

## Master's Thesis in Mathematical Finance

Youssef Raad

KU-id: zfw568

## The Black-Scholes Option Pricing Model A Markov-Switching Extension

Date: 22-12-2025

Supervisor: Rolf Poulsen

Institute: Institute for Mathematical Science  
Department Department of Mathematical Sciences  
Author(s): Youssef Mahmoud Raad  
Title and subtitle: The Black-Scholes Option Pricing Model: A Markov-Switching Extension  
Description: This Master's thesis explores extensions to the Black-Scholes option pricing model as a frequent critique of said model is its lack of adaptivity to economical turbulent environments. The extensions are via continuous and discrete state-space models often called Markov regime-switching models. The analysis is done using S&P 500 data ( $\hat{GSPC}$  and  $\hat{TR}$ ).  
Advisor: Rolf Poulsen  
Date: 9<sup>th</sup> January 2026

### Acknowledgements

A special thanks to Théo Michelot, Assistant Professor of Statistics at Dalhousie University, for his help throughout the process answering, at time, extremely trivial questions.

I would like to thank Rolf Poulsen, Professor of Mathematical Finance at the University of Copenhagen, for his supervision during the preparation project and this thesis. Thank you for the opportunity and for the chance to collaborate as a teaching assistant.

Above all and without comparison, my deepest thanks go to my mother. I would not have accomplished anything without her steady belief, patient encouragement and constant support. She has been the light at the end of every bleak day, my reason to continue and any success I have rests on her sacrifices and her power of will. Nothing I do will ever fully repay that debt—and I would not have it any other way.

This one is for you.

### Note

We only provide proofs which are not well known across both fields of statistics and mathematical finance and/or give rise to meaningful matters of learning. Furthermore, we will list definitions if they are rarely encountered upon or subject to many varying forms of definitions. Results and definitions that we use to prove main results will be listed in [Appendix A.2](#).

## Abstract

The Black-Scholes model (BSM) has long been a cornerstone of financial theory; however, its assumption of constant drift and volatility fails to capture the time-varying nature of asset returns, such as volatility clustering and abrupt regime shifts. This thesis investigates whether extending the BSM to include dynamic parameter evolution improves out-of-sample forecasting performance. Two extensions are proposed and implemented: a discrete Black-Scholes Hidden Markov Model (BS-HMM) and a continuous Black-Scholes State-Space Model (BS-SSM).

Using daily S&P 500 data from 1927 to 2025, the models are calibrated via Maximum Likelihood Estimation. In-sample analysis reveals that the extended models, particularly a 4-state BS-HMM and a factor-loaded continuous state-space model ( $BS-SSM_\beta$ ), provide a superior fit to historical data compared to the static BSM. These models successfully identify distinct market phases, distinguishing between tranquil “bull” markets and high-volatility “crisis” regimes, such as the 1929 Crash and the 2008 Financial Crisis.

Despite the richer descriptive power and improved in-sample fit, out-of-sample evaluation on a hold-out period (2020–2025) indicates that the regime-switching extensions offer negligible gains in one-step-ahead point forecasting accuracy (MSE and RMSE) relative to the constant-parameter BSM. While the extended models produce more realistic, horizon-dependent forecast densities, the findings suggest that the added complexity of latent state inference does not translate into superior short-term predictive power for point forecasts.

## CONTENTS

1	Data (I of II) . . . . .	5
2	Theory & Methodology . . . . .	8
2.1	The Black–Scholes Model . . . . .	8
2.1.1	Likelihood Formulation & Parameter Estimation . . . . .	9
2.1.2	Simulation . . . . .	10
2.2	Hidden Markov Models . . . . .	12
2.2.1	Independent Mixture Models . . . . .	12
2.2.2	Markov Chains & Hidden Markov Models . . . . .	14
2.2.3	State-Dependent Distributions . . . . .	20
2.2.4	Likelihood Formulation & Parameter Estimation . . . . .	21
2.2.5	Standard Errors & Confidence Intervals . . . . .	28
2.2.6	Forecasting, Decoding and State Prediction . . . . .	30
2.2.7	Number of States . . . . .	34
2.2.8	Simulation . . . . .	36
2.3	Continuous State-Space Models . . . . .	37
2.3.1	Autoregressive Processes . . . . .	38
2.3.2	Likelihood Formulation & Parameter Estimation . . . . .	46
2.4	Model Selection Criteria & Assessment . . . . .	52
2.4.1	Information Criteria: AIC & BIC . . . . .	52
2.4.2	Pseudo-Residuals . . . . .	53
3	Data (II of II) . . . . .	57
4	Empirical Data Application . . . . .	62
4.1	Model Selection & Assessment . . . . .	62
4.2	Model Presentation . . . . .	66
4.3	Forecast . . . . .	74
5	Discussion . . . . .	77
6	Conclusion . . . . .	78
	Bibliography . . . . .	80
	Appendix . . . . .	85
A.1	Code . . . . .	85
A.2	Definitions, Derivations & Proofs . . . . .	86
A.3	Figures . . . . .	90
A.4	Tables . . . . .	98

# List of Symbols, Notation & Abbreviations

Symbol/Notation	Description
$\mathbb{N}$	Set of all positive integers
$\mathbb{R}$	Set of all real numbers
$\mathbb{P}$	Historical probability measure
$\mathbb{Q}$	Equivalent Martingale measure
$W_t^{\mathbb{P}}$	Brownian motion under a measure (here the measure is exemplified with $\mathbb{P}$ )
$\Omega$	Sample space
$\mathcal{F}$	Event space
$\{\mathcal{F}\}_{t \geq 0}$	Filtration
$(\Omega, \mathcal{F}, \{\mathcal{F}\}_{t \geq 0}, \mathbb{P})$	Filtered probability space
$C_t$	State occupied by Markov chain at time- $t$
$p_i$	Density function in state $i$
$\mathbf{P}(r)$	Diagonal matrix with $i$ th diagonal element $p_i$
$I_N$	$N \times N$ -dimensional diagonal matrix with $i$ element 1 (identity matrix)
$S_t$	Random variable at time- $t$ (asset price)
$X_t$	Random variable at time- $t$ (log-return)
$\mathbf{c}^{(-t)}$	$(c_1, \dots, c_{t-1}, c_{t+1}, \dots, c_T)$ (and similarly for $\mathbf{X}$ , $\mathbf{S}$ and $\mathbf{V}$ )
$\mathbf{C}^{(t)}$	$(C_1, C_2, \dots, C_t)$ (and similarly for $\mathbf{X}$ , $\mathbf{S}$ and $\mathbf{V}$ )
$\mathbf{C}_t^T$	$(C_t, C_{t+1}, \dots, C_T)$ (and similarly for $\mathbf{X}$ , $\mathbf{S}$ and $\mathbf{V}$ )
$\alpha_t$	Forward probability, i.e. $\mathbb{P}(\mathbf{X}^{(t)} = \mathbf{x}^{(t)}, C_t = i)$
$\boldsymbol{\alpha}_t$	(Row) vector of forward probabilities
$\beta_t$	Backward probability, i.e. $\mathbb{P}(\mathbf{X}_{t+1}^T = \mathbf{x}_{t+1}^T \mid C_t = i)$
$\boldsymbol{\alpha}_t$	(Row) vector of forward probabilities
$\Gamma$	Transition probability matrix of a Markov chain
$\gamma_{ij}$	$(i, j)$ 'th element in $\Gamma$ ; probability of transitioning from state $i$ to state $j$ in a Markov chain
$\delta$	Stationary distribution of a Markov chain
$\mathbf{1}_N$	$N$ -dimensional vector of 1's
$\mathbf{0}_N$	$N$ -dimensional row vector of 0's
$\mathbf{1}_{N \times N}$	$N \times N$ -dimensional matrix filled with 1's
$e_i$	$(0, \dots, 0, 1, 0, \dots, 0)$ i.e. a (row) vector of dimension $T$ with a 1 in the $t$ 'th entry
$T$	Number of observations
$N$	Number of states
$\phi$	Normalized vector of forward probabilities
$\psi$	Predicted state probabilities
$\mu$	Drift of the Black-Scholes model
$\sigma$	Volatility of the Black-Scholes model
$\rho$	Autoregressive parameter
$\sigma_{\varepsilon}^2$	variance of the innovations $\varepsilon$ of the AR(1) process
$\Delta$	Time-increment between some observations
$\mathcal{C}$	State space
$H$	Hessian matrix
$\text{SE}(\cdot)$	Standard Error of some estimator .
$\mathcal{L}_T$	Likelihood function of $T$ observations
$\ell_T$	Log-likelihood function of $T$ observations
$\xrightarrow{P}$	Convergence in probability
$\xrightarrow{D}$	Convergence in distribution
$\mathcal{N}(\mu, \sigma^2)$	Normal distribution with mean $\mu$ and variance $\sigma^2$
$\mathcal{U}[a, b]$	Uniform distribution function over the range $a$ to $b$
$1_{\{A\}}$	Indicator function of some set $A$
Abbreviation	Description
BS	Black-Scholes
BSM	Black-Scholes model
HMM	Hidden Markov model (sometimes we use the combination BS-HMM)
SSM	State space model (sometimes we use the combination BS-SSM)
TR	Total Return
NTR	Net total return
AR(1)	Autoregressive process of order 1
t.p.m	Transition probability matrix
EMM	Equivalent Martingale measure
nlm	Non-linear maximization
AIC	Akaike information criterion
BIC	Bayesian information criterion
SDE	Stochastic differential equation
a.s.	Almost surely
CI	Confidence interval
LLN	Law of large numbers
CLT	Central limit theorem
HAC	Heteroskedasticity- and autocorrelation-consistent
DAG	Directed acylic graph

## Introduction

The Black-Scholes model (BSM) of asset prices has long been a cornerstone of financial theory and practice. It models the asset price  $S_t$  as a geometric Brownian motion with constant expected rate of return  $\mu$  and volatility  $\sigma$ , an assumption that leads to elegant analytical solutions for option pricing. Furthermore, the asset prices are log-normally distributed. This simplicity, however, comes at the cost of realism. In practice, asset returns exhibit time-varying volatility and occasional jumps or regime shifts that the log-normal Black–Scholes (BS) framework cannot capture. A proposed model to circumvent such jumps is Merton’s Jump-Diffusion Model [37]. This model superimposes a jump component on a diffusion component of the asset price process. Formally, identifying jumps is challenging because large discrete returns can arise either from rare discontinuities or extreme diffusive shocks, making the two statistically indistinguishable in finite samples. Moreover, high-frequency data introduce microstructure noise and volatility clustering effects, which can mimic jumps and bias inference even in sophisticated econometric tests ([1], [3]). Indeed, [12] shows that jumps in financial asset prices are often erroneously identified and, in reality, are rare events that account for only a very small share of total price variation. Empirical returns often have heavier tails and more abrupt changes than the BSM assumes, indicating that the constant- $\sigma$  assumption is too restrictive [44]. Indeed, it is often found that a single set of fixed parameters is inadequate to describe market dynamics across all periods. As market conditions evolve, the static BSM tends to lose predictive accuracy unless its parameters are frequently re-calibrated [19]. This need for continual calibration of  $\mu$  and  $\sigma$  underscores a key limitation of the classical model. That is, it is not well suited for forecasting in environments where volatility and other characteristics change over time.

To address the BSM’s limitations in forecasting, this thesis considers two extensions that relax the assumption of constant parameters. The first is a BSM with a Hidden Markov Model (BS-HMM) for its parameters. In this regime-switching extension, the asset price still follows a diffusion as in BS, but its drift and/or volatility can switch between a finite set of states (regimes). These regime changes are governed by a hidden Markov chain. For instance, the market might alternate between “low-volatility” and “high-volatility” regimes, each with its own  $\sigma$  and possibly different  $\mu$ , with probabilistic transitions between regimes. Such an HMM-based approach can capture phenomena like bull and bear market regimes or sudden volatility shifts that the classical model would miss. By allowing discrete shifts in parameters, the BS-HMM can dynamically adapt to structural changes in the data. We hypothesize that a BS-HMM will improve forecast accuracy by accounting for the possibility of regime shifts (e.g. an upcoming turbulent period) that a single-regime model would underestimate.

The second extension is a BSM with a Continuous State-Space Model (BS-SSM) for the parameters. In this approach, the drift and volatility are treated as latent continuous-time state variables that evolve according to their own stochastic process, specifically, an autoregressive pro-

cess of order 1. The observable data, the prices, are linked to these hidden states, forming a state-space model. Essentially, this is akin to a SVM. Volatility is no longer fixed, but changes over time in a continuous manner and we infer its path from the observed prices. Such state-space formulations are very flexible, allowing  $\sigma_t$  and  $\mu_t$  to vary at each time step and capturing gradual shifts or cyclical patterns in volatility that a BS-HMM might not fully reflect. Conceptually, the BS-SSM treats the BS equation as having time-dependent parameters  $\mu_t$ ,  $\sigma_t$  governed by an underlying dynamical system. This continuous adaptation may better capture the nuanced evolution of market risk over time. Like the BS-HMM, the BS-SSM relaxes the constant parameter assumption, but it does so in a way that allows for more gradual, continuous changes rather than abrupt switches. We expect that by tracking a latent volatility factor, the BS-SSM can forecast future price distributions more accurately during periods of slowly changing market conditions or volatility trends.

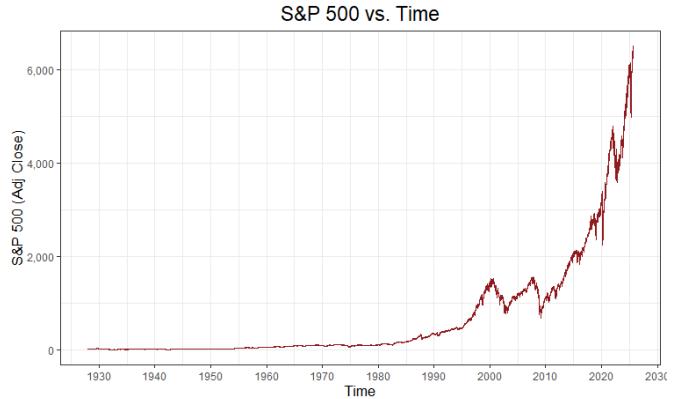
The overarching research question addressed in this thesis is: Do these extended BSMs deliver superior out-of-sample forecasting performance compared to the BSM? In other words, we will test whether incorporating either discrete regime shifts or continuous stochastic parameter evolution leads to more accurate predictions of future asset prices than the traditional model with fixed  $\mu$  and  $\sigma$ . This question is of both academic interest and practical importance. If the extended models can demonstrably improve forecast accuracy, it would suggest a pathway to better risk management and derivative pricing by acknowledging and modeling the non-stationarity in asset dynamics. Conversely, if the extensions do not improve forecasts, that finding is also informative as it would imply that the added complexity does not yields a sufficient performance gain for prediction and the basic BSM, perhaps with frequent recalibration, remains hard to beat.

To answer this question, we conduct an empirical comparison using historical data on the S&P 500 index. The methodology involves estimating each model (the classical BS, the BS-HMM and the BS-SSM) on a training sample and then evaluating their predictive performance on a hold-out sample. The models' parameters are fitted via maximum likelihood estimation, ensuring that each model is calibrated to the historical dynamics of the index. Once calibrated, each model generates out-of-sample forecasts of the S&P 500's price (or rather, log-return distribution) for the evaluation period. The key focus is on out-of-sample forecasting performance. By evaluating the predictions on data not used for estimation, we ensure a fair test of whether the additional flexibility of the HMM or state-space formulations translates into better predictive power. In short, we wish to extend the BSM by state-space models to examine if regime-switching does solve known issues of the BSM.

