival
 measure

- survival in importance sampling, *see* importance sampling, survival measure
- T-forward, XXXVII, 29, 174
 - domestic and foreign, 180
 - terminal, 176
- min-max volatility swap, 222, 938
 - capped, 940
 - semi-static replication, 939
- moment explosion, 323, 343, 344, 361, 759, 760
 - impact on convexity adjustment, *see* convexity adjustment, moment explosion
- SABR model, *see* SABR model, moment explosion
- stochastic volatility model, *see* stochastic volatility model, moment explosion
- SV model with general variance process, *see* SV model with general variance process, moment explosion
- moment matching, 887, 919–923
 - Asian option, 920
 - basket option, 922
- moment-generating function, 13
- Monte Carlo, 95–165
 - A-stable scheme, 110
 - Asian option, 107
 - Asian option on basket, 107
 - average rate option, *see* Monte Carlo, Asian option
 - barrier option, 124–128
 - adjusting barrier for sampling frequency, 128
 - double-barrier knock-out, 124
 - bias, 122
 - bias/standard error trade-off, 123
 - Brownian motion, *see* Brownian motion
 - calibration by stochastic optimization method, 953
 - central limit theorem, 96
 - convergence rate, 97
 - discretization bias, 426
 - efficiency, 144
 - Euler scheme, 110, 111
 - linear SDE, 112
 - region of stability, 111
 - weak convergence order, 111
- Euler-Maruyama scheme, *see* Monte Carlo, Euler scheme
- Heun scheme, 116
- higher-order schemes, 116
- implicit Euler scheme, 113
 - region of stability, 114
- implicit Milstein scheme, 390
- log-Euler scheme, 112, 113
- lookback option, 125
- low-discrepancy sequence, *see* Monte Carlo, random number generation, quasi-random
- lower bound for American option, 34, 35, 164
- parametric, 159, 161
- regression-based, 161
- mean-square error, 123
- Milstein scheme, 119, 121
 - multi-dimensional, 121
- modified trapezoidal scheme, *see* Monte Carlo, Heun scheme
- optimal root-mean-square error, 123
- perfect foresight bias, *see* American/Bermudan option, perfect foresight bias
- predictor-corrector, 115, 116
 - convergence order, 116
- random number generation, 97
 - acceptance-rejection method, 99–101
 - Box-Muller method for Gaussian distribution, 99
 - composition method, 101–102
 - conditional Gaussian, 1066
 - correlated Gaussian, 103
 - correlated Gaussian by Cholesky decomposition, 103
 - correlated Gaussian by eigenvalue decomposition, 104
 - inverse transform method, 98
 - linear congruential generator, 97
 - Marsaglia polar method for Gaussian distribution, 99
 - Mersenne twister, 98
 - period, 98
 - pseudo-random, 97, 130
 - quasi-random, 129

- Sobol, 129
 region of stability, 110
 Richardson extrapolation, 122, 468
 sample mean, 96
 sampling extremes, 124–128
 adjusting barrier for sampling frequency, 128, 937, 970
 with Brownian bridge, 125
 SDE discretization, 108
 second-order scheme, 119, 121
 seed, 97
 standard error, 97, 122
 for digital option, 133
 for greeks, 132, 135
 strong convergence order, 111
 strong law of large numbers, 96
 strongly consistent, 109
 third-order scheme, 468
 upper bound for American option, 34–36, 163, 164
 variance reduction, *see* variance reduction
 weak convergence, 109
 weak convergence order, 110
 weakly consistent, 109
 most likely path, *see* volatility, implied,
 most likely path approximation
 multi-rate vanilla derivative, 763–813
 copula method, *see* copula method
 Libor market model, 807
 observation lag, 782
 stochastic volatility, *see* multi-stochastic volatility model
 term structure models, 806
 multi-stochastic volatility model, 800–806, 1149
 correlation impact, 803
 measure change by CMS caplet calibration, 802
 measure change by drift adjustment, 801
 Monte Carlo
 Quadratic-Exponential scheme, 803
 multi-rate vanilla derivative, 800–806
 multi-tranche, *see* CLE, multi-tranche
 non-central chi-square distribution, 284
 asymptotics, 392
 CDF, 102, 319
 in CEV model, 283
 in delta-gamma VaR/cVaR, 998
 in LS model, 517
 two-dimensional, 517
 Normal model, XXXVIII, 283
 CMS spread, 774
 vega to swaptions, 775
 numeraire, 10, 171
 change of numeraire, 12
 Girsanov's theorem, *see* Girsanov's theorem
 discrete money market account, XXXVIII, 175
 money market account, XXXVIII, 22, 28, 172
 OIS, *see* overnight index swap
 one-dimensional integral for spread option, 787
 operator calculus, 998–999
 OTC market, *see* over-the-counter market
 out-of-model adjustment, 951–971
 adjusters method, 954–956
 algorithm, 955
 as control variate, 955
 volatility adjustment, 956
 asset-based adjustment, 963–964
 CMS spread option, 964
 coupon calibration, 952–954
 delta-adjustment method, 956
 extended calibration, 953
 fee adjustment method, 967–969
 additive, 968
 blended, 968
 impact on derivatives, 968
 multiplicative, 968
 issues, 961, 964
 mapping function adjustment, 965
 market adjustment, 965
 path re-weighting method, 956–961
 as control variate, 961
 Boltzman-Gibbs distribution, 959
 Boltzman-Gibbs weights, 959
 dual, 961