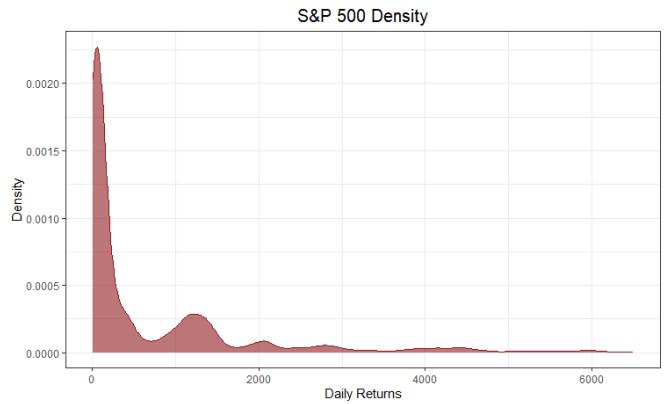
## 1 Data (I of II)

The Standard and Poor's 500 (S&P 500) will be used throughout for analysis. The S&P 500 is a free-float-adjusted, value-weighted index of large-capitalization U.S. equities maintained by S&P Dow Jones Indices. By construction it spans major sectors of the U.S. economy and concentrates on firms with substantial public float, liquidity and operating history, yielding a diversified portfolio in which aggregate—rather than idiosyncratic—risk dominates variation. Its methodology and eligibility criteria (e.g., float adjustment, sector classification, profitability and size thresholds and scheduled reconstitutions) are transparent and stable over time, producing a well-documented data-generating mechanism that supports reproducible empirical work. We will work with the price version of the index, using daily data from Yahoo Finance (`^GSPC`). This series excludes cash dividends. S&P also publishes total return (TR) and net total return (NTR) variants (`^SP500TR`, `^SP500NTR`) that reinvest dividends gross or net of withholding taxes, respectively (see [52, 51, 49, 58]).

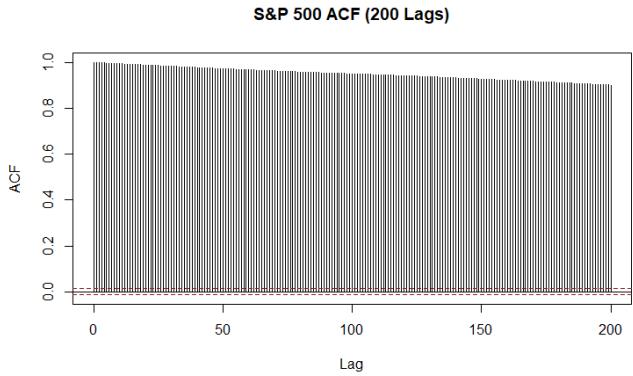
For asset-price analysis, the index offers several practical advantages. It captures a large share of total U.S. equity market capitalization and trading volume, which enhances the signal-to-noise ratio of return observations and reduces the impact of microstructure frictions at daily horizons. The series is long, clean and consistently adjusted for corporate actions, enabling inference across multiple macroeconomic episodes without survivorship bias tied to current constituents.



**Figure 1:** S&P 500 index time series.



**Figure 2:** Kernel density of the S&P 500 index.



**Figure 3:** Autocorrelation of closing prices (200 lags).

The first observation of the data set falls on the date 1927-12-30 at market close and the last on the end of the trading day 2025-09-05. We train the considered models to the 2019-12-31 and forecast on the remaining.

The time series is shown in Figure 1. The density in Figure 2 places most mass near the origin: first quartile = 24.86, median = 103.21 and third quartile = 1076.22. The smallest observation is

4.40 and the largest is 6502.80. The shape appears multi-modal (roughly six modes, with the first three most prominent), which we keep in mind when discussing state-number selection. Finally, [Figure 3](#) shows strong persistence in levels over 200 lags, as expected for nonstationary prices.

**Dividend Treatment and Construction of a Daily Dividend–Yield Series** In this thesis, we adopt the Black–Scholes specification with constant parameters  $(\mu, \sigma, q)$ . With a continuous dividend yield  $q$ , the ex-dividend price index  $S_t$  satisfies

$$\frac{dS_t}{S_t} = (\mu - q) dt + \sigma dW_t^{\mathbb{P}}, \quad d \ln S_t = (\mu - q - \frac{1}{2}\sigma^2) dt + \sigma dW_t^{\mathbb{P}},$$

so estimation on the price index identifies the capital-gains drift  $\mu_{\text{cap}} := \mu - q$  rather than the total-return drift  $\mu_{\text{tot}} := \mu_{\text{cap}} + q$ . By contrast, the S&P total-return ( $\hat{\text{TR}}$ ) index reinvests ordinary cash dividends on the ex-date according to S&P’s index mathematics; the net total return (NTR) variant additionally accounts for withholding taxes [\[50, 53\]](#).

**Post-1988 (Daily): TR–PR Differencing Used to Estimate a Constant  $q$**  Let  $P_t$  denote the S&P 500 price index level (e.g.,  $\hat{\text{GSPC}}$ ) and  $T_t$  the corresponding total-return level (e.g.,  $\hat{\text{TR}}$ ). Over one trading day, for  $t = 2, 3, \dots, T$ ,

$$r_t^{PR} = \frac{P_t}{P_{t-1}} - 1, \quad r_t^{TR} = \frac{T_t}{T_{t-1}} - 1, \quad 1 + r_t^{TR} = (1 + r_t^{PR})(1 + r_t^{div}).$$

Hence the dividend simple return is

$$r_t^{div} = \frac{1 + r_t^{TR}}{1 + r_t^{PR}} - 1,$$

and the log (continuous) dividend yield is

$$q_t^{(\log)} = \ln \frac{T_t}{T_{t-1}} - \ln \frac{P_t}{P_{t-1}}.$$

Let  $m_q := \mathbb{E}[q_t^{(\log)}]$  denote the daily mean log dividend yield and define the continuous annualised dividend yield as  $q := 252 m_q$ . We estimate  $m_q$  by the sample mean

$$\bar{q}^{(\log)} := \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} q_t^{(\log)}, \quad \hat{q} := 252 \bar{q}^{(\log)},$$

where  $\mathcal{T}$  denotes the full set of trading days in the sample. For descriptive purposes, we also report calendar-year averages of  $q_t^{(\log)}$ , annualised in the same way, to illustrate the time variation in dividend yields; these yearly statistics are not used in the actual model estimation.

Statistical inference for the constant-dividend estimator  $\hat{q}$  relies on heteroskedasticity- and

autocorrelation-consistent (HAC) standard errors of Newey–West type, with a naive i.i.d. standard error used as a benchmark. The construction of these standard errors and their extension to regime-specific dividend yields  $\hat{q}_i$ , is detailed in [Section 3](#).

**Pre-TR Period (Monthly): Backfill from Shiller** Before daily TR series are readily available, we employ Shiller’s long-run monthly S&P price and dividend data [47, 46]. Let  $P_m$  be the month-end price index and  $D_m^{TTM}$  the trailing-twelve-month dividend total reported for month  $m$ ; the implied monthly flow is  $d_m := D_m^{TTM}/12$ . The monthly log dividend yield is

$$q_m^{(\log)} = \ln(1 + d_m/P_m).$$

To obtain a daily series that preserves each month’s total, distribute  $q_m^{(\log)}$  evenly across the  $N_m$  business days in month  $m$ :

$$q_t^{(\log)} = \frac{q_m^{(\log)}}{N_m} \quad (t \in m),$$

so that  $\sum_{t \in m} q_t^{(\log)} = q_m^{(\log)}$  by additivity of log returns. This daily backfill is used only to compute the constant estimator  $\hat{q}$  over samples that include pre-1988 months. As a cross-check, S&P’s Dividend Points indices track cumulative dividends in index points and reset annually [48]. Note that this method understates the variance of daily dividends in the pre-1988 sample compared to the post-1988 sample as it creates an artificial “step function” in the pre-1988 dividend volatility. We stress that  $q$  is treated as an exogenous variable estimated via sample means a posteriori to any model maximization

## 2 Theory & Methodology

### 2.1 The Black–Scholes Model

**Model under  $\mathbb{P}$**  The Black and Scholes model [6] (or Black, Scholes and Merton model [36]) assumes that there is a riskless asset with interest rate  $r$  such that the bank/money account  $B$  has time-varying dynamics

$$dB_t = rB_t dt,$$

and that the dynamics of the price of the underlying asset under the historical probability measure  $\mathbb{P}$  are

$$dS_t = (\mu - q)S_t dt + \sigma S_t dW_t^{\mathbb{P}}, \quad S_0 = s_0 \quad (1)$$

It is assumed in the Black-Scholes model that  $r$ ,  $\mu$  and  $\sigma$ , where  $q \geq 0$  denotes a constant dividend yield and that, for valuation purposes, that  $\mu$  is an  $\mathcal{F}$ -adapted process. The solution of [Equation 1](#) can easily be found by using the transformation  $Z_t = \log(S_t)$ , where we assume that a solution exists and that  $S_t$  is a (strictly positive) solution. Itô's formula [5, Thm. 4.19] gives

$$\begin{aligned} dZ_t &= \frac{1}{S_t} dS_t + \frac{1}{2} \left( -\frac{1}{S_t^2} \right) (dS_t)^2 \\ &= \frac{1}{S_t} ((\mu - q)S_t dt + \sigma S_t dW_t^{\mathbb{P}}) + \frac{1}{2} \left( -\frac{1}{S_t^2} \right) \sigma^2 S_t^2 dt \\ &= ((\mu - q)dt + \sigma dW_t^{\mathbb{P}}) - \frac{1}{2} \sigma^2 dt. \end{aligned}$$

This leaves us with the equation

$$dZ_t = \left( \mu - q - \frac{\sigma^2}{2} \right) dt + \sigma dW_t^{\mathbb{P}}, \quad Z_0 = \log s_0. \quad (2)$$

Note that no counts of the r.v.  $Z_t$  is on the RHS of [Equation 2](#). By implication, integrating yields

$$Z_t = \log(s_0) + \left( \mu - q - \frac{\sigma^2}{2} \right) t + \sigma W_t^{\mathbb{P}},$$

which means we obtain the solution to [Equation 1](#) by reversing the log-transformation

$$S_t = s_0 \exp \left( (\mu - q)t + \sigma W_t^{\mathbb{P}} - \frac{\sigma^2}{2} t \right). \quad (3)$$

From this point forward, we shall only consider a fixed time horizon  $[0, \infty]$ .

**Distributional Properties** By normality of the Wiener process increments [5, Def. 4.1], zero-mean property [5, Prop. 4.5] and Itô isometry [5, Prop. 4.5] the distribution of  $Z_t$  (and equivalently  $S_t$ ) is

$$\begin{aligned} Z_t &\sim \mathcal{N}\left(\log(s_0) + \left(\mu - q - \frac{\sigma^2}{2}\right)t, \sigma^2 t\right) \\ \iff S_t &\sim \text{Lognormal}\left(\log(s_0) + \left(\mu - q - \frac{\sigma^2}{2}\right)t, \sigma^2 t\right) \end{aligned}$$

Now, define the continuously compounded log-return over the time horizon  $[t, t-1]$  as

$$X_t = \log\left(\frac{S_t}{S_{t-1}}\right) \tag{4}$$

Using our solution [Equation 3](#) for [Equation 1](#) and by normality of the Wiener process increments [5, Def. 4.1] yields

$$X_t = \left(\mu - q - \frac{\sigma^2}{2}\right)t + \sigma(W_t^{\mathbb{P}} - W_{t-1}^{\mathbb{P}}) \sim \mathcal{N}\left(\left(\mu - q - \frac{\sigma^2}{2}\right)\Delta, \sigma^2\Delta\right), \tag{5}$$

such that the corresponding probability density function is

$$f_{X_t}(x_t) = \frac{1}{\sqrt{2\pi\sigma\Delta}} \exp\left(-\frac{(x_t - (\mu - q - \frac{\sigma^2}{2})\Delta)^2}{2\sigma^2\Delta}\right), \quad x_t \in \mathbb{R} \tag{6}$$

where  $\Delta := t_i - t_{i-1}$  is used to denote the time-increment. Note that if we have  $T$  observations of the asset price we will have  $T-1$  observations of the log-returns.

### 2.1.1 Likelihood Formulation & Parameter Estimation

Let  $\{X_t\}_{t=1}^T$  denote  $T$  observed log-returns with fixed step size  $\Delta > 0$ . From [Equation 5](#), under  $\mathbb{P}$  the returns are i.i.d.

$$X_t \sim \mathcal{N}\left((\mu - q - \frac{1}{2}\sigma^2)\Delta, \sigma^2\Delta\right),$$

where  $\mu$  and  $\sigma > 0$  are constant drift and volatility parameters of the BSM. The BSM with parameter vector  $(\mu, q, \sigma)$  is obviously unidentifiable in terms of the parameters  $\mu$  and  $q$ . As such, for MLE we estimate the *capital-gains* drift  $\mu_{\text{cap}}$  and retrieve the *total-return* drift via  $\mu_{\text{total}} = \mu_{\text{cap}} + \hat{q}$ , where  $\hat{q}$  is the estimated continuous dividend yield (see [Section 1](#)). Henceforth we work explicitly with  $\mu := \mu_{\text{cap}}$  unless otherwise specified

For a realization  $\mathbf{x}^{(T)} = (x_1, \dots, x_T)$  of log-returns with step  $\Delta$ , the likelihood and log-likelihood

are

$$\begin{aligned} \mathcal{L}_T(\mu, \sigma) &= \prod_{t=1}^T \frac{1}{\sqrt{2\pi\sigma^2\Delta}} \exp\left\{-\frac{[x_t - (\mu - \frac{1}{2}\sigma^2)\Delta]^2}{2\sigma^2\Delta}\right\} \\ &\iff \\ \ell_T(\mu, \sigma) &= -\frac{T}{2} \log(2\pi\sigma^2\Delta) - \frac{1}{2\sigma^2\Delta} \sum_{t=1}^T [x_t - (\mu - \frac{1}{2}\sigma^2)\Delta]^2. \end{aligned} \quad (7)$$

### 2.1.2 Simulation

To simulate the standard Black-Scholes model (and the extended models), we need to develop a discretization scheme to simulate the continuous time process defined in [Equation 1](#). To simulate paths for  $(S_t)$  at discrete times  $\mathcal{T} = \{t_i\}_{i=1}^T$ , we generate random samples of  $(S_{t+\Delta})$  given  $(S_t)$  for any increment  $\Delta$ . Repeatedly appending increments constructs the complete path  $(S_t)_{t \in \mathcal{T}}$ . We derive the Euler discretization scheme, often attributed to the work of [\[28\]](#).

**Discretization Scheme** Let  $(\Omega, \mathcal{F}, \{\mathcal{F}\}_{t \geq 0}, \mathbb{P})$  be a filtered probability space. Assume some r.v.  $X_t$  is driven by the SDE

$$dX_t = \mu(X_t, t)dt + \sigma(X_t, t)dW_t^{\mathbb{P}}, \quad (8)$$

where  $W_t^{\mathbb{P}}$  is a Wiener process under  $\mathbb{P}$ . Equally-spaced time increments are used for notational convenience, allowing us to write  $t_i - t_{i-1} := \Delta$ . Integrating [Equation 8](#) from  $t$  to the incremented distance  $t + \Delta$  yields

$$X_{t+\Delta} = X_t + \int_t^{t+\Delta} \mu(X_u, u)du + \int_t^{t+\Delta} \sigma(X_u, u)dW_u^{\mathbb{P}}. \quad (9)$$

At time- $t$ ,  $\hat{X}_t$  is known. We aim to obtain the incremented,  $\hat{X}_{t+\Delta}$ . Euler scheme approximates the integrals using the left end-point rule, such that the deterministic integral of [Equation 9](#) is approximated as the product of the integrand at time- $t$  and the integration range  $\Delta$

$$\int_t^{t+\Delta} \mu(X_u, u)du \approx \mu(X_t, t) \int_t^{t+\Delta} du = \mu(X_t, t)\Delta.$$

Left end-points is a natural candidate as at time- $t$  the value of  $\mu(X_t, t)$  is known. Now, let  $Z^{\mathbb{P}} \sim \mathcal{N}(0, 1)$ . The stochastic integral is approximated as

$$\int_t^{t+\Delta} \sigma(X_u, u)dW_u^{\mathbb{P}} \approx \sigma(X_t, u) \int_t^{t+\Delta} dW_u^{\mathbb{P}} = \sigma(X_t, u)(W_{t+\Delta}^{\mathbb{P}} - W_t^{\mathbb{P}}) = \sigma(X_t, u)\sqrt{\Delta}Z^{\mathbb{P}},$$

because  $W_{t+\Delta}^{\mathbb{P}} - W_t^{\mathbb{P}}$  and  $\sqrt{\Delta}Z^{\mathbb{P}}$  are identically distributed [5, Def. 4.3]. Assembling the results yields the general form of the Euler discretization scheme:

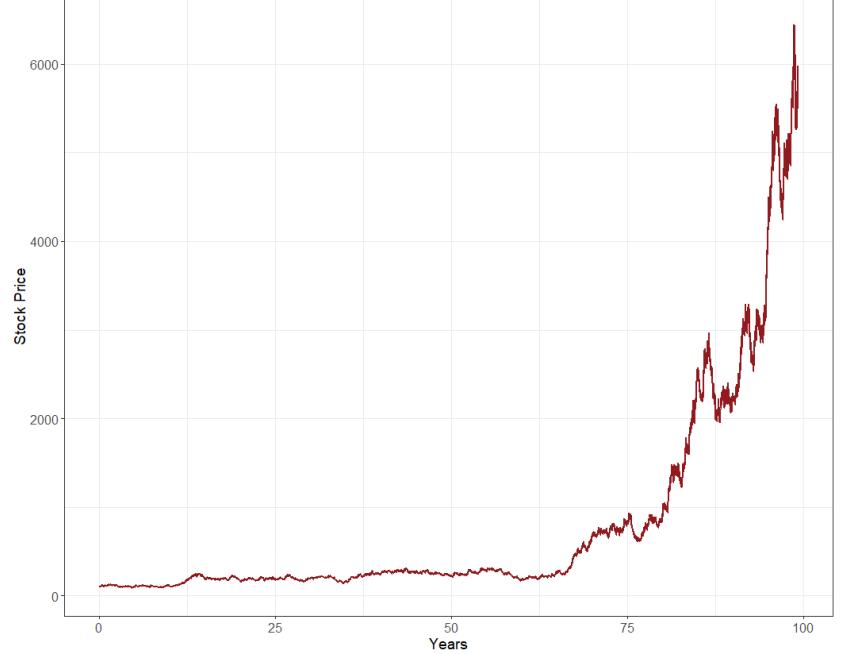
$$\hat{X}_{t+\Delta} = \hat{X}_t + \mu(X_t, t)\Delta + \sigma(X_t, t)\sqrt{\Delta}Z^{\mathbb{P}}. \quad (10)$$

Applying Euler discretization to  $dS_t$  in equation [Equation 1](#) by substituting the diffusion and drift of  $dr_t$  into [Equation 10](#) yields the final discretization Euler scheme of the Black-Scholes' model dynamics

$$\hat{S}_{t+\Delta} = \hat{S}_t + \mu\hat{S}_t\Delta + \sigma\hat{S}_t\sqrt{\Delta}Z^{\mathbb{P}}. \quad (11)$$

The discretization scheme does induce a discretization error to the continuous time process which is highly dependent on parameter subsets and discretization grid roughness as dictated by  $\Delta$ . For a thorough examination and comments see the well-known work of [28, 8, 32].

**Simulating** We simulate the BSM with parameters  $\mu = 0.05$ ,  $\sigma = 0.15$ ,  $S_0 = 100$  and  $n = 25000$  over one realization of the stock price with daily observations, implying  $\Delta = 1/252$  corresponding to approximately  $25000/252 \approx 99$  years. Furthermore, we use the log-return  $X_t$  formulation in the implementation as given in [Equation 4](#). The simulated stock price path is seen in [Figure 4](#) where the used method of discretization was that developed in [Section 2.1.2](#). The estimated values along side the true values are seen in [Table 1](#). Parameter estimates were found using the `nls`-function (Non-Linear Minimization) [14] in R for consistency, although closed-form solutions are available. `nls` is extremely popular for HMMs (see [55], [39], [34], [41], [59]) and provide Hessians along side extremely fast function evaluations via a Newtonian-style algorithm, as compared to i.e. the Nelder–Mead technique which is a heuristic search method. As is quite evident, noting the Euler discretization scheme error, the estimates are accurate.



**Figure 4:** A simulated stock price path  $S_t$  in the BSM.

Parameter	True Value	Estimated Value	Relative Error (%)
$\mu$	0.05000	0.05224	4.47
$\sigma$	0.15000	0.15002	0.0107

**Table 1:** True vs. maximum likelihood estimated BSM parameters,  $\mu$  and  $\sigma$ . MLEs were found using direct numerical maximization using the `nls`-function.

## 2.2 Hidden Markov Models

### 2.2.1 Independent Mixture Models

A basic but useful method to dealing with overdispersed data is with that of a mixture model. Mixture models capture unobserved population heterogeneity: the overall population comprises latent groups, each with its own distribution for the observed variable.

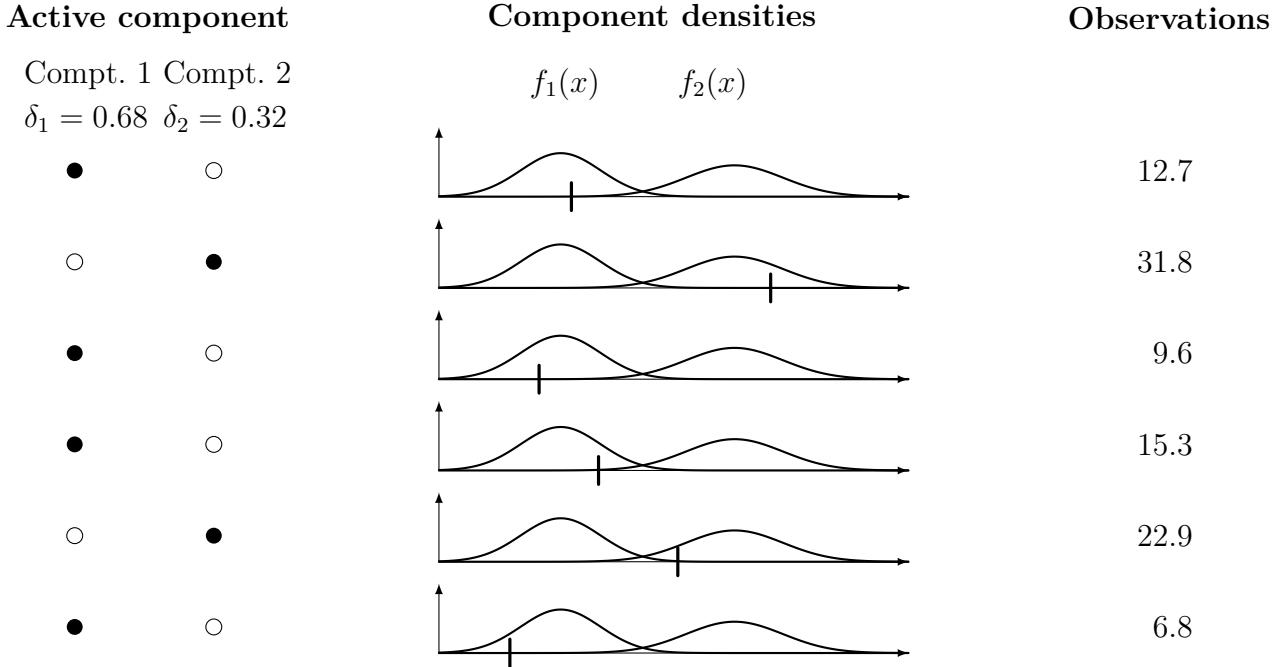
A independent mixture distribution consists of  $N \in \mathbb{N} \setminus \{1\}$  component distributions and a mixing distribution which choses between the different components. Let  $\delta_1, \dots, \delta_N$  denote the probabilities of the  $N$  different components and  $f_1, \dots, f_N$  denote their probability density functions. Let  $X$  denote the continuous r.v. that has mixture distribution and  $C$  the discrete r.v that performs the mixing, i.e.

$$\mathbb{P}(C = k) = \begin{cases} \delta_1, & k = 1, \\ \vdots & \vdots \\ \delta_N, & k = N, \\ 0, & \text{otherwise,} \end{cases} \quad \text{with } \delta_k \geq 0, \sum_{k=1}^N \delta_k = 1.$$

The density function of  $X$  is then

$$f_X(x) = \sum_{i=1}^N f_{X|C}(x | i) \mathbb{P}(C = i) = \sum_{i=1}^N \delta_i f_{X,i}(x).$$

[Figure 5](#) illustrates the generative story of an independent finite mixture with two components. For each observation  $x_j$ , a selector  $C_j \in \{1, 2\}$  is first drawn with  $P(C_j = i) = \delta_i$ . In each row, the filled dot in the left panel marks the realized  $C_j$  (component 1 or 2). The middle panel shows the component densities  $f_1$  and  $f_2$  (stylized), a local  $x$ -axis and a short vertical tick at the location of the realized  $x_j$ . The right panel lists the observation labels. Across rows the selectors  $C_j$  are i.i.d. with mixing weights  $(\delta_1, \delta_2) = (0.68, 0.32)$ ; conditional on  $C_j$ , the  $x_j$  are drawn from the corresponding component distribution.



**Figure 5:** A independent mixture model with 2 components,  $\delta_1 = 0.68$  with corresponding density  $f_1(x)$  and  $\delta_2 = 0.32$  with corresponding density  $f_2(x)$ . The black vertical line on the first axis represents the observation. A filled dot under a component denotes if it is active.

**Parameter Estimation** Let  $\zeta_1, \dots, \zeta_N$  be the parameter vectors of the  $N$ -component distributions,  $\delta_1, \dots, \delta_N$  the mixing parameters where  $\delta_k \geq 0$ ,  $\sum_{k=1}^N \delta_k = 1$  and  $x_1, \dots, x_T$  be the  $T$  observations. The likelihood of a mixture model with  $N$  components is then given by

$$\mathcal{L}_T(\zeta_1, \dots, \zeta_N, \delta_1, \dots, \delta_N) = \prod_{j=1}^T \sum_{i=1}^N \delta_i f_{X_j, i}(x_j). \quad (12)$$

Thus, if the components are specified only by a single parameter,  $2m - 1$  independent parameters have to be estimated by the component sum constraint.

**Unbounded Likelihood in Mixtures** The beforementioned theory relating to the independent mixture distribution could easily be expanded to a discrete case by use probability masses instead of densities. However, one key difference arises; it can happen that in the vicinity of certain parameter subsets, the likelihood is unbounded. In a Gaussian mixture, the likelihood can be made arbitrarily large by assigning a component mean to coincide with an observed data point and letting that component's variance approach zero. In the case of a unbounded likelihood it is often argued in the literature that the MLEs do not exist [45, p. 4630].

In this case, one can somewhat circumvent the complication via a discrete likelihood approxi-

mation of [Equation 12](#)

$$\mathcal{L}_T^{\text{discrete}}(\zeta_1, \dots, \zeta_N, \delta_1, \dots, \delta_N) = \prod_{j=1}^T \sum_{i=1}^N \delta_i \int_{a_j}^{b_j} f_{X_j, i}(x_j),$$

where the interval  $(a_j, b_j)$  consists of those values which would be recorded as  $x_j$ , if observed. For a set of r.v.'s  $X_1, X_2, \dots, X_T$ , the discrete likelihood is of the form  $\mathbb{P}(a_t < X_t < b_t)$  for all  $t$ . An alternative remedy is to enforce a positive lower bound on the component variances and then seek the best local maximum under this constraint.

### 2.2.2 Markov Chains & Hidden Markov Models

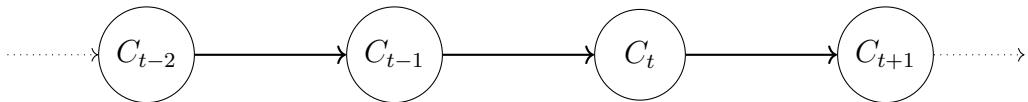
Let  $\{C_t\}_{t \in \mathbb{N}}$  be a sequence of discrete r.v.'s.  $\{C_t\}_{t \in \mathbb{N}}$  is said to be a discrete-time Markov chain if, for all  $t \in \mathbb{N}$ , it satisfies the Markov property

$$\mathbb{P}(C_{t+1} | C_t, C_{t-1}, \dots, C_1) = \mathbb{P}(C_{t+1} | C_t).$$

In other words, condition on the history up to and including time- $t$  only depends on time- $t$ . We will when convenient use the notation  $\mathbf{C}^{(t)} = (C_1, C_2, \dots, C_t)$  such that

$$\mathbb{P}(C_t = j | C_{t-1} = i, \dots, C_1 = k) = \mathbb{P}(C_t = j | C_{t-1} = i),$$

where  $C_t \in \mathcal{C}$  is the state at time  $t = 1, 2, 3, \dots, T$  and  $\mathcal{C}$  is the state space. The Markov property is a first relaxation of the assumption of independence and can be seen visualized in [Figure 6](#).



**Figure 6:** A (first-order) Markov chain for the sequence of discrete r.v.'s  $\{C_t\}_{t \in \mathbb{N}}$ .

A  $N$ -dimensional (or  $N$ -state) hidden Markov model (HMM)  $\{X_t\}_{t \in \mathbb{N}}$  is a dependent mixture model that assumes that the distribution of the observed response variable  $X_t$  depends exclusively on a hidden state  $C_t \in \mathcal{C}$ , where  $\mathcal{C} = \{C_t : t = 1, 2, \dots, T\}$  is modelled by a discrete time  $N$ -state Markov chain, meaning,  $C_t$  satisfies the Markov property. Summarized, the model is

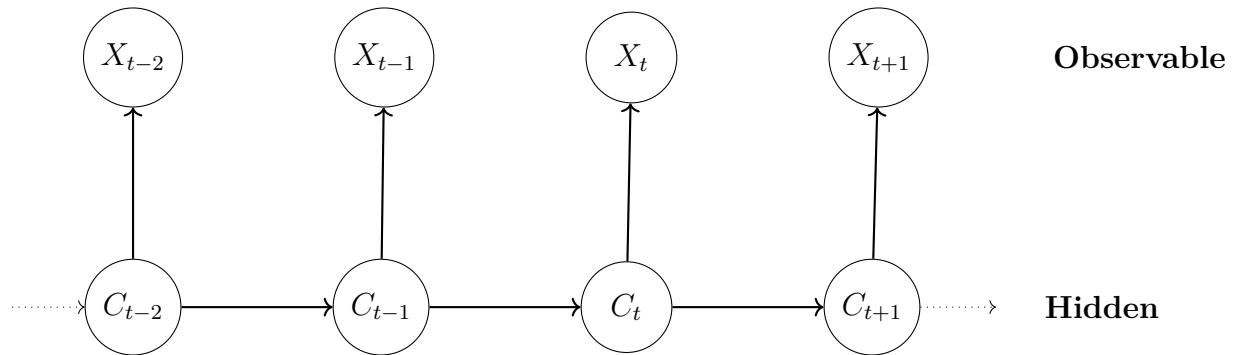
$$\begin{aligned} \mathbb{P}(C_t | \mathbf{C}^{(t-1)}) &= \mathbb{P}(C_t | C_{t-1}), \quad t = 2, 3, \dots, \\ \mathbb{P}(X_t | \mathbf{X}^{(t-1)}, \mathbf{C}^{(t)}) &= \mathbb{P}(X_t | C_t), \quad t \in \mathbb{N}. \end{aligned}$$

The model thus consists of a unobserved/hidden parameter process  $\{C_t\}_{t \in \mathbb{N}}$  satisfying the Markov property and a state-dependent process  $\{X_t\}_{t \in \mathbb{N}}$  in which the distribution of  $X_t$  depends exclusively on the time- $t$  state,  $C_t$ .