- inappropriate use, 958
- partition function, 958
- risk sensitivities, 961
- PDE for coupon values, 953
- proxy model method, 961
- spread adjustment method, 966
- strike adjustment method, 969–971
 - impact on derivatives, 970
- over-the-counter market, 193
- overhedge, 1023
- overlay curve, *see* yield curve, overlay curve
- overnight index swap, 193, 200, 266

- P&L, 696, 991–995
- P&L analysis, 986
- P&L attribution, *see* P&L explain
- P&L explain, 993–995
 - bump-and-do-not-reset explain, *see* P&L explain, waterfall explain
 - bump-and-reset explain, 994–995
 - waterfall explain, 993–994
- P&L explanation, *see* P&L explain
- P&L of hedged book, 987–990
- P&L predict, *see* P&L prediction analysis
- P&L prediction analysis, 258, 991–993
 - first-order, 991
 - second-order, 991
 - unpredicted P&L, 991
- par rate, *see* forward swap rate
- parameter averaging, *see* calibration, time averaging
- partial differential equation, *see* PDE
- partition function, 958
- pathwise delta approximation, *see* pathwise differentiation method, 135–139, 1035–1060
- adjoint method, 1056
 - computational complexity, 1053, 1057
- barrier option, *see* barrier option, pathwise differentiation method
- Bermudan swaption, *see* Bermudan swaption greeks, pathwise differentiation method

- CLE, *see* CLE greeks, pathwise differentiation method
- computational complexity, 1052, 1053
- discontinuous payoff, 1042, 1061
- European option, 1054
- gamma, 1050, 1056
- generalized trigger product, *see* generalized trigger product, pathwise differentiation method
- Libor market model, *see* Libor market model, pathwise differentiation method
- money market account, 1046
- Monte Carlo models, 1051–1060
- pathwise delta approximation, 1059
- PDE models, 1044–1050
- sensitivity path generation, 138–139
- TARN, *see* TARN, pathwise differentiation method
- vega, 1050, 1056
- payoff smoothing, 1001–1034
 - adaptive integration, 1006
 - adding singularity to grid, 78, 1007
 - barrier option, 1074–1077
 - benefits, 1012
- Bermudan swaption, *see* CLE greeks, tube Monte Carlo
- box smoothing, 1015–1018
 - multiple dimensions, 1020
 - on discrete grid, 1015
- by importance sampling, 1065–1077
- CLE, *see* CLE greeks, tube Monte Carlo
- continuity correction, 59, 1012
- fuzzy logic, 1028
- gamma, 1019
- grid shifting, 1007
- hat smoothing, 1019
- integration, 1012
- Lagrange basis, 59, 1019
- locality, 1019
- Monte Carlo, 1022–1030
- moving average, 1012, 1013
 - choice of window, 1014
- multiple dimensions, 1019–1022
 - box smoothing, 1020
 - dominant dimension, 1022
- one dimension, 1014

partial analytical integration, 76–78, 1010
 partial coupons, 1028
 PDE, 1012
 piecewise smooth function on a grid, 1016
 singularity removal, 1009
 TARN, *see* TARN, payoff smoothing; TARN, tube Monte Carlo
 tube Monte Carlo, *see* tube Monte Carlo
 PCA, *see* Principal Components Analysis
 PDE, 18, 43–93
 A-stable scheme, 55
 ADI scheme, 43, 82–85
 boundary conditions, 85
 Asian option, 70
 backward induction, 51
 Black-Scholes, *see* Black model, PDE
 boundary conditions
 for barrier options, 64
 from PDE itself, 385, 424
 linear at boundary, 48
 log-linear at boundary, 48
 Cauchy problem, 18, 44
 centering, 561
 conditional stability, 55
 consistent scheme, 56
 convection-dominated, 61–64
 convergent scheme, 56
 coupon-paying, 67
 Craig-Sneyd scheme, *see* PDE, predictor-corrector scheme
 Crank-Nicolson scheme, 50
 American options, 69
 not strongly A-stable, 55
 oscillations, 55, 58
 Dirichlet problem, 44, 64
 space discretization, 46
 dividends, 67, 68
 domain truncation, 44
 stability of greeks, 1002
 Douglas-Rachford scheme, 85, 91
 boundary conditions, 85
 early exercise, 69
 exponentially fitted schemes, 63

extra state variable method, *see* PDE, path-dependent options
 for implied volatility, *see* volatility, implied, PDE for
 forward equation, *see* Kolmogorov forward equation
 fully implicit scheme, 50
 greeks off grid, 1005
 L-stable scheme, 55
 Lagrange basis, 58, 59
 Lax equivalence theorem, 56
 local volatility model, *see* local volatility model, PDE
 mesh refinement, 73, 79
 equidistant blocks, 74
 non-equidistant, 75
 multi-dimensional, 92
 multi-exercise, 71
 multi-level time-stepping, 58
 non-equidistant discretization, 56
 Nyquist frequency, 59
 odd-even effect, 59
 operator splitting, 82
 orthogonalization, 86
 drawbacks, 88
 partial analytical integration,
 see payoff smoothing, partial analytical integration
 path-dependent options, 69, 71, 868, 870, 896, 899, 932, 934
 Peaceman-Rachford scheme, 84
 boundary conditions, 85
 predictor-corrector scheme, 89–92
 quantization error, 59
 Rannacher stepping, 58–61, 67, 457
 semi-Lagrangian methods, 64
 Shannon Sampling Theorem, 59
 similarity reduction, 71
 sinh transform, 384
 smoothing, 58–61
 continuity correction, 59
 grid dimensioning, 1002
 grid shifting, 60, 1002
 space discretization, 45
 stable scheme, 53
 strongly A-stable scheme, 55
 time discretization, 49
 theta scheme, 50
 two-dimensional, 80

- two-dimensional with mixed derivatives, 86, 89
 upwinding, 62
 variable transform, 44
 von Neumann method, 53–56
 amplification factor, 54
 stability criterion, 54
 well-posed, 56
 Poisson distribution, 102
 portfolio replication, *see* Bermudan swaption greeks, portfolio replication for hedging power Gaussian copula, *see* copula, power Gaussian predictor-corrector, 89, 115, 382, 641
 Monte Carlo, *see* Monte Carlo, predictor-corrector PDE, *see* PDE, predictor-corrector scheme present value of a basis point, *see* swap, annuity principal component, 105 Principal Components Analysis, 105, 106, 498, 602–604 principal factor, 105 product integral, 484 Profit-And-Loss, *see* P&L pseudo-Gaussian model, *see* quasi-Gaussian model pseudo-random number generator, *see* Monte Carlo, random number generation, pseudo-random put-call parity, 24 PVBP, *see* swap, annuity QG model, *see* quadratic Gaussian model qG model, *see* quasi-Gaussian model quadratic covariation, XXXVII, 7 quadratic Gaussian model, 441, 518–533
 as affine model, 519
 benchmark rate parameterization, 525 Bermudan swaption, *see* Bermudan swaption calibration, local projection method, quadratic Gaussian model bond dynamics, 521
 bond reconstitution formula, 520 calibration, 531–532
 multi-pass bootstrap, 531 CLE, *see* CLE calibration, local projection method, quadratic Gaussian model CMS spread option, *see* CMS spread option, quadratic Gaussian model curve factor, 523 European swaption, 526–531
 approximations, 528
 exact, 527
 Fourier integration, 529
 rank-2 approximation, 530 Fourier integration, 530 mean-reverting state variables, 519 moment-generating function, 529 Monte Carlo, 533 one-factor, 441 parameterization, 523–526 PDE, 533 quadratic approximation to swap rate, 529 short rate, 519 short rate in SV form, 525 short rate state distribution
 in annuity measure, 526
 in forward measure, 521 short rate state dynamics, 441, 519
 in forward measure, 521
 in annuity measure, 526 smile generation, 523–524 spanned stochastic volatility, 523, 532 TARN, *see* TARN, local projection method, quadratic Gaussian model volatility factor, 523 volatility smile, 531 volatility swap, *see* volatility swap, quadratic Gaussian model quadratic variation, XXXVII, 7 quadratic volatility model, 287–291
 European call option value, 290
 European put option value, 290, 291 Markovian projection, 1137 measure change, 289 small-noise expansion, 308