Specifically for this thesis, the observed response variable is a state-dependent process,  $\{X_t\}_{t \in \mathbb{N}}$ , the log returns. The process is a noisy observation process in the sense that it is assumed to be produced by a underlying unobserved hidden state process,  $\{C_t\}_{t \in \mathbb{N}}$ . The distribution of  $X_t$ , which is Gaussian, is conditionally independent of previous observations and all states except the current hidden state  $i \in \mathcal{C}$ :

$$f_{X_t|\mathbf{x}^{(t-1)}, \mathbf{C}^{(t)}}(x_t | \mathbf{x}^{(t-1)}, \mathbf{c}^{(t)}) = f_{X_t|C_t}(x_t | i), \quad t = 1, 2, 3, \dots, T,$$

where  $f$  denotes a probability density function. Note that we do not say that  $S_t | S_s$ ,  $t > s$  are unconditionally independent. The structure of a regular hidden Markov model can be seen in [Figure 7](#), where the conditional independence can be intuitively understood.



**Figure 7:** A hidden Markov Model of order 1.

The Markov chain induces dependence in the state-dependent process, meaning, the observations are independent of each other within states.

We will assume time-homogeneity of the Markov chain throughout this paper. The assumption of time-homogeneity of the Markov chain gives rise to the *state transition probabilities* in the  $N \times N$  transition probability matrix (t.p.m.)  $\Gamma$  as

$$\Gamma = \begin{bmatrix} \gamma_{11} & \cdots & \gamma_{1N} \\ \vdots & \ddots & \vdots \\ \gamma_{N1} & \cdots & \gamma_{NN} \end{bmatrix}, \quad \gamma_{ij} = \mathbb{P}(C_{t+1} = j | C_t = i) \in [0, 1], \quad \sum_{j \in \mathcal{C}} \gamma_{ij} = 1, \quad (13)$$

for all  $i, j \in \mathcal{C}$ , where  $\gamma_{ij}$  denotes the probability of transitioning from state  $i$  at time- $t$  to state  $j$  at time- $t + 1$ , where the assumption of time-homogeneity is seen in action by the fact that the transition probabilities do not depend on the time index. We define two key concepts for a Markov chain.

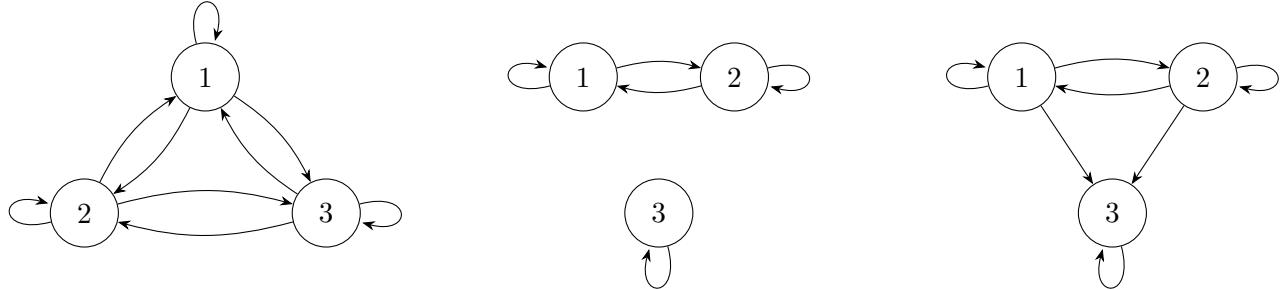
**Definition 2.1.** A Markov chain  $\{C_t\}_{t=0,1,\dots}$  with state space  $\mathcal{C}$  is said to be irreducible if

$$\forall i, j \in \mathcal{C}, \exists t < \infty : \mathbb{P}(C_{n+t} = j | C_n = i) > 0$$

Equivalently, every state can be reached from every other state with positive probability in some finite number of steps. In this case, we say that all states communicate.

The concept of aperiodicity can be seen intuitively in Figure 8.

(a) Irreducible      (b) Reducible: two classes      (c) Reducible: absorbing



**Figure 8:** Three finite Markov chains: (a) irreducible; (b) reducible with two communicating classes; (c) reducible with an absorbing class.

**Definition 2.2.** For a state  $i \in \mathcal{C}$ , define its period as

$$d(i) := \gcd\{ t \geq 1 : (\boldsymbol{\Gamma}^t)_{ii} > 0 \}.$$

A state  $i$  is said to be aperiodic if  $d(i) = 1$ . A Markov chain is called aperiodic if all its states are aperiodic.

The concept of aperiodicity is visualized in Figure 9.

(a) Periodic ( $d = 2$ )      (b) Aperiodic ( $d = 1$ )



**Figure 9:** Periodic vs aperiodic Markov chains. Adding a self-loop breaks periodicity.

The unconditional probabilities of the state process refer to the probability of the process being in state  $i$  at time- $t$ —unconditional of all previous states of the process. These are summarized in the row vector of probabilities

$$\boldsymbol{\delta}^{(t)} = \underbrace{\left[ \mathbb{P}(C_t = 1) \quad \dots \quad \mathbb{P}(C_t = N) \right]}_{1 \times N}, \quad (14)$$

where the number of probabilities equals the number of states of the Markov chain. We let  $\boldsymbol{\delta}^{(1)}$  denote the *initial distribution* of the Markov chain, which provides the probabilities of the process being in the different states at time-1. This allows for a convenient and surprisingly useful result.

**Theorem 2.1.** Let  $\boldsymbol{\delta}^{(t)}$  be defined as in [Equation 14](#). All future distributions of the Markov chain can then be found by

$$\boldsymbol{\delta}^{(t+1)} = \boldsymbol{\delta}^{(t)}\boldsymbol{\Gamma} = \boldsymbol{\delta}^{(1)}\boldsymbol{\Gamma}^{(t)} = \boldsymbol{\delta}\boldsymbol{\Gamma}^t.$$

*Proof.* The first equality simply follows from the law of total probability:

$$\begin{aligned} \delta_i^{(t+1)} &= \mathbb{P}(C_{t+1} = i) \\ &= \sum_{j \in \mathcal{S}} \underbrace{\mathbb{P}(C_t = j)}_{\delta_i^{(t)}} \underbrace{\mathbb{P}(C_{t+1} = j \mid C_t = i)}_{\gamma_{ij}} \\ &\Rightarrow \\ \boldsymbol{\delta}^{(t+1)} &= \begin{bmatrix} \delta_1^{(t+1)} & \dots & \delta_N^{(t+1)} \end{bmatrix} \\ &= \begin{bmatrix} \delta_1^{(t)} & \dots & \delta_N^{(t)} \end{bmatrix} \begin{bmatrix} \gamma_{11} & \dots & \gamma_{1N} \\ \vdots & \ddots & \vdots \\ \gamma_{N1} & \dots & \gamma_{NN} \end{bmatrix} \\ &= \boldsymbol{\delta}^{(t)}\boldsymbol{\Gamma}. \end{aligned}$$

Lastly, equality two and three follows from:

$$\boldsymbol{\delta}^{(t+1)} = \boldsymbol{\delta}^{(t)}\boldsymbol{\Gamma} = \boldsymbol{\delta}^{(t-1)}\boldsymbol{\Gamma}\boldsymbol{\Gamma} = \boldsymbol{\delta}^{(t-2)}\boldsymbol{\Gamma}\boldsymbol{\Gamma}\boldsymbol{\Gamma} = \dots = \boldsymbol{\delta}^{(1)}\boldsymbol{\Gamma}^{t-1}.$$

□

We now turn our attention to the *stationary distribution*. A Markov chain with a t.p.m.  $\boldsymbol{\Gamma}$  is said to have stationary distribution  $\boldsymbol{\delta}$ , a row vector with non-negative elements, if

$$\boldsymbol{\delta}\boldsymbol{\Gamma} = \boldsymbol{\delta}, \quad \boldsymbol{\delta}\mathbf{1}^\top = 1, \tag{15}$$

where  $\mathbf{1}_N$  is a  $N$ -dimensional vector with entries 1. The first of the requirements in [Equation 15](#) expresses the stationarity, i.e. moving forward in time is independent of the t.p.m.,  $\boldsymbol{\Gamma}$ . The second is the requirement that  $\boldsymbol{\delta}$  is indeed a probability distribution. To see why this holds, note that

$$\boldsymbol{\delta}\boldsymbol{\Gamma} = \boldsymbol{\delta} \iff \boldsymbol{\delta} - \boldsymbol{\delta}\boldsymbol{\Gamma} = \mathbf{0}_N \iff \boldsymbol{\delta}(\mathbf{I}_N - \boldsymbol{\Gamma}) = \mathbf{0}_N,$$

where  $\mathbf{0}_N$  is an  $N$ -dimensional row vector of zeros. Now, note that

$$\begin{aligned} \sum_i \delta_i = 1 &\iff \begin{bmatrix} \delta_1 & \dots & \delta_N \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} = 1 \\ &\iff \begin{bmatrix} \delta_1 & \dots & \delta_N \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix} = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \\ &\iff \boldsymbol{\delta} \mathbf{1}_{N \times N} = \mathbf{1}_N. \end{aligned}$$

Adding the two equations, factoring out  $\boldsymbol{\delta}$  and transposing, then yields the desired result

$$\begin{aligned} \boldsymbol{\delta}(\mathbf{I}_N - \boldsymbol{\Gamma} + \mathbf{1}_{N \times N}) = \mathbf{1}_N &\iff (\mathbf{I}_N - \boldsymbol{\Gamma} + \mathbf{1}_{N \times N})^\top \boldsymbol{\delta}^\top = \mathbf{1}_N^\top \\ &\iff \left( \mathbf{I}_N - \boldsymbol{\Gamma} + \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix} \right)^\top \begin{bmatrix} \delta_1 \\ \vdots \\ \delta_N \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \end{aligned}$$

Consequently, a Markov chain started from its stationary distribution will continue to have that distribution at all subsequent time points and we shall refer to such a process as a stationary Markov chain [59, p. 17]. Intuitively, a stationary distribution reflects the long-term proportion of time the model spends in each state.

To find the stationary distribution, one can obtain an explicit expression by solving the following system of equations [59, p. 18]

$$(\mathbf{I}_N - \boldsymbol{\Gamma} + \mathbf{1}_{N \times N})^\top \boldsymbol{\delta}^\top = \mathbf{1}_N^\top, \quad (16)$$

where  $\mathbf{I}_N$  is a  $N$ -dimensional identity matrix,  $\mathbf{1}_{N \times N}$  is a  $N \times N$ -dimensional matrix filled with 1's and  $\mathbf{1}_N$  is a vector filled with 1's.

**Proposition 2.1.** *Let  $\boldsymbol{\Gamma} \in \mathbb{R}^{N \times N}$  be row-stochastic, i.e.  $\boldsymbol{\Gamma} \mathbf{1}_N = \mathbf{1}_N$  and define  $J = \mathbf{1}_N \mathbf{1}_N^\top$ . A row vector  $\boldsymbol{\delta} \in \mathbb{R}^{1 \times N}$  satisfies*

$$\boldsymbol{\delta} \boldsymbol{\Gamma} = \boldsymbol{\delta}, \quad \boldsymbol{\delta} \mathbf{1}_N = 1$$

*if and only if*

$$(\mathbf{I}_N - \boldsymbol{\Gamma} + J)^\top \boldsymbol{\delta}^\top = \mathbf{1}_N.$$

*If  $\boldsymbol{\Gamma}$  is irreducible, then  $\mathbf{A} := (\mathbf{I}_N - \boldsymbol{\Gamma} + J)^\top$  is nonsingular and*

$$\boldsymbol{\delta}^\top = \mathbf{A}^{-1} \mathbf{1}_N.$$

*Proof.* ( $\Rightarrow$ ) If  $\delta\Gamma = \delta$ , then

$$(\mathbf{I}_N - \Gamma^\top)\delta^\top = 0.$$

Since  $\delta\mathbf{1}_N = 1$ , we have

$$\mathbf{J}\delta^\top = (\delta\mathbf{1}_N)\mathbf{1}_N = \mathbf{1}_N.$$

Adding gives

$$(\mathbf{I}_N - \Gamma^\top + \mathbf{J})\delta^\top = (\mathbf{I}_N - \Gamma^\top)\delta^\top + \mathbf{J}\delta^\top = \mathbf{1}_N,$$

equivalently  $(\mathbf{I}_N - \Gamma + \mathbf{J})^\top\delta^\top = \mathbf{1}_N$ .

( $\Leftarrow$ ) Assume  $(\mathbf{I}_N - \Gamma^\top + \mathbf{J})\delta^\top = \mathbf{1}_N$ . Left-multiplying by  $\mathbf{1}_N^\top$  and using  $\mathbf{1}_N^\top\Gamma^\top = \mathbf{1}_N^\top$  and  $\mathbf{1}_N^\top\mathbf{J} = N\mathbf{1}_N^\top$  yields

$$N(\delta\mathbf{1}_N) = \mathbf{1}_N^\top\mathbf{1}_N = N,$$

so  $\delta\mathbf{1}_N = 1$ . Substituting this back gives

$$(\mathbf{I}_N - \Gamma^\top)\delta^\top = \mathbf{1}_N - \mathbf{J}\delta^\top = \mathbf{1}_N - \mathbf{1}_N = 0,$$

so  $\Gamma^\top\delta^\top = \delta^\top$ , i.e.  $\delta\Gamma = \delta$ . Suppose  $(\mathbf{I}_N - \Gamma^\top + \mathbf{J})x = 0$ . Left-multiplying by  $\mathbf{1}_N^\top$  gives

$$\mathbf{1}_N^\top(\mathbf{I}_N - \Gamma^\top + \mathbf{J})x = N\mathbf{1}_N^\top x = 0,$$

so  $\mathbf{1}_N^\top x = 0$ . The equation then reduces to

$$(\mathbf{I}_N - \Gamma^\top)x = 0 \implies \Gamma^\top x = x.$$

For irreducible  $\Gamma$ , the eigenspace of eigenvalue 1 is one-dimensional and spanned by the strictly positive stationary vector. Any nonzero such eigenvector has strictly positive sum, contradicting  $\mathbf{1}_N^\top x = 0$ . Thus  $x = 0$ , so  $\mathbf{I}_N - \Gamma^\top + \mathbf{J}$  is invertible. Therefore

$$\delta^\top = (\mathbf{I}_N - \Gamma^\top + \mathbf{J})^{-1}\mathbf{1}_N.$$

□

When the transition probabilities are time-varying (i.e. functions of covariates), the stationary distribution does not exist [39, p. 14]. However, for fixed covariate values, a single transition probability matrix can be determined, allowing for the computation of a stationary distribution. Throughout we will assume stationarity of the Markov chain. This is adequate as the considered data has a extremely long run time. Furthermore, computational cost is currently extremely demanding which is alleviated by assuming stationarity as it will be evident in [Proposition 2.2](#). We use [Equation 16](#) throughout the code after fitting to find the stationary distribution to ease

computational drag.

### 2.2.3 State-Dependent Distributions

The state-dependent distributions are the probability density functions of  $X_t$  given some state  $i \in \mathcal{C}$  at time- $t$  given by<sup>1</sup>

$$f_{i,X_i}(x_t) := f_{X_t|C_t}(x_t | i).$$

If the state process is stationary, the unconditional distribution of  $S_t$  can be given by

$$f_{X_t}(x_t) \stackrel{\dagger}{=} \sum_{i \in \mathcal{C}} f_{X_t,i}(x_t, i) = \sum_{i \in \mathcal{C}} f_{X_t|C_t}(x_t | i) f(C_t = i) = \sum_{i \in \mathcal{C}} \delta_i^{(t)} f_{i,X_t}(x_t) \stackrel{\dagger\dagger}{=} \sum_{i \in \mathcal{C}} \delta_i f_{i,X_t}(x_t), \quad (17)$$

where  $\dagger$  follows from the law of total probability and  $\dagger\dagger$  by stationarity.