smile dynamics, 350
strict supermartingale, 288
time-dependent, 308

Quadratic-Exponential scheme, *see*
square-root process, Monte Carlo,
Quadratic-Exponential scheme
multi-dimensional, *see* multi-
stochastic volatility model,
Monte Carlo, Quadratic-
Exponential scheme
quadrature, 531, 786
 Gauss-Hermite, 531, 787
 Gauss-Legendre, 786
 Gauss-Lobatto, 786
quanto CMS, 744–748
 annuity mapping function, 748
 convexity adjustment, 747–748
 copula method, 747
 quanto adjustment, 745
 replication method, 746
quasi-Gaussian model, 537–587
 Bermudan swaption, *see* Bermudan
 swaption calibration, local pro-
 jection method, quasi-Gaussian
 model
 bond reconstitution formula, 538
 calibration, 581
 CEV local volatility, 545
 CLE, *see* CLE calibration, local pro-
 jection method, quasi-Gaussian
 model
 CMS convexity adjustment, *see*
 CMS, convexity adjustment,
 quasi-Gaussian model
 density approximation, 583
 direct integration, 558, 583
 Libor-with-delay, *see* Libor-with-
 delay, quasi-Gaussian model
 linear local volatility, 545–548
 calibration, 548
 European swaption, 547
 for swaption strip, 547
 swap rate dynamics, 546
 swap rate inter-temporal correla-
 tion, 555
 swap rate variance ratio, 553
 Markovian projection, 541, 564, 577,
 1139
 mean reversion, *see* mean reversion

mean reversion calibration, *see*
 mean reversion calibration
Monte Carlo, 563
 Euler scheme, 563
multi-factor, 572–583
 benchmark rate correlations, 582
 benchmark rate parameterization,
 574
 bond reconstitution formula, 574
 calibration to spread options, 582
 correlation smile, 582
 loadings, 582
 local volatility, 574
 Monte Carlo, 583
 PDE, 582
 short rate state distribution in
 annuity measure, 577
 short rate state dynamics, 573
 stochastic volatility, 574–583
 swap rate dynamics, 576–581
 swap rate dynamics by Markovian
 projection, 577
one-factor local volatility, 539
 short rate state dynamics, 539
PDE, 560–563
 convection-dominated, 561
 domain truncation, 562
 space discretization, 561
short rate state distribution, 559
short rate state dynamics, 538
 in annuity measure, 542, 543
 in forward measure, 583
single-state approximation, 563–567
small-time asymptotics, 559
stochastic volatility, 567–572
 bond reconstitution formula, 568
 calibration, 570–571
 Monte Carlo, 572
 non-zero correlation, 572
 PDE, 572
 swap rate dynamics, 568–570
 unspanned, 568
swap rate dynamics, 540–545, 549
 approximate, 541–545
 approximate linear, 542
 approximate quadratic, 545
swap rate variance, 544
swap rate volatility, 540

- TARN, *see* TARN, local projection
 method, quasi-Gaussian model
 volatility swap, *see* volatility swap,
 quasi-Gaussian model
- Radon-Nikodym derivative, 9, 1067
 range accrual, 211
 CMS, 211
 CMS spread, 211, 764
 curve cap, 212, 764
 dual, 212, 764
 floating, 764
 product-of-ranges, 212
 ratchet swap, 212
 relative entropy, 957
 replication method, 337, 722
 CMS, *see* CMS, convexity adjustment, replication method
 European option, *see* European-style option, replication method
 Libor-in-arrears, *see* Libor-in-arrears, replication method
 Libor-with-delay, *see* Libor-with-delay, replication method
 semi-static, 939
 reserve, 986
 rho, 980
 Riccati, 364
 Riemann zeta function, 128
 risk limit, 986
 risk measure, 996
 coherent, 996
 risk sensitivities, 1093
 common definitions, 980
 delta, *see* delta
 grid dimensioning for stability, 1002
 grid shifting for stability, 1002
 Jacobian method, 254–258, 985, 986,
 1105, 1106, 1111, 1118, 1119,
 1121
 off PDE grid, 1005
 perturbation approach, 1050
 vega, *see* vega
 root search, 99
 Newton-Raphson method, 99, 116,
 235
 secant method, 235
 Runge-Kutta method, 116, 365, 432,
 434, 514
 running maximum, 124
 running minimum, 124
 SABR model, 343–345, 357, 951, 1121
 ad-hoc improvements, 703
 density tail, 760
 moment explosion, 344
 volatility smile expansion, 345
 SALI tree, *see* tree, SALI
 sausage Monte Carlo, *see* tube Monte Carlo
 SDE, *see* diffusion, SDE
 SDE discretization, *see* Monte Carlo,
 SDE discretization
 Sharpe ratio, 22
 shifted log-normal model, *see* displaced
 log-normal model
 short rate, 169
 short rate model, 172
 affine, *see* affine short rate model
 affine one-factor, *see* affine short
 rate model, one-factor
 Black-Derman-Toy, *see* Black-Derman-Toy model
 calibration to yield curve, 455
 forward induction, 456
 forward-from-backward induction,
 458
 Cox-Ingersoll-Ross, *see* Cox-Ingersoll-Ross model
 Dybvig parameterization, 461–463,
 466
 HJM representation, 462
 econometric, 449
 empirical estimation, 449
 forward volatility impact on
 Bermudan swaption, 876
 Gaussian approximation, 1064
 Gaussian model for basis spread,
 681
 Gaussian short rate, *see* Gaussian
 short rate model
 Hagan and Woodward parameterization, 463–466
 Ho-Lee, *see* Ho-Lee model
 importance sampling, 1063–1065
 log-normal, 443–449
 issues, 445

- Sandmann-Sondermann transform, 446
 Monte Carlo, 467–469
 Euler scheme, 467
 Milstein scheme, 467
 payoff construction issues, 468
 SDE discretization, 467
 variance reduction, 468
 multi-factor, 477
 path independence, 444
 PDE, 454–455
 domain truncation, 454
 power-type, 449
 quadratic Gaussian, *see* quadratic Gaussian model
 quasi-Gaussian, *see* quasi-Gaussian model
 time-stationary, 416
 volatility calibration, 459–461
 multi-pass bootstrap, 461
 shout option, 935
 on capped coupon, 935
 optimal stopping time, 936
 similarity reduction, 71, 869
 CLE, *see* CLE valuation, PDE, similarity reduction
 PDE, *see* PDE, similarity reduction
 single-rate vanilla derivative, 695–762
 approximately single-rate, 707
 cap, *see* cap
 CMS cap, *see* CMS cap
 CMS floor, *see* CMS floor
 CMS swap, *see* CMS swap
 ED future, *see* ED future
 European swaption, *see* European swaption
 futures contract, *see* ED future
 Libor-in-arrears, *see* Libor-in-arrears
 Libor-with-delay, *see* Libor-with-delay
 range accrual, *see* range accrual
 singular value, 860
 singular value decomposition, *see* CLE regression, SVD decomposition
 truncated, *see* CLE regression, truncated SVD decomposition
 singularity removal, *see* payoff smoothing, singularity removal
 skew vega, *see* vega, skew vega
 smile vega, *see* vega, smile vega
 snowball, *see* CLE, snowball
 snowbear, 213
 snowrange, 213
 snowstorm, 213
 Sonia, 193, 200
 spline, 230, 270–275
 Catmull-Rom, 238, 240, 271, 272
 cubic C^2 , 273–274
 cubic smoothing, 248
 exponential tension spline, 243
 Hermite cubic, 238, 270–273
 interpolating, 248
 Kochanek-Bartels, 272
 least-squares regression, 248
 natural, 241
 natural cubic, 273
 shape preserving, 275
 smoothing, 234
 TCB, *see* spline, Kochanek-Bartels tension, 240, 243, 244, 246, 247, 250, 272, 274–275
 convergence to piecewise linear, 275
 tension factor, 243
 spot Libor measure, *see* measure, spot
 spot rate, *see* short rate
 square-root process, 315
 $E(\sqrt{z})$, 1153, 1155
 basic properties, 318–320
 boundary behavior, 319
 conditional CDF, 319
 conditional moments, 319
 Feller condition, 319
 moment-generating function, 322, 342, 364, 372
 time-dependent parameters, 364
 moments, 375
 Monte Carlo, 388–394
 Euler scheme, 389
 exact simulation, 388
 full truncation scheme, 389
 higher-order schemes, 389
 log-normal approximation, 390
 moment-matching schemes, 390
 Quadratic-Exponential scheme, 392, 394
 truncated Gaussian scheme, 391
 multi-dimensional, 1152