As the log-returns follow a normal distribution, we can directly specify the densities in [Equation 17](#). Namely

$$f_{i,X_t}(x_t) = \frac{1}{\sqrt{2\pi \sigma_i^2 \Delta}} \exp\left(-\frac{(x_t - (\mu_i - \frac{1}{2}\sigma_i^2)\Delta)^2}{2\sigma_i^2 \Delta}\right), \quad x_t \in \mathbb{R}$$

**Parameter Count** As all the parameters are now defined for the BS-HMM we proceed to count the number of parameters to be estimated. The state process is characterized by  $\boldsymbol{\delta}$  and  $\boldsymbol{\Gamma}$ . The latter has  $N \times (N - 1) = N^2 - N$  free parameters due to the row-sum constraint (last equality of [Equation 13](#)). For a stationary Markov chain, we need not estimate the initial distribution as this equals the stationary distribution, which would otherwise yield  $N$  additional parameters (see [Equation 14](#)). As previously stated, we simply use [Equation 16](#) to obtain the stationary distribution after estimating the transition probabilities.

Under conditional independence, the state-dependent process is governed by the state-dependent distributions. In the Black–Scholes setting these are parameterized by  $(\mu, \sigma)$ , so if both are state-dependent we require  $2N$  parameters. The  $N \times N$  transition probabilities. However, as the row-constraint states that the sum of transition probabilities in row  $i$  has to equal 1 this leaves us with  $N \times (N - 1)$  parameters for the T.P.M. The stationary distribution is estimated from eq. ?? which then requires no extra parameters to be estiamted. In total we therefore estimate

$$\# \text{Parameters}_2 = N^2 - N + 2N = N^2 + N.$$

If, instead, exactly one of  $\mu$  or  $\sigma$  is state-dependent (the other being globally state-independent),

---

<sup>1</sup>The notation of the state-dependent density functions should not be confused with a joint density function. We will explicitly write  $i$  as a lower and first index when state-dependent densities are intended and not joint density functions.

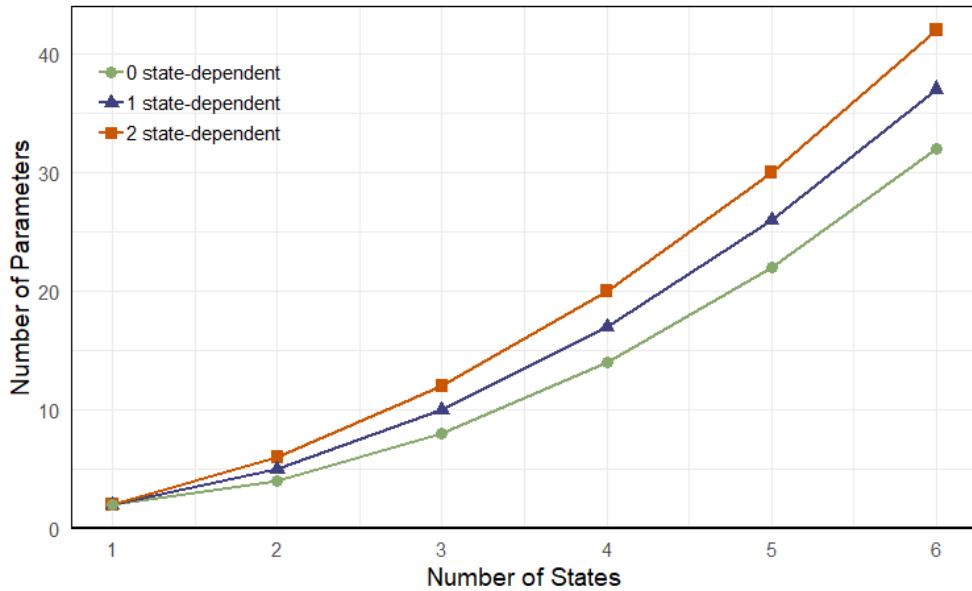
we replace the  $2N$  by  $N + 1$ , yielding

$$\#\text{Parameters}_1 = N^2 - N + (N + 1) = N^2 + 1.$$

Lastly, if *neither*  $\mu$  nor  $\sigma$  is state-dependent (both globally state-independent), we only need to estimate one count of both  $\mu$  and  $\sigma$ , yielding

$$\#\text{Parameters}_0 = N^2 - N + 2.$$

For a visualization of the relation between state-dependent parameters and number of states see [Figure 10](#).



**Figure 10:** Number of parameters as a function of the number of states  $N$  when 0, 1, or 2 of the BS parameters ( $\mu, \sigma$ ) are modeled as state-dependent.

#### 2.2.4 Likelihood Formulation & Parameter Estimation

The likelihood of a hidden Markov model has a convenient recursive form which is seen in the next result. Throughout, let the vector of parameters for estimation be denoted by  $\zeta = (\Gamma, \mu, \sigma)^\top$ .

**Proposition 2.2.** *Let  $\{C_t\}_{t=1}^T$  be a homogeneous, finite  $N$ -state Markov chain on  $\mathcal{C}$ , with transition probability  $\Gamma = (\gamma_{ij})_{i,j=1}^N$ . Consider a observation process  $\{X_t\}_{t=1}^T$ . The joint likelihood of observing  $\{X_t\}_{t=1}^T$  is then given by*

$$\mathcal{L}_T(\zeta) = \delta^{(1)} \mathbf{P}(x_1) \Gamma \mathbf{P}(x_2) \Gamma \mathbf{P}(x_3) \cdots \Gamma \mathbf{P}(x_T) \mathbf{1}^\top,$$

where  $\delta^{(1)}$  is the initial distribution,  $\mathbf{P}(x)$  is the diagonal matrix with the state-dependent distribution  $f_{1,X}(x), f_{2,X}(x), \dots, f_{N,X}(x)$  given in [Equation 17](#) as elements and  $\Gamma$  is the t.p.m.. If  $\delta^{(1)}$

is the stationary distribution  $\boldsymbol{\delta}$  of the Markov chain, then in addition

$$\mathcal{L}_T(\boldsymbol{\zeta}) = \boldsymbol{\delta} \boldsymbol{\Gamma} \mathbf{P}(x_1) \boldsymbol{\Gamma} \mathbf{P}(x_2) \boldsymbol{\Gamma} \mathbf{P}(x_3) \cdots \boldsymbol{\Gamma} \mathbf{P}(x_T) \mathbf{1}_N^\top.$$

*Proof.* Note that

$$\begin{aligned} \mathcal{L}_T(\boldsymbol{\zeta}) &= f_{\mathbf{X}^{(T)}}(\mathbf{x}^{(T)}) \\ &= \sum_{c_1, \dots, c_T=1}^N f_{\mathbf{X}^{(T)} | \mathbf{C}^{(T)}}(\mathbf{x}^{(T)} | \mathbf{c}^{(T)}) \mathbb{P}(\mathbf{C}^{(T)} = \mathbf{c}^{(T)}) \end{aligned}$$

and by Lemma A.2.1

$$\begin{aligned} \mathcal{L}_T(\boldsymbol{\zeta}) &= f_{\mathbf{X}^{(T)}}(\mathbf{x}^{(T)}) \\ &= \mathbb{P}(C_1) \prod_{k=2}^T \mathbb{P}(C_t | C_{t-1}) \prod_{k=1}^T f_{X_k | C_k}(x_k | C_k). \end{aligned}$$

It then follows that

$$\begin{aligned} \mathcal{L}_T(\boldsymbol{\zeta}) &= \sum_{c_1, c_2, \dots, c_T=1}^N (\delta_{c_1} \gamma_{c_1, c_2} \gamma_{c_2, c_3} \cdots \gamma_{c_{T-1}, c_T}) (f_{c_1, X_1}(x_1) f_{c_2, X_2}(x_2) \cdots f_{c_T, X_T}(x_T)) \\ &= \sum_{c_1, c_2, \dots, c_T=1}^N \delta_{c_1} f_{c_1, X_1}(x_1) \gamma_{c_1, c_2} f_{c_2, X_2}(x_2) \gamma_{c_2, c_3} \cdots \gamma_{c_{T-1}, c_T} f_{c_T, X_T}(x_T) \\ &= \boldsymbol{\delta} \mathbf{P}(x_1) \boldsymbol{\Gamma} \mathbf{P}(x_2) \cdots \boldsymbol{\Gamma} \mathbf{P}(x_T) \mathbf{1}_N^\top. \end{aligned}$$

The last equality exploits the fact that a multiple sum of terms having a certain simple multiplicative form can in general be written as a matrix product (see [59, Ex. 7(b)] or Lemma A.2.5). If  $\boldsymbol{\delta}$  is the stationary distribution of the Markov chain, we simply have

$$\boldsymbol{\delta} \mathbf{P}(x_1) = \boldsymbol{\delta} \boldsymbol{\Gamma} \mathbf{P}(x_1).$$

□

The recursive nature of the likelihood in Proposition 2.2 enables computationally efficient evaluation through numerical optimization. The likelihood is maximized using direct numerical methods, leveraging the forward algorithm (which we define in a second).

**Forward Probabilities** The forward algorithm utilizes the forward probabilities which for  $t = 1, 2, \dots, T$  and  $j \in \mathcal{C}$  are given as

$$\alpha_t(j) = f_{\mathbf{X}^{(t)}, C_t}(\mathbf{x}^{(t)}, j) = \mathbb{P}(C_t = j) f_{\mathbf{X}^{(t)}|C_t=j}(\mathbf{x}^{(t)}), \quad \boldsymbol{\alpha}_t = [\alpha_t(1) \dots \alpha_t(N)]. \quad (18)$$

In other words, the forward probabilities contain information on the likelihood of the observations up to and including time- $t$ . Also note that from the definition of  $\boldsymbol{\alpha}_t$  that, for  $t = 1, 2, \dots, T - 1$ ,  $\boldsymbol{\alpha}_{t+1} = \boldsymbol{\alpha}_t \boldsymbol{\Gamma} \mathbf{P}(x_{t+1})$  which can be written in scalar form as

$$\alpha_{t+1}(j) = \left( \sum_{i \in \mathcal{C}} \alpha_t(i) \gamma_{ij} \right) f_{j, X_{t+1}}(x_{t+1}). \quad (19)$$

We are now able to prove the following result to justify their description as probabilities by utilizing the recursive form in [Equation 19](#) and [Lemma A.2.1](#).

**Proposition 2.3.** *For  $t = 1, 2, \dots, T$  and  $j \in \mathcal{C}$ ,*

$$\alpha_t(j) = f_{\mathbf{X}^{(t)}, C_t}(\mathbf{x}^{(t)}, j).$$

*Proof.* Firstly, since  $\boldsymbol{\alpha}_1 = \boldsymbol{\delta} \mathbf{P}(x_1)$  it follows that for  $t = 1$

$$\alpha_1(j) = \delta_j f_{j, X_1}(x_1) = \mathbb{P}(C_1 = j) f_{X_1|C_1}(x_1 | j) = f_{X_1, C_1}(x_1, j).$$

For some  $t \in \mathbb{N}$  we then show it holds for  $t + 1$ :

$$\begin{aligned} \alpha_{t+1}(j) &\stackrel{\dagger}{=} \sum_{i \in \mathcal{C}} \alpha_t(i) \gamma_{ij} f_{j, X_{t+1}}(x_{t+1}) \\ &= \sum_{i \in \mathcal{C}} f_{\mathbf{X}^{(t)}, C_t}(\mathbf{x}^{(t)}, i) \mathbb{P}(C_{t+1} = j | C_t = i) f_{X_{t+1}|C_{t+1}}(x_{t+1} | j) \\ &= \sum_{i \in \mathcal{C}} f_{\mathbf{X}^{(t)}, C_t, C_{t+1}}(\mathbf{x}^{(t)}, i, j) f_{X_{t+1}|C_{t+1}}(x_{t+1} | j) \\ &= \sum_{i \in \mathcal{C}} f_{\mathbf{X}^{(t+1)}, C_t, C_{t+1}}(\mathbf{x}^{(t+1)}, i, j) \\ &\stackrel{\ddagger}{=} f_{\mathbf{X}^{(t+1)}, C_{t+1}}(\mathbf{x}^{(t+1)}, j), \end{aligned}$$

where  $\dagger$  is the scalar forward recursion (matrix–vector form  $\boldsymbol{\alpha}_{t+1} = \boldsymbol{\alpha}_t \boldsymbol{\Gamma} \mathbf{P}(x_{t+1})$ ) and we used the HMM conditional independences  $X_{t+1} \perp (\mathbf{X}^{(t)}, C_t) | C_{t+1}$  and  $\mathbf{X}^{(t)} \perp C_{t+1} | C_t$ . Finally,  $\ddagger$  is marginalization over the discrete r.v.  $C_t$ : summing over  $i$  yields  $f_{\mathbf{X}^{(t+1)}, C_{t+1}}(\mathbf{x}^{(t+1)}, j)$ .  $\square$

Consequently, [Equation 18](#) allows us to write the likelihood from [Proposition 2.2](#) as

$$\mathcal{L}_t(\zeta) = f_{\mathbf{X}^{(t)}}(\mathbf{x}^{(t)}) \stackrel{\dagger}{=} \sum_{j \in \mathcal{C}} f_{\mathbf{X}^{(t)}, C_t}(\mathbf{x}^{(t)}, j) = \sum_{j \in \mathcal{C}} \alpha_t(j).$$

where  $\dagger$  follows from the law of total probability. The probability of the Markov chain occupying state- $j \in \mathcal{C}$  at different times  $t$ , is its proportion of the forward probability at time- $t$  for state  $j$ :

$$f_{C_t | \mathbf{X}^{(t)}}(j | \mathbf{x}^{(t)}) = \frac{f_{C_t, \mathbf{X}^{(t)}}(j, \mathbf{x}^{(t)})}{f_{\mathbf{X}^{(t)}}(\mathbf{x}^{(t)})} = \frac{\alpha_t(j)}{\sum_{i \in \mathcal{C}} \alpha_t(i)}.$$

We can then state the (row) vector of forward probabilities for  $t = 1, 2, \dots, T$  as

$$\boldsymbol{\alpha}_t = \boldsymbol{\delta} \mathbf{P}(x_1) \boldsymbol{\Gamma} \mathbf{P}(x_2) \cdots \boldsymbol{\Gamma} \mathbf{P}(x_t) = \boldsymbol{\delta} \mathbf{P}(x_1) \prod_{s=2}^t \boldsymbol{\Gamma} \mathbf{P}(x_s),$$

following the convention that an empty product is the identity matrix [59, p. 38]. Assembling, [Proposition 2.2](#) states that  $\mathcal{L}_T(\zeta) = \boldsymbol{\alpha}_T \mathbf{1}_N^\top$  and for  $t \geq 2$  we defined  $\boldsymbol{\alpha}_t = \boldsymbol{\alpha}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t)$ . This allows us to define the *forward algorithm*:

$$\begin{aligned} \boldsymbol{\alpha}_1 &= \boldsymbol{\delta} \mathbf{P}(x_1); \\ \boldsymbol{\alpha}_t &= \boldsymbol{\alpha}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t), \quad \text{for } t = 2, 3, \dots, T; \\ \mathcal{L}_T &= \boldsymbol{\alpha}_T \mathbf{1}_N^\top. \end{aligned}$$

Note, for a  $N$ -state HMM,  $\boldsymbol{\delta}$  has  $N$  elements,  $\mathbf{P}(x)$  has  $N$  elements (all in the diagonal) and  $\boldsymbol{\Gamma}$  has  $N \times N$  elements. For the forward algorithm, this implies that  $\boldsymbol{\alpha}_t$  is a sum of  $N$  products consisting of a previous iteration,  $\boldsymbol{\alpha}_{t-1}$ , a transition probability  $\gamma_{ij}$  and a state-dependent probability  $f_{i, X_t}(x_t)$ ,  $i \in \mathcal{C}$ . Hence, for each  $t \in \{1, 2, \dots, T\}$ , there are  $N$  elements to be computed of  $\boldsymbol{\alpha}_t$ . Finally, this implies that the number of operations to calculate the likelihood of  $T$  observations is of order  $TN^2$ .