- PDF, 1153, 1156
 - stationary distribution, 320, 383
- static replication, 210, 717
 - CMS, *see* CMS, convexity adjustment, replication method
 - European option, *see* European-style option, replication method
 - Libor-in-arrears, *see* Libor-in-arrears, replication method
 - Libor-with-delay, *see* Libor-with-delay, replication method
- stochastic optimization method, 953
- stochastic volatility model, 315–402,
 - 569, 570, 1140
 - as interpolation rule, 701
 - ATM volatility, 348
 - basket option, *see* Markovian projection, basket option in SV model
 - calibration, 701–702
 - calibration norm, 702
 - normalization, 702
 - caplet calibration, 705
 - CEV type, *see* SABR model
 - CMS convexity adjustment, 738
 - correlation, 347
 - dampening constant, 325
 - delta, 697
 - effective skew, 373
 - effective volatility, 371, 372
 - effective volatility of variance, 375
 - European option, 327
 - control variate, 328
 - volatility mixing, 339
 - explicit solution, 320
 - for CMS rate, 738–742
 - dynamics in forward measure, 739
 - Fourier integration, 324–339
 - arbitrary European payoffs, 336, 338
 - convolution, 325
 - direct integration, 330
 - discrete, 330
 - FFT, 330
 - for variance, 339–343
 - integration bounds, 330
 - strip of convergence, 329
 - with control variate, 328, 330
- hedging, 353–358
- level parameter, 317
- link between forward and annuity measures, 739
- LSV, *see* local stochastic volatility model
- martingale property, 320
- mean reversion speed, 316, 317, 348
 - half-life, 318
- measure change, 322
- moment explosion, 323
- moment-generating function, 321, 324, 327
 - branch cut, 330
 - singularities, 329
 - time-dependent parameters, 364
- Monte Carlo, 387–397
 - Broadie-Kaya scheme, 394
 - Broadie-Kaya simplified scheme, 396
 - exact scheme, 394
 - martingale correction, 397
 - Taylor-type schemes, 396
- variance process, *see* square-root process, Monte Carlo
- multi-dimensional, *see* multi-stochastic volatility model
- PDE, 381–387
 - boundary conditions for stochastic variance, 385
 - boundary conditions from PDE itself, 385
 - discretizing spot, 387
 - discretizing stochastic variance, 383
 - for forward Kolmogorov equation, 386
 - predictor-corrector, 382
 - quadratic discretization, 384
 - range for spot, 386
 - range for stochastic variance, 382
 - sinh transform, *see* PDE, sinh transform
 - sinh-quadratic discretization, 384
 - variable transform, 383, 384
- process for variance, *see* square-root process
- skew, 317, 346
- smile dynamics, 347–349, 351, 353, 354

- SV volatility, 317
 time-dependent, 363–402
 asymptotic expansion, 366–370
 averaging, *see* calibration, time averaging
 Fourier integration, 363, 366
 volatility of variance, 316, 317, 346
 volatility of volatility, 318
 stopping time, 6
 straddle, 223
 strategy, 7
 doubling, 10
 gains process, 8
 permissible, 9
 replicating, 11
 self-financing, 8, 17
 Stratonovich integral, 5
 strike price, 24
 structured note, *see* exotic swap
 structured swap, *see* exotic swap
 Student's *t*-distribution, 101
 Monte Carlo, 101
 survival measure
 Bermudan swaption, *see* Bermudan swaption, survival measure
 importance sampling, *see* importance sampling, survival measure
 SV model, *see* stochastic volatility model
 SV model with general variance process, 359–361
 martingale properties, 360
 moment explosion, 361
 properties, 359
 stationary distribution, 360
 strict supermartingale, 360
 SVD, *see* CLE regression, SVD decomposition
 SVI model, *see* volatility smile, SVI
 swap, 197
 accreting, 200
 amortizing, 200
 annuity, XXXVIII, 199
 annuity factor, 170
 averaging, *see* averaging cash flow
 cash-settled, 744
 CMS, *see* CMS swap
 effective date, 225
 fixed-floating, 198, 199, 230, 231
 valuation formula, 199
 fixing dates, 198
 legs, 197
 Libor-in-arrears, *see* Libor-in-arrears
 Libor-with-delay, *see* Libor-with-delay
 par rate, *see* forward swap rate
 payer, 203
 payment dates, 198
 receiver, 203
 swap rate, *see* forward swap rate
 swap market model, 617, 675–677
 swap measure, *see* measure, annuity
 swap rate, *see* forward swap rate
 swaption grid, *see* European swaption, swaption grid
 Tanaka extension of Ito's lemma, 7, 26, 294, 1131
 targeted redemption note, *see* TARN
 TARN, 217, 218, 925–933
 cap at trigger, 219
 global model, 927
 impact of inter-temporal correlation, *see* inter-temporal correlation, impact on TARNs
 importance sampling, 1068–1077
 one-step survival conditioning, 1069
 removing first digital, 1068
 leverage, 927
 Libor market model, 927
 lifetime cap, *see* TARN, cap at trigger
 lifetime floor, *see* TARN, make whole
 local projection method, 928–931
 Gaussian short rate model, 929
 Markov-functional model, 931
 quadratic Gaussian model, 931
 quasi-Gaussian model, 931
 make whole, 219
 Markov-functional model, 473
 multi-factor quasi-Gaussian model, 927
 partial analytical integration, 1011
 pathwise differentiation method, 1044

- payoff smoothing, 1011, 1029, 1068–1077
 PDE, 931–933
 cap at trigger, 933
 make whole, 933
 Monte Carlo pre-simulation, 933
 upper bound for extra state variable, 932
 tube Monte Carlo, 1029
 valuation formula, 218
 volatility smile, 927, 929–931
 tenor structure, XXXVIII, 170
 index function, 591
 tension spline, *see* spline, tension
 term parameters, 378
 term structure model, 202, 277
 terminal swap rate model, 707–714
 annuity mapping function, 708, 713, 722, 724–725, 728, 730, 732
 as conditional expected value, 724–725
 calibration to market, 728
 forward swap rate condition, 733
 forward value condition, 732
 in measure change, 735
 linear approximation, 728
 LM model, *see* Libor market model, annuity mapping function
 mean reversion, *see* CMS, convexity adjustment, impact of mean reversion
 multi-rate, 765
 swap rate squared condition, 733
 CMS convexity adjustment, *see* CMS, convexity adjustment, linear TSR model
 consistency condition, 708
 exponential TSR model, 712–713
 Libor-with-delay, *see* Libor-with-delay, swap-yield TSR model
 linear TSR model, 709
 CMS convexity adjustment, *see* CMS, convexity adjustment, linear TSR model
 forward CMS straddle, 941
 mean reversion parameterization, 710
 swap rate distribution in forward measure, 736, 737
 vega hedging, 712
 loading from Gaussian model, 712
 no-arbitrage condition, 708
 PDF of swap rate in forward measure, 737
 from CMS caplets, 737
 reasonableness, 708
 swap rate distribution in forward measure, 736
 swap-yield TSR model, 713–714
 CMS convexity adjustment, *see* CMS, convexity adjustment, swap-yield TSR model
 theta, 980, 992
 rolling yield curve, 992
 Tikhonov regularization, *see* CLE regression, Tikhonov regularization
 time decay, 52
 time value, 27
 “tip-top”, *see* “flip-flop”
 tower rule, *see* iterated conditional expectations
 tree, 423
 binomial, 444, 456
 SALI, 78
 trinomial, 51, 456
 truncated Gaussian scheme, *see* square-root process, Monte Carlo, truncated Gaussian scheme
 TSR model, *see* terminal swap rate model
 tube Monte Carlo, 1022–1030
 barrier option, *see* barrier option, tube Monte Carlo
 Bermudan swaption, *see* CLE greeks, tube Monte Carlo
 CLE, *see* CLE greeks, tube Monte Carlo
 digital option, 1024
 discrete knock-in barrier, 1028
 generalized trigger product, *see* barrier option, tube Monte Carlo
 partial coupons, 1028
 TARN, *see* TARN, tube Monte Carlo
 underhedge, 1023
 uniform distribution, XXXVII, 768
 universal law of volatility, 1137

upwinding, *see* PDE, upwinding

value-at-risk, 499, 975, 996–998
 conditional, 996
 delta VaR, 998
 delta-gamma VaR/cVaR, 998
 Gaussian, 997
 historical, 996

vanilla derivative, 695–813
 multi-rate, *see* multi-rate vanilla derivative
 single-rate, *see* single-rate vanilla derivative

vanilla model, 202, 277, 315, 1121, 1129
 for multi-rate derivative, *see* multi-rate vanilla derivative
 for single-rate derivative, *see* single-rate vanilla derivative

local volatility model, *see* local volatility model
 stochastic volatility model, *see* stochastic volatility model

vanna, 980

VaR, *see* value-at-risk

variance reduction, 143–158
 antithetic variates, 144
 efficiency, 145
 non-Gaussian, 145

common random number scheme, 132, 134

conditional Monte Carlo, 127

control variate, *see* control variate from hedging strategy, *see* control variate, dynamic

importance sampling, *see* importance sampling

moment matching, 146

systematic sampling, 145

Vasicek model, 411–413
 bond reconstitution formula, 412
 bond volatility, 413
 forward rate volatility, 413
 short rate distribution, 411
 short rate dynamics, 411
 yield curve shapes, 412

vega, 355, 980, 1095–1125
 additivity, 1103

Bermudan swaption, 1114

bucketed shocks, 1099

CMS spread option, 1116, 1120
 constant Libor correlations, 1120
 constant Libor correlations, 1115, 1120

constant term swap correlations, 1116, 1118–1120

cumulative shocks, 1099

direct method, 1098–1102, 1110
 Bermudan swaption, 1103
 European swaption, 1102
 second-order effects, 1111

European swaption, 1113, 1114

flat shock, 1099

forward swaption straddle, 948
 “good”, 1102–1105

hybrid method, 1111–1113
 algorithm, 1112
 Bermudan swaption, 1114
 CMS spread option, 1116
 European swaption, 1113, 1114

in LM model
 coverage, 884

indirect method, 1105–1111, 1121
 Bermudan swaption, 1109
 European swaption, 1108
 least-squares problem, 1106
 locality, 1107
 smoothing, 1107

Jacobian method, *see* vega, indirect method; risk sensitivities, Jacobian method