**Backwards Probabilities** Define

$$\boldsymbol{\beta}_t^\top = \boldsymbol{\Gamma} \mathbf{P}(x_{t+1}) \boldsymbol{\Gamma} \mathbf{P}(x_{t+2}) \cdots \boldsymbol{\Gamma} \mathbf{P}(x_T) \mathbf{1}_N^\top = \left( \prod_{s=t+1}^T \boldsymbol{\Gamma} \mathbf{P}(x_s) \right) \mathbf{1}_N^\top,$$

with the convention that an empty product is the identity matrix, The cast  $t = T$  yields  $\boldsymbol{\beta}_T = \mathbf{1}_N$ . We thne show that  $\beta_t(j)$ , the  $j$ th component of  $\boldsymbol{\beta}_t$ , can be identified as the the conditional

probability

$$f_{\mathbf{X}_{t+1}^T | C_t}(\mathbf{x}_{t+1}^T | i)$$

It then follows that for  $t = 1, 2, \dots, T$ ,

$$\alpha_t(j)\beta_t(j) = f_{\mathbf{X}^{(t)}, C_t}(\mathbf{x}^{(t)}, j)$$

From the definition of  $\boldsymbol{\beta}_t$  it follows that  $\boldsymbol{\beta}_t^\top = \boldsymbol{\Gamma}\mathbf{P}(x_{t+1})\boldsymbol{\beta}_{t+1}^\top$  and hence the name *backwards* probabilities. We now identify the backward probabilities as probabilities by the proposition below.

**Proposition 2.4.** *Assume  $\mathbb{P}(C_t = i) > 0$ . For  $t = 1, 2, \dots, T-1$  and  $i \in \mathcal{C}$ ,*

$$\beta_t(i) = f_{\mathbf{X}_{t+1}^T, C_t}(\mathbf{x}_{t+1}^T | i)$$

*Proof.* We proof the proposition by induction. For  $t = T-1$

$$\beta_{T-1}(i) = \sum_j \mathbb{P}(C_T = j | C_{T-1} = i) f_{X_T, C_T}(x_T | j), \quad (\dagger)$$

since  $\boldsymbol{\beta}_{T-1}^\top = \boldsymbol{\Gamma}\mathbf{P}(x_T)\mathbf{1}_N^\top$ . Furtermore, by Lemma A.2.3

$$\begin{aligned} \mathbb{P}(C_T = j | C_{T-1}=i) f_{X_T|C_T}(x_T | j) &= \mathbb{P}(C_T = j | C_{T-1} = i) f_{X_t|C_{T-1}, C_T}(x_T | i, j) \\ &= f_{X_T, C_{T-1}, C_t}(x_T, i, j) / \mathbb{P}(C_{T-1} = i). \end{aligned} \quad (\dagger\dagger)$$

Substituting  $(\dagger\dagger)$  into  $(\dagger)$  gives

$$\begin{aligned} \beta_{T-1}(i) &= \frac{1}{\mathbb{P}(C_{T-1} = i)} \sum_j f_{X_T, C_{T-1}, C_t}(x_T, i, j) \\ &= f_{X_T, C_{T-1}}(x_t, i) / \mathbb{P}(C_{T-1} = i) \\ &= f_{X_T|C_{T-1}}(x_t | i) \end{aligned}$$

as required.

To demonstrate that validity at time  $t+1$  implies validity at time  $t$ , we begin by observing that the recursive definition of  $\beta_t$ , in combination with the inductive hypothesis, gives

$$\beta_t(i) = \sum_j \gamma_{ij} f_{X_{t+1}|C_{t+1}}(x_{t+1}, j) f_{\mathbf{X}_{t+2}^T | C_{t+1}}(\mathbf{x}_{t+2}^T | j). \quad (\dagger\dagger\dagger)$$

However, Lemma A.2.2 and Lemma A.2.3 imply that

$$f_{X_{t+1}|C_{t+1}}(x_{t+1}, j) f_{\mathbf{X}_{t+2}^T}(\mathbf{x}_{t+2}^T \mid j) = f_{\mathbf{X}_{t+1}^T|C_t, C_{t+1}}(\mathbf{x}_{t+1}^T \mid i, j). \quad (\dagger\dagger\dagger\dagger)$$

Substitute from  $(\dagger\dagger\dagger\dagger)$  into  $(\dagger\dagger\dagger)$  which yields

$$\begin{aligned} \beta_t(i) &= \sum_j \mathbb{P}(C_{t+1} = j \mid C_t = i) f_{\mathbf{X}_{t+1}^T|C_t, C_{t+1}}(\mathbf{x}_{t+1}^T \mid i, j) \\ &= \frac{1}{\mathbb{P}(C_t = 1)} \sum_j f_{\mathbf{X}_{t+1}^T, C_t, C_{t+1}}(\mathbf{x}_{t+1}^T, i, j) \\ &= \frac{f_{\mathbf{X}_{t+1}^T, C_t}(\mathbf{x}_{t+1}^T, i)}{\mathbb{P}(C_t = i)} \\ &= f_{\mathbf{X}_{t+1}^T|C_t}(\mathbf{x}_{t+1}^T \mid i) \end{aligned}$$

which is the required conditional probability.  $\square$

**Scaling the Likelihood** Let  $\mathcal{L}_t(\zeta)$  denote the likelihood of the observations up to time  $t$  under a fixed parameter specification  $\zeta$  of a HMM. Then, under suitable regularity conditions, there exists a constant  $h \in \mathbb{R}$  such that (see [30])

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \mathcal{L}_t(\zeta) = h, \quad \text{a.s.}$$

In particular,

- if  $h < 0$ , the likelihood  $\mathcal{L}_t(\zeta)$  converges to 0 exponentially fast as  $t \rightarrow \infty$ ;
- if  $h > 0$ , the likelihood  $\mathcal{L}_t(\zeta)$  diverges to  $\infty$  exponentially fast as  $t \rightarrow \infty$ .

I.e. the likelihood approaches either 0 or  $\infty$  a.s., exponentially fast. This is highly problematic as our model is already susceptible to numerical overflow complications.

As such, observe firstly from Proposition 2.2, that the HMM likelihood is a product of matrices and not scalars. Consequently, it is not possible to circumvent numerical underflow by computing the logarithm of the likelihood as the sum of logarithms of its factors. Therefore, we adapt the method used by [59, p. 48] (although heavily inspired by [17, p. 78]): For  $t = 1, \dots, T$  define the standardised vector of forward probabilities at time- $t$  as:

$$\boldsymbol{\phi}_t = \frac{\boldsymbol{\alpha}_t}{\boldsymbol{\alpha}_t \mathbf{1}_N^\top} = \frac{\boldsymbol{\alpha}_t}{\sum_{j \in \mathcal{C}} \alpha_t(j)}, \quad \boldsymbol{\phi}_t = [\phi_t(1) \dots \phi_t(N)], \quad \sum_{j \in \mathcal{C}} \phi_t(j) = 1.$$

This yields the normalized forward probabilities, which are far less susceptible to numerical underflow:

For  $t = 1$  :

$$\boldsymbol{\phi}_1 = \frac{\boldsymbol{\alpha}_1}{\boldsymbol{\alpha}_1 \mathbf{1}_N^\top} = \frac{\boldsymbol{\delta}_0 \mathbf{P}(x_1)}{\boldsymbol{\delta}_0 \mathbf{P}(x_1) \mathbf{1}_N^\top}.$$

For  $t = 2, \dots, T$  :

$$\boldsymbol{\phi}_t = \frac{\boldsymbol{\alpha}_t}{\boldsymbol{\alpha}_t \mathbf{1}_N^\top} = \frac{\boldsymbol{\alpha}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t)}{\boldsymbol{\alpha}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t) \mathbf{1}_N^\top} = \frac{\boldsymbol{\alpha}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t) / (\boldsymbol{\alpha}_{t-1} \mathbf{1}_N^\top)}{\boldsymbol{\alpha}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t) \mathbf{1}_N^\top / (\boldsymbol{\alpha}_{t-1} \mathbf{1}_N^\top)} = \frac{\boldsymbol{\phi}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t)}{\boldsymbol{\phi}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t) \mathbf{1}_N^\top}.$$

I.e. scalar multiplication as opposed to matrix multiplication. To see why this is actually the case, we derive the likelihood,  $\mathcal{L}_T(\boldsymbol{\zeta})$ , in terms of  $\boldsymbol{\phi}$  instead of  $\boldsymbol{\alpha}$ .

Firstly, using  $\boldsymbol{\alpha}_0 = \boldsymbol{\delta}$ , note that

$$\mathcal{L}_T(\boldsymbol{\zeta}) = \boldsymbol{\alpha}_T \mathbf{1}^\top = \frac{\boldsymbol{\alpha}_1 \mathbf{1}^\top}{\boldsymbol{\alpha}_0 \mathbf{1}^\top} \frac{\boldsymbol{\alpha}_2 \mathbf{1}^\top}{\boldsymbol{\alpha}_1 \mathbf{1}^\top} \cdots \frac{\boldsymbol{\alpha}_T \mathbf{1}^\top}{\boldsymbol{\alpha}_{T-1} \mathbf{1}^\top} = \prod_{t=1}^T \frac{\boldsymbol{\alpha}_t \mathbf{1}^\top}{\boldsymbol{\alpha}_{t-1} \mathbf{1}^\top}, \quad (20)$$

where  $\frac{\boldsymbol{\alpha}_t \mathbf{1}^\top}{\boldsymbol{\alpha}_{t-1} \mathbf{1}^\top} \in \mathbb{R}$ . This allows us to find the log-likelihood function using [Equation 20](#)

$$\begin{aligned} \ell_T(\boldsymbol{\zeta}) &= \log \mathcal{L}_T(\boldsymbol{\zeta}) \\ &= \log \prod_{t=1}^T \frac{\boldsymbol{\alpha}_t \mathbf{1}_N^\top}{\boldsymbol{\alpha}_{t-1} \mathbf{1}_N^\top} \\ &= \sum_{t=1}^T \log \left( \frac{\boldsymbol{\alpha}_t \mathbf{1}_N^\top}{\boldsymbol{\alpha}_{t-1} \mathbf{1}_N^\top} \right) \\ &= \log \left( \frac{\boldsymbol{\alpha}_1 \mathbf{1}_N^\top}{\boldsymbol{\alpha}_0 \mathbf{1}_N^\top} \right) + \sum_{t=2}^T \log \left( \frac{\boldsymbol{\alpha}_t \mathbf{1}_N^\top}{\boldsymbol{\alpha}_{t-1} \mathbf{1}_N^\top} \right) \\ &= \log \left( \frac{\boldsymbol{\delta} \mathbf{P}(x_1) \mathbf{1}_N^\top}{\boldsymbol{\delta} \mathbf{1}_N^\top} \right) + \sum_{t=2}^T \log \left( \frac{\boldsymbol{\alpha}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t) \mathbf{1}_N^\top}{\boldsymbol{\alpha}_{t-1} \mathbf{1}_N^\top} \right) \\ &= \log (\boldsymbol{\delta} \mathbf{P}(x_1) \mathbf{1}_N^\top) + \sum_{t=2}^T \log (\boldsymbol{\phi}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t) \mathbf{1}_N^\top), \end{aligned}$$

which is exactly stating that the log-likelihood is a sum of logarithmic-values.

Furthermore, we implement working parameters to address constraints of positivity for the CIR model parameters and row sums equaling one in the t.p.m..

For the rest of this paragraph, denote by  $\hat{\cdot}$  the estimator of some parameter  $\cdot$ .

Assume, for example,  $N = 3$ . Firstly, set the *working parameters*  $\eta_i = \log \lambda_i$  for some parameter  $\lambda_i$ . After we have maximized the likelihood with respect to the unconstrained parameters, the constrained parameter estimates can be obtained by transforming back:  $\hat{\lambda}_i = e^{\hat{\eta}_i}$ . Next, start by

defining the matrix with entries  $\tau_{ij} \in \mathbb{R}$

$$\mathbf{T} = \begin{pmatrix} - & \tau_{12} & \tau_{13} \\ \tau_{21} & - & \tau_{23} \\ \tau_{31} & \tau_{32} & - \end{pmatrix},$$

and  $g : \mathbb{R} \rightarrow \mathbb{R}^+$  (strictly increasing) function  $e^x$ . Define

$$\nu_{ij} = \begin{cases} g(\tau_{ij}) & \text{for } i \neq j \\ 1 & \text{for } i = j \end{cases},$$

and

$$\gamma_{ij} = \frac{\nu_{ij}}{\sum_{k=1}^N \nu_{ik}}, \quad i, j = 1, 2, \dots, N,$$

and  $\boldsymbol{\Gamma} = (\gamma_{ij})_{i,j=1}^N$ . We perform the calculation of the likelihood-maximizing parameters in two steps:

- I.** Maximize  $\mathcal{L}_T$  with respect to the working parameters  $\mathbf{T} = \{\tau_{ij}\}$  and  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_N)^\top$  which are all unconstrained by construction.
- II.** Transform the estimates of the working parameters to estimates of the natural parameters:

$$\hat{\mathbf{T}} \rightarrow \hat{\boldsymbol{\Gamma}}, \quad \hat{\boldsymbol{\eta}} \rightarrow \hat{\boldsymbol{\lambda}}.$$

Consider  $\boldsymbol{\Gamma}$  for the case  $g(x) = \exp(x)$  and general  $N$ . Here we have

$$\gamma_{ij} = \frac{\exp(\tau_{ij})}{1 + \sum_{k \neq i} \exp(\tau_{ik})}, \quad i \neq j,$$

and the diagonal elements of  $\boldsymbol{\Gamma}$  follow from the row sums of 1. The transformation in the opposite direction is

$$\tau_{ij} = \log \left( \frac{\gamma_{ij}}{1 - \sum_{k \neq i} \gamma_{ik}} \right) = \log \left( \gamma_{ij}/\gamma_{ii} \right), \quad i \neq j.$$

### 2.2.5 Standard Errors & Confidence Intervals

Unfortunately, relatively little is known about the properties of the MLEs of HMMs [59, p. 56]. However, asymptotic results are available but requires the estimation of the variance-covariance matrix of the estimators of the parameters. It is possible to estimate the standard errors from

the Hessian of the log-likelihood evaluated at the maximum, however, this causes issues when parameters are on the boundary of their parameter space. This phenomenon does unfortunately happen quite often [59, p. 56]. Another method is that of parametric bootstrapping but such methods are heavily reliant on computational power. As is becoming quite evident, the thesis is already extremely reliant on computational power and as such we restrict ourself to the former method.

**Standard Errors via the Hessian** Point estimates of  $\hat{\Theta} = (\hat{\Gamma}, \hat{\lambda})$  are not difficult to compute. However, exact intervals are not available. However, under certain regularity conditions, the MLEs of a HMM parameters are consistent, asymptotically normal and efficient [11, Chapter 12]. Thus, estimating the standard errors of the MLEs can then be used to find approximate confidence intervals by the property of asymptotic normality. It is fairly well known in the litterature, that for (independent) mixture models the sample size has to be very large and that mixtures with small weights or too many components (overfitting) can be of great practical concern [35, p. 68], [23, p. 53].

In order to estimate the standardr errors of the MLEs of an HMM, we use the approximate Hessian of the minus log-likelihood at the minimum supplied by our `nlm`-optimizer in R. This Hessian is inverted to estimat the asymptotic variance-covariance matrix of the estimators of the parameters. However, as the parameters have been transformed, the Hessian is that of the working parameters  $\eta_i$  and not the original natural parameters  $\zeta_i$ . In other words, we have the Hessian at the minimum of  $-\ell$  with respect to the working parameters

$$\mathbf{H}_w = - \left( \frac{\partial^2 \ell}{\partial \eta_i \partial \eta_j} \right),$$

but we are interested in the Hessian with respect to the natural parameters

$$\mathbf{H}_n = - \left( \frac{\partial^2 \ell}{\partial \zeta_i \partial \zeta_j} \right).$$

From [40, p. 247], the following relation holds between  $\mathbf{H}_w$  and  $\mathbf{H}_n$  at the minimum (and for the rest of the paragraph, every matrix is evaluated at the minimum)

$$\mathbf{H}_w = \mathbf{M} \mathbf{H}_n \mathbf{M}^\top \quad \text{and} \quad \mathbf{H}_n^{-1} = \mathbf{M} \mathbf{H}_w \mathbf{M}^\top, \tag{21}$$

where  $\mathbf{M}$  has entries  $m_{ij} = \partial \zeta_j / \partial \eta_i$ . Using Equation 21 and that  $\mathbf{M}$  is readily available, we deduce  $\mathbf{H}_n^{-1}$  from  $\mathbf{H}_w^{-1}$  and then use  $\mathbf{H}_n^{-1}$  to find standard errors for the natural parameters, provided such parameters are not on the boundary of the parameter space.