Libor market model, 1095–1125
 bootstrap calibration, 1111, 1112
 multi-factor, 1115
 projection, 1123

local projection method, 867

local vs. global, 1097

locality, 1104
 benchmark set locality, 1104
 exotic locality, 1104
 full set locality, 1104

market vega, 984, 1096, 1110

model vega, 984, 1096, 1124–1125

pathwise differentiation method, *see* pathwise differentiation method, vega

projection, 1122–1124

relationship to gamma, 981

- row shocks, 1099
- running cumulative shocks, 1099
- scaling, 1103
- skew vega, 1113–1115
- smile vega, 1113–1115
- volatility, 27
 - average convexity, 307
 - Bachelier, *see* volatility, Normal
 - basis point, *see* volatility, Normal
 - Black, XXXVIII, 204
 - bp, *see* volatility, Normal
 - CEV, 280, 623
 - Dupire's, *see* Dupire local volatility
 - factor volatility, 499
 - forward volatility of Libor rate, 817
 - Gaussian, *see* volatility, Normal
 - implied, 278
 - as average of realized, 989
 - effects of mis-specification, 987
 - most likely path approximation, 990
 - PDE for, 296
 - local, *see* Dupire local volatility
 - Normal, 204, 283, 623
 - Normal for CMS spread option, 774
 - separable, 300
 - small-noise expansion, 307
 - spanned stochastic volatility, 452
 - spot volatility, 817
 - spread, 774
 - strike-dependent, 775
 - stochastic, *see* stochastic volatility model
 - unspanned stochastic volatility, 443
 - "volatility squeeze", 422
- volatility cube, *see* European swaption, volatility cube
- volatility derivative, *see* forward volatility derivative
- volatility skew, 279
- volatility smile, 279, 315
 - ATM backbone, 699, 700
 - backbone, 696
 - adjustable, 697–700
 - curvature, 1138
 - dynamics, 279, 348, 696–700, 818
 - sticky delta, 350, 697
 - sticky strike, 352, 697
 - forward skew, 944
- Gaussian backbone, 698
- impact on forward volatilities, *see* forward volatility, impact of volatility smile
- impact on inter-temporal correlation, *see* inter-temporal correlation, impact of volatility smile
- probability density from, 278
- SABR, *see* SABR model
- shadow delta hedging, 697
- skew vega, 1114
- skew-dominated, 352
- slope, 279
- smile vega, 1114
- SVI, 703, 951, 1121
- upward sloping, 281
- vega, 1114
- volatility structure, 815
- volatility swap, 220, 221, 933–945
 - capped, 937
 - CMS spread, 221
 - copula method, *see* copula method, volatility swap
 - fixed-expiry, 221, 940
 - fixed-tenor, 221, 940
 - impact of forward volatility, 944
 - impact of volatility smile dynamics, 941
 - Libor market model, 933, 934
 - local projection method, 934
 - min-max, *see* min-max volatility swap
 - PDE, 934
 - quadratic Gaussian model, 941
 - quasi-Gaussian model, 941
 - with barrier, 222
 - with shout, 221, 935
 - volga, 980
 - Volterra integral equation, 436
 - vomma, 980
- Wiener process, *see* Brownian motion
- year fraction, 224
- yield curve, 191, 230, 231, 233
 - base index curve, 268
 - basis risk, 270
 - benchmark set, 230

- forecasting curve, *see* yield curve, index curve
- index curve, 261, 267, 677
- index-discounting basis, 197, 261
- instantaneous forward curve, 233
- joint evolution of discount and forward curves, 677
- multi-index curve group, 267–270
- overlay curve, 259
- perturbation locality, 230, 251–253, 258
- Principal Components Analysis, *see* Principal Components Analysis
- ringing, 235, 242, 243, 252
- smooth, 258
- spread curve, 269, 884
- tenor basis, 230, 267
- TOY effect, 258
- yield curve construction, 229–275
 - benchmark set, 231
 - bootstrapping, 234
 - flat forward, 236
 - linear yield, 235
 - constrained optimization, 248
 - cross-currency, 259
 - cross-currency arbitrage, 260
 - cubic spline C^2 , 240–243
 - problems, 242
- curve overlays, 258
- FX forwards, 259
- Hermite spline, 238–240
 - iterative solution, 239
- Jacobian rebuild, 256
- multi-index curve group, 230, 265
- non-parametric fitting, 245–250
 - norm specification, 245
 - optimization algorithm, 245
- separate discount and forward curves, 260
- spline, *see* spline
- spline fitting, 234–244
- tension spline, 243–244
- yield curve risk, 250–258
 - cumulative shifts, 256, 257
 - forward rate approach, 252
 - Jacobian method, *see* risk sensitivities, Jacobian method
 - par-point approach, 251
 - rolling for theta, 992
 - waterfall approach, *see* yield curve risk, cumulative shifts
- yield curve spread option, *see* CMS spread option
- zero-coupon bond, *see* discount bond
- zero-coupon bond option, 185