### 2.2.6 Forecasting, Decoding and State Prediction

In this section,  $\boldsymbol{\delta}$  denotes the initial distribution, but every result is identical if it were to be the stationary distribution.

**Conditional Densities** Using the HMM likelihood formulation discussed in Section ??, we obtain for  $t = 2, 3, \dots, T$  that

$$\begin{aligned} f_{X_t|\mathbf{X}^{(-t)}}(x_t \mid \mathbf{x}^{(-t)}) &= \frac{\boldsymbol{\delta}\mathbf{P}(x_1)\mathbf{B}_2 \cdots \mathbf{B}_{t-1}\boldsymbol{\Gamma}\mathbf{P}(x_t)\mathbf{B}_{t+1} \cdots \mathbf{B}_T\mathbf{1}_N^\top}{\boldsymbol{\delta}\mathbf{P}(x_1)\mathbf{B}_2 \cdots \mathbf{B}_{t-1}\boldsymbol{\Gamma}\mathbf{B}_{t+1} \cdots \mathbf{B}_T\mathbf{1}_N^\top} \\ &\propto \boldsymbol{\alpha}_{t-1}\boldsymbol{\Gamma}\mathbf{P}(x_t)\boldsymbol{\beta}_t^\top, \end{aligned} \quad (22)$$

where  $\mathbf{B}_t = \boldsymbol{\Gamma}\mathbf{P}(x_t)$ ,  $\boldsymbol{\alpha}_t = \boldsymbol{\delta}\mathbf{P}(x_1)\mathbf{B}_2 \cdots \mathbf{B}_t$  and  $\boldsymbol{\beta}_t^\top = \mathbf{B}_{t+1} \cdots \mathbf{B}_T\mathbf{1}_N^\top$ .

For  $t = 1$ , we similarly have

$$\begin{aligned} f_{X_1|\mathbf{X}^{(-1)}}(x_1 \mid \mathbf{x}^{(-1)}) &= \frac{\boldsymbol{\delta}\mathbf{P}(x_1)\mathbf{B}_2 \cdots \mathbf{B}_T\mathbf{1}_N^\top}{\boldsymbol{\delta}\mathbf{I}_N\mathbf{B}_2 \cdots \mathbf{B}_T\mathbf{1}_N^\top} \\ &\propto \boldsymbol{\delta}\mathbf{P}(x_1)\boldsymbol{\beta}_1^\top. \end{aligned} \quad (23)$$

This ratio represents the conditional density of  $x_t$ , where the numerator corresponds to the joint likelihood with the observation at time  $t$  replaced by  $x_t$ , while the denominator is the full likelihood of the observed data, treating  $x_t$  as missing.

These conditional densities can be expressed as mixtures of the state-dependent densities. Since  $\mathbf{P}(x) = \text{diag}(f_1(x), \dots, f_N(x))$  is diagonal, both Equation 22 and Equation 23 yield

$$f_{X_t|\mathbf{X}^{(-t)}}(x_t \mid \mathbf{x}^{(-t)}) \propto \sum_{i \in \mathcal{C}} d_i(t) f_{i,X_t}(x_t),$$

where for Equation 22,  $d_i(t)$  equals the product of the  $i$ 'th entry of  $\boldsymbol{\beta}_t$  and the  $i$ 'th entry of  $\boldsymbol{\alpha}_{t-1}\boldsymbol{\Gamma}$ , while for Equation 23, it is the product of the  $i$ 'th entry of  $\boldsymbol{\beta}_1$  and the  $i$ 'th entry of  $\boldsymbol{\delta}$ . Normalizing these weights gives

$$f_{X_t|\mathbf{X}^{(-t)}}(x_t \mid \mathbf{x}^{(-t)}) = \sum_{i \in \mathcal{C}} w_i(t) f_{i,X_t}(x_t), \quad w_i(t) = \frac{d_i(t)}{\sum_{j \in \mathcal{C}} d_j(t)}.$$

Here,  $w_i(t)$  are mixing weights that depend on the model parameters and on the remaining observations  $\mathbf{x}^{(-t)}$ .

**Forecast Density** Forecast densities are a special case of conditional densities. Let  $h \in \mathbb{Z}^+$  denote the forecast horizon. For continuous-valued observations, the  $h$ -step-ahead forecast density

$f_{X_{T+h}|\mathbf{X}^{(T)}}(x_{T+h} \mid \mathbf{x}^{(T)})$  is obtained analogously to [Equation 22](#):

$$\begin{aligned} f_{X_{T+h}|\mathbf{X}^{(T)}}(x_{T+h} \mid \mathbf{x}^{(T)}) &= \frac{f_{\mathbf{X}^{(T)}, X_{T+h}}(\mathbf{x}^{(T)}, x_{T+h})}{f_{\mathbf{X}^{(T)}}(\mathbf{x}^{(T)})} \\ &= \frac{\delta \mathbf{P}(x_1) \mathbf{B}_2 \cdots \mathbf{B}_T \Gamma^h \mathbf{P}(x_{T+h}) \mathbf{1}_N^\top}{\delta \mathbf{P}(x_1) \mathbf{B}_2 \cdots \mathbf{B}_T \mathbf{1}_N^\top} \\ &= \frac{\alpha_T \Gamma^h \mathbf{P}(x_{T+h}) \mathbf{1}_N^\top}{\alpha_T \mathbf{1}_N^\top} \\ &= \phi_T \Gamma^h \mathbf{P}(x_{T+h}) \mathbf{1}_N^\top, \quad \phi_T = \frac{\alpha_T}{\alpha_T \mathbf{1}_N^\top}. \end{aligned}$$

Thus, the forecast density is also a mixture of the  $N$  state-dependent densities:

$$f_{X_{T+h}|\mathbf{X}^{(T)}}(x_{T+h} \mid \mathbf{x}^{(T)}) = \sum_{i \in \mathcal{C}} \psi_i(h) f_{i, X_{T+h}}(x_{T+h}),$$

where  $\psi_i(h)$  is the  $i$ 'th entry of  $\phi_T \Gamma^h$ . Since the full forecast distribution is available, it is possible to construct both point and full interval forecasts. As the forecast horizon  $h$  increases, the predictive density converges to the marginal stationary density of the HMM, i.e.

$$\lim_{h \rightarrow \infty} f_{X_{T+h}|\mathbf{X}^{(T)}}(x_{T+h} \mid \mathbf{x}^{(T)}) = \lim_{h \rightarrow \infty} \phi_T \Gamma^h \mathbf{P}(x_{T+h}) \mathbf{1}_N^\top = \delta^* \mathbf{P}(x_{T+h}) \mathbf{1}_N^\top, \quad (24)$$

where  $\delta^*$  is the stationary distribution of the Markov chain. The limit follows from the fact that for any nonnegative (row) vector  $\vartheta$  whose entries sum to 1, the vector  $\vartheta \Gamma^h$  approaches  $\delta^*$  as  $h \rightarrow \infty$ , provided that the chain is irreducible and aperiodic [[21](#), p. 394].

**Decoding** We turn to determining the states of the Markov chain that are most likely to have given rise to the observation sequence under the fitted model.

Consider again the vectors of forward  $\alpha_t$  and backward probabilities  $\beta_t$ . For the derivation of the most likely state of the Markov chain at time  $t \in \{1, 2, \dots, T\}$ , we remind of the equation

$$\alpha_t(i) \beta_t(i) = f_{\mathbf{X}^{(t)}, C_t}(\mathbf{x}^{(t)}, i)$$

We can then rewrite the conditional distribution of  $C_t$  given the observations, for  $i \in \mathcal{C}$  as

$$\begin{aligned} \mathbb{P}(C_t = i \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) &= \frac{f_{\mathbf{X}^{(T)}, C_t}(\mathbf{x}^{(T)}, i)}{f_{\mathbf{X}^{(T)}}(\mathbf{x}^{(T)})} \\ &= \frac{\alpha_t(i) \beta_t(i)}{\mathcal{L}_T} \end{aligned}$$

For each  $t \in \{1, \dots, T\}$ , given the observations, the most probable state  $i_t^*$ , is defined as

$$i_t^* = \operatorname{argmax}_{i=1, \dots, N} \mathbb{P}\left(C_t = i \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}\right). \quad (25)$$

In other words, this approach determines the most likely state separately for each time- $t$  by maximizing the conditional probability. We refer to this as local decoding.

However, we are most interested in the most likely sequence of hidden states. As such, we are interested in the quantity

$$(i_1^*, \dots, i_T^*) = \operatorname{argmax}_{(i_1, \dots, i_T) \in \mathcal{C}^T} \mathbb{P}\left(\mathbf{C}^{(T)} = \mathbf{C}^{(T)} \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}\right). \quad (26)$$

Finding the solution of [Equation 26](#) over all possible state sequences involves  $N^T$  function evaluations which is not feasible. A feasible solution is the so called Viterbi algorithm ([\[56\]](#), [\[22\]](#)). Define

$$\xi_{1i} = f_{X_1, C_1}(x_1, i) = \delta_i f_{i, X_1}(x_1),$$

and for  $t = 2, 3, \dots, T$ ,

$$\xi_{ti} = \max_{c_1, c_2, \dots, c_{t-1}} f_{\mathbf{X}^{(t)}, C_t, \mathbf{C}^{(t-1)}}(\mathbf{x}^{(t)}, i, \mathbf{c}^{(t-1)})$$

This leads us to the recursion of  $\xi$ .

**Proposition 2.5.** *For  $t = 2, 3, \dots, T$  and  $j \in \mathcal{C}$ , it follows that*

$$\xi_{ij} = \left( \max_i (\xi_{t-1, i} \gamma_{ij}) \right) f_{j, X_t}(x_t).$$

*Proof.* Fix  $t \geq 2$  and  $j \in \mathcal{C}$ . For any  $(c_1, \dots, c_{t-1}) \in \mathcal{C}^{t-1}$ , by the HMM conditional independences,

$$\begin{aligned} f_{\mathbf{X}^{(t)}, C_t, \mathbf{C}^{(t-1)}}(\mathbf{x}^{(t)}, j, \mathbf{c}^{(t-1)}) &= f_{j, X_t}(x_t) \mathbb{P}(C_t = j \mid C_{t-1} = c_{t-1}) \\ &\quad \times f_{\mathbf{X}^{(t-1)}, C_{t-1}, \mathbf{C}^{(t-2)}}(\mathbf{x}^{(t-1)}, c_{t-1}, \mathbf{c}^{(t-2)}) \\ &= f_{j, X_t}(x_t) \gamma_{c_{t-1} j} f_{\mathbf{X}^{(t-1)}, C_{t-1}, \mathbf{C}^{(t-2)}}(\mathbf{x}^{(t-1)}, c_{t-1}, \mathbf{c}^{(t-2)}). \end{aligned}$$

Maximizing over  $\mathbf{c}^{(t-1)}$  and extracting the factors that do not depend on the maximization variables gives

$$\xi_{tj} = f_{j, X_t}(x_t) \max_{c_{t-1} \in \mathcal{C}} \left\{ \gamma_{c_{t-1} j} \max_{c_1, \dots, c_{t-2} \in \mathcal{C}} f_{\mathbf{X}^{(t-1)}, C_{t-1}, \mathbf{C}^{(t-2)}}(\mathbf{x}^{(t-1)}, c_{t-1}, \mathbf{c}^{(t-2)}) \right\}.$$

By the definition of  $\xi_{t-1,i}$ ,

$$\max_{c_1, \dots, c_{t-2}} f_{\mathbf{X}^{(t-1)}, C_{t-1}, \mathbf{C}^{(t-2)}}(\mathbf{x}^{(t-1)}, i, \mathbf{c}^{(t-2)}) = \xi_{t-1,i},$$

so

$$\xi_{tj} = f_{j, X_t}(x_t) \max_{i \in \mathcal{C}} (\gamma_{ij} \xi_{t-1,i}),$$

which is the desired recursion. The initialization  $\xi_{1i} = \delta_i f_{i, X_1}(x_1)$  follows from  $f_{X_1, C_1}(x_1, i) = f_{i, X_1}(x_1) \mathbb{P}(C_1 = i)$ .  $\square$

The required maximizing sequence of states  $\{i\}_{i=1}^T$  can then be determined recursively from

$$i_T = \operatorname{argmax}_{i=1, \dots, N} \xi_{Ti},$$

and for  $t = T - 1, T - 2, \dots, 1$  from

$$i_t = \operatorname{argmax}_{i=1,2,\dots,N} (\xi_{ti} \gamma_{i, i_{t+1}}).$$

Because the global-decoding objective is a product of probabilities, it's convenient to maximize its logarithm to avoid numerical underflow; the Viterbi recursions translate directly to the log domain. As an alternative, you can use likelihood-style scaling by normalizing each time- $t$  row of the matrix  $\{\xi_{ti}\}$  so that the entries sum to 1. The Viterbi algorithm applies to both stationary and non-stationary (time-inhomogeneous) Markov chains; the initial distribution need not be the stationary distribution.

**State Prediction** We turn our attention to finding conditional distributions of  $C_t$  when  $t > T$ , i.e. state prediction.

**Proposition 2.6.** *Given observations  $x_1, \dots, x_T$ , it follows that*

$$\mathcal{L}_T \mathbb{P} \left( C_t = i \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)} \right) = \begin{cases} \boldsymbol{\alpha}_T \mathbf{\Gamma}^{t-T} \mathbf{e}_i^\top, & \text{for } t > T \quad (\text{state prediction}), \\ \alpha_T(i), & \text{for } t = T \quad (\text{filtering}), \\ \alpha_t(i) \beta_t(i), & \text{for } 1 \leq t < T \quad (\text{smoothing}). \end{cases}$$

*The filtering and smoothing parts (for present or past states) are identical to the state probabilities and could be combined as  $\beta_T(i) = 1$ . The state prediction formula is therefore a generalization to*

$t > T$  and can be restated as

$$\mathbb{P}(C_{T+h} = i \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) = \boldsymbol{\alpha}_T \boldsymbol{\Gamma}^h \mathbf{e}_i^\top / \mathcal{L}_T = \boldsymbol{\phi}_T \boldsymbol{\Gamma} \mathbf{e}_i^\top$$

*Proof.* Fix  $h \geq 1$  and  $i \in \{1, \dots, N\}$ . By the law of total probability over the current state  $C_T$ ,

$$\mathbb{P}(C_{T+h} = i \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) = \sum_{j=1}^N \mathbb{P}(C_{T+h} = i \mid C_T = j, \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) \mathbb{P}(C_T = j \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}).$$

By the (time-homogeneous) Markov property of  $\{C_t\}$  and the HMM conditional-independence structure, the future of the chain depends on the past observations only through  $C_T$ , hence

$$\mathbb{P}(C_{T+h} = i \mid C_T = j, \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) = \mathbb{P}(C_{T+h} = i \mid C_T = j) = (\boldsymbol{\Gamma}^h)_{ji}.$$

Therefore,

$$\begin{aligned} \mathbb{P}(C_{T+h} = i \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) &= \sum_{j=1}^N (\boldsymbol{\Gamma}^h)_{ji} \mathbb{P}(C_T = j \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) \\ &= \boldsymbol{\phi}_T \boldsymbol{\Gamma}^h \mathbf{e}_i^\top. \end{aligned}$$

Using  $\boldsymbol{\phi}_T = \boldsymbol{\alpha}_T / \mathcal{L}_T$  yields the equivalent form  $\mathbb{P}(C_{T+h} = i \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) = \boldsymbol{\alpha}_T \boldsymbol{\Gamma}^h \mathbf{e}_i^\top / \mathcal{L}_T$ .

For the special cases: when  $h = 0$ ,  $\boldsymbol{\Gamma}^0 = \mathbf{I}$ , so  $\mathbb{P}(C_T = i \mid \mathbf{X}^{(T)}) = \phi_T(i) = \alpha_T(i) / \mathcal{L}_T$  (filtering). For  $t < T$ , the standard forward-backward identity  $\mathbb{P}(C_t = i \mid \mathbf{X}^{(T)}) = \alpha_t(i) \beta_t(i) / \mathcal{L}_T$  holds, with  $\beta_T(i) = 1$  by definition of the backward variables, giving the stated smoothing expression.  $\square$

Note that as  $h \rightarrow \infty$ ,  $\boldsymbol{\phi}_T \boldsymbol{\Gamma}^h \rightarrow \boldsymbol{\delta}$  and so  $\mathbb{P}(C_{T+h} = i \mid \mathbf{X}^{(T)} = \mathbf{x}^{(T)}) \rightarrow \delta_i$ .

### 2.2.7 Number of States

HMMs are prone to overfitting [26, p. 2]. That is, HMMs are not well suited for order estimation as small variations in the data are known to cause such models to overestimate the number of groups as well as the frequency of transitions when the number of states is unknown. This raises the question; How does one adequately choose the number of states in a HMM?

In HMMs, the number of states must be specified a priori to the analysis rather than estimated during model fitting by the above-mentioned reason. However, this decision can be challenging, as standard model selection criteria like AIC and BIC often favour a large number of states, which can reduce interpretability<sup>2</sup>. In particular, AIC tends to select more states as increased model flexibility allows for better data fitting, though this can come at the expense of generalizability and interpretability. In other words, we might misspecify some variation in the data as an extra

---

<sup>2</sup>This will become evident in Section 2.4.1.

false state- $i$ , as it gives a better model fitting, even though it might just be a (large) variation within a true state- $j$ . [42] and [39, p. 4] highlighted this issue, recommending that the choice of states should be guided by domain expertise and model validation rather than relying solely on selection criteria. As such, our information criteria, AIC and BIC, will also be utilized for model selection, but not exclusively. We describe the information criteria in [Section 2.4.1](#).

A proposed solution to the a priori number of state selection, is the heuristic method of counting modes in the distribution of the data. However, this can be severely problematic. For example, assume the known number of states is two. If the means are approximately equal but the variance differ it is virtually impossible by visual examination to determine that the number of states is one or some larger integer.

Following [42] and [39], we determine the number of states based on domain expertise only. However, we note that the application of HMMs to financial contexts—and especially to interest rates—remains extremely limited. Consequently, we must develop our own arguments to justify the chosen number of states based on domain expertise.

In the context of equity market modeling, selecting between 1 and 5 states provides a meaningful balance between model complexity, interpretability and macro-financial relevance. A two-state model may capture broad bull and bear market regimes. A third state can correspond to recovery or neutral phases in market cycles. Higher state counts may reflect nuanced market phases, such as asset bubbles, mild corrections, or crashes.

The number of states affects the parameter estimation in a regime-switching Black-Scholes model. In particular, the volatility parameter  $\sigma$  becomes state-dependent. If too many states are included, temporary fluctuations in volatility may be mistaken for persistent regime shifts. Conversely, too few states may obscure significant differences in market regimes, such as between stable bull markets and volatile rallies.

We interpret (or rather hypothesize) the possible values of  $N \in \{1, 2, 3, 4, 5\}$  as follows:

- $N = 1$ : The trivial case—no regime-switching. The standard Black-Scholes model with constant drift and volatility.
- $N = 2$ : A dichotomy of bull and bear markets, capturing broad upturn and downturn market regimes.
- $N = 3$ : Extension to classical business cycle phases: expansion, recession and recovery.
- $N = 4$ : Potential to differentiate mild versus severe market states, e.g., modest bull markets vs. overheated bubbles, or shallow vs. deep downturns. a 4-state system could also be nuances of a bull and bear market, i.e. two bull and two bear market states that allow for varying degrees of severity.
- $N = 5$ : Allows capturing even finer distinctions, such as neutral/stagnant markets or extreme panic phases during financial crises.

For example, during a moderate downturn, the drift  $\mu$  may slightly decline and volatility  $\sigma$  rise modestly, reflecting controlled risk aversion. In contrast, in an extreme crisis,  $\mu$  becomes strongly negative and  $\sigma$  spikes due to panic-driven trading, credit contractions and liquidity crises. Capturing both with the same state may underestimate risk in crisis scenarios or overstate it in moderate corrections.

In summary, based on economic reasoning and business cycle theory, we restrict our analysis to  $N \in \{1, 2, 3, 4, 5\}$ . This reflects plausible macroeconomic regimes and aims to avoid overfitting while maintaining explanatory power and interpretability in financial modeling.

### 2.2.8 Simulation

We simulate the BS-HMM with parameters  $n = 25000$  and daily observations  $\Delta = 1/252$  (approximately  $25000/252 \approx 99$  years). The model parameters are seen in [Table 2](#). To simulate from a  $N$ -state hidden Markov model, we extend the Euler discretized version of the BS SDE given in [Equation 11](#) to include a state sequence from a simulated Markov chain, simply by using the `sample()`-function in base R. Combining the simulated Markov chain with the discretized BS SDE, we achieve the state-dependent Euler-discretized version of the BS SDE

$$\hat{S}_{t+\Delta} = \hat{S}_t + \mu_i \hat{S}_t \Delta + \sigma_i \hat{S}_t \sqrt{\Delta} Z^{\mathbb{P}}, \quad i \in \mathcal{C}$$

We limit ourselves to the most daunting and computationally dragging case, which is the 5-state BS-HMM. The results are seen in [Table 2](#). The simulated price path is shown in figure [??](#). We use `nlm` in R to maximize the likelihood given in [Equation 20](#) for the 5-state BS-HMM with both  $\mu$  and  $\sigma$  state-dependent.

Parameter	True Values					Estimated Values					Relative Error (%)				
$\mu$	0.02000					0.02168					8.40				
	0.04000					0.04511					12.8				
	0.05000					0.02251					54.9				
	0.08000					0.08053					0.656				
	0.10000					0.09927					0.730				
$\sigma$	0.00500					0.00510					2.00				
	0.01000					0.01039					3.90				
	0.01500					0.01419					5.40				
	0.02000					0.01842					7.90				
	0.02500					0.02424					3.04				
$\delta$	0.20000					0.21346					6.73				
	0.20000					0.20233					1.16				
	0.20000					0.15377					23.1				
	0.20000					0.17229					13.9				
	0.20000					0.25815					29.1				
$\Gamma$	0.93174 0.01707 0.01707 0.01707 0.01707					0.92380 0.01792 0.01986 0.01009 0.02833					0.850 4.97 16.3 40.9 66.0				
	0.01707 0.93174 0.01707 0.01707 0.01707					0.01530 0.87999 0.02855 0.04281 0.03336					10.3 5.56 67.3 151 95.4				
	0.01707 0.01707 0.93174 0.01707 0.01707					0.02949 0.00005 0.84818 0.03968 0.05340					72.8 99.7 8.98 132 213				
	0.01707 0.01707 0.01707 0.93174 0.01707					0.02303 0.02058 0.04967 0.87150 0.03522					35.0 20.5 191 6.47 106				
	0.01707 0.01707 0.01707 0.01707 0.93174					0.01962 0.01225 0.03552 0.00732 0.92529					14.9 28.2 108 57.1 0.690				

**Table 2:** True, estimated and relative error (%) for the Black–Scholes 5-state hidden Markov model parameters where  $\mu$  and  $\sigma$  are modeled AS state-dependent.

The issue with the MLEs in the 5-state BS-HMM with both  $\mu$  and  $\sigma$  state-dependent is that the Euler discretization error is exaggerated by the nature of the Markov-switching. A transition from state- $i$  to state- $j$  will cause the convergence of the Euler discretized stock price path to be throttled. A change of state will inevitably reset the convergence and thus induce bias. As such, it is difficult and perhaps a thesis or Ph.d. of its own, to examine the convergence of discretizations in a Markov-switching model.

## 2.3 Continuous State-Space Models

The structure of states in HMMs is often not well suited to problems at hand. As discussed in [Section 2.2.7](#), the number of states is often not known a priori. As such, one has to rely on domain expertise to specify the number of states. In some instances, the choice of the number of states can be determined by model selection criteria and examination of residuals. Furthermore, the states can be intuitive and visited a *reasonably* number of times. However, in fact, most of the time the number of states remains difficult in practice when the underlying data generating process is unknown. Furthermore, as the number of states can possibly be extremely large, the number of parameters rise extremely fast. If we for example have 10 states and 2 state-dependent variables and no state-independent variables, we need to estimate a total of  $10^2 + 10 = 110$  parameters for the 10-state HMM with 2 state-dependent parameters and no state-independent parameters. In these cases it can be advantageous to consider alternative models formulations where the state process is continuous-valued as opposed to discrete-valued and is relatively parsimonious in terms of the number of parameters for estimation.

### 2.3.1 Autoregressive Processes

For the S&P 500 data, it would be intuitive to assume that the rate of occurrence is continuous-valued, as that would allow for gradual change over the years. As means of investment, methods and behavior is constantly changing. A simple model that would capture such changes could be formulated as:

(BS-SSHM):

$$\begin{aligned} C_t &= \rho C_{t-1} + \varepsilon_t, & \varepsilon_t &\sim \mathcal{N}(0, \sigma_\varepsilon^2), \\ \mu_t &= \mu e^{C_t}, \\ \sigma_t &= \sigma e^{C_t}, \end{aligned} \tag{27}$$

$$X_t | C_t \sim \mathcal{N}(\mu_t - \frac{1}{2}\sigma_t^2) \Delta, \sigma_t^2 \Delta,$$

or by factor loading the latent states by constants  $\beta_\mu, \beta_\sigma \in \mathbb{R}$ :

(BS-SSM $_\beta$ ):

$$\begin{aligned} C_t &= \rho C_{t-1} + \varepsilon_t, & \varepsilon_t &\sim \mathcal{N}(0, \sigma_\varepsilon^2), \\ \mu_t &= \mu + \beta_\mu C_t, \\ \sigma_t &= \sigma \exp(\beta_\sigma C_t), \end{aligned} \tag{28}$$

$$X_t | C_t \sim \mathcal{N}((\mu_t - \frac{1}{2}\sigma_t^2) \Delta, \sigma_t^2 \Delta).$$

for  $t = 1, 2, \dots$  and with the recursion initiated in  $C_0 = C \in \mathbb{R}_+$ <sup>3</sup>. The autoregressive parameter is  $\rho \in \mathbb{R}$  and the innovations  $\varepsilon_t$  are independently and identically distributed with a normal distribution with mean zero and variance  $\sigma_\varepsilon^2$ <sup>4</sup>. In other words  $\varepsilon_t$  are i.i.d.  $\mathcal{N}(0, \sigma_\varepsilon^2)$ . This is called an autoregressive process of order 1. It follows that  $\mathbb{E}[C_t | C_{t-1}] = \rho C_{t-1}$  while  $\mathbb{V}[C_t | C_{t-1}] = \sigma_\varepsilon^2$  and therefore the time-dependence, or dynamics, is modelled through the conditional mean of  $C_t$  given the past.

The dynamics of the AR(1) process is clearly seen to be dependent on the autoregressive parameter,  $\rho$ . The simple recursion for  $C_t$  can be written as

$$C_t = \rho^t C + \sum_{i=0}^{t-1} \rho^i \varepsilon_{t-i}.$$

In particular,  $C_t$  is normally distributed with time-varying parameters

$$\begin{aligned} \mathbb{E}[C_t] &= \rho^t C \\ \mathbb{V}[C_t] &= (1 + \rho^2 + \rho^4 + \dots + \rho^{2(t-1)}) \sigma_\varepsilon^2 \end{aligned}$$

---

<sup>3</sup>Note that we index from 0 for the AR(1) theory rather than 1 as the HMM theory. This is to adhere to the general literature and is a simple index shift.

<sup>4</sup>We use the subscript  $\varepsilon$  in  $\sigma_\varepsilon$  to avoid confusion with the BSM parameter  $\sigma$ .

If  $|\rho| < 1$ ,  $\mathbb{E}[C_t] \rightarrow 0$  as  $\rho^t \rightarrow 0$  for  $t \rightarrow \infty$ . Concludingly, for  $|\rho| < 1$ , as  $t \rightarrow \infty$ ,  $C_t$  will resemble the so called linear-process,

$$C_t^* = \sum_{i=0}^{\infty} \rho^i \varepsilon_{t-i},$$

in terms of the sequence  $\{\varepsilon_t\}_{t=\dots,-1,0,1,\dots}$  of i.i.d.  $\mathcal{N}(0, \sigma_\varepsilon^2)$  variables. The process  $\{C_t^*\}_{t=0,1,\dots}$  is then Gaussian distributed with

$$\begin{aligned}\mathbb{E}[C_t^*] &= 0, \\ \mathbb{V}[C_t^*] &\stackrel{\dagger}{=} \frac{\sigma_\varepsilon^2}{(1 - \rho^2)},\end{aligned}$$

as

$$\frac{1 - \rho^{2t}}{1 - \rho^2} \sigma_\varepsilon^2 \xrightarrow{\dagger\dagger} \frac{\sigma_\varepsilon^2}{(1 - \rho^2)}, \quad t \rightarrow \infty,$$

where  $\dagger$  and  $\dagger\dagger$  follows from [Lemma A.2.6](#). The distribution of  $C_t^*$  clearly is clearly independent of time and is an example of a stationary process. Thus, if  $|\rho| < 1$ ,  $C_t$  is asymptotically stationary in the sense that it resembles the stationary process  $C_t^*$  for  $t \rightarrow \infty$ <sup>5</sup>.

**Stationarity & Distributions** Proceeding, we formalize the notion of stationarity and examine distributional properties of the AR(1) process by introducing the notion of a drift function.

**Definition 2.3.** *The process  $\{C_t\}_{t=0,1,\dots}$  is said to be a staionary process if for all  $t, h \geq 0$ , the joint distribution of  $(C_t, \dots, C_{t+h})$  does not depend on  $t \geq 0$ .*

By Definition 2.3, note that for a stationary process with well-defined second order moments,  $\mathbb{E}[C_t]$  and  $\mathbb{V}[C_t]$  are constant and that the covariance between  $C_t, C_{t+h}$ , i.e.  $\text{Cov}[C_t, C_{t+h}]$  depends only on  $h$  and not  $t$ .

The definition of stationarity comments only on the joint distribution of the variables but nothing about dependence over time. Assume  $\{C_t\}_{t \in \mathbb{Z}}$  is a stationary process that is dependent over time with finite second order moment  $\mathbb{E}[|C_t|^2] < \infty$ . A often used indicator to detect dependence is the auto-correlation. For a stationary process  $C_t \in \mathbb{R}$ , the autocovariance function is given by

$$v(h) = \text{Cov}[C_t, C_{t+h}],$$

---

<sup>5</sup>We will formalize the concepts in the next paragraph.

and autocorrelation function (ACF) defined by,

$$\text{ACF}(h) = \text{Corr}[C_t, C_{t+h}] = \frac{\text{Cov}[C_t, C_{t+h}]}{\sqrt{\mathbb{V}[C_t]\mathbb{V}[C_{t+h}]}} \stackrel{\dagger}{=} \frac{v(h)}{v(0)},$$

where  $\dagger$  holds by stationarity. The functions for various  $h$  describe the correlation and hence indicate dependence over time.

In general, "mixing" (i.e., asymptotic independence) captures that the dependence between  $C_t, C_{t+h}$  vanishes as  $h \rightarrow \infty$ . This idea is crucial for time series and replaces the concept of independence. The idea is that a stationary process  $\{C_t\}_{t=0,1,\dots}$  is said to be mixing<sup>6</sup> (or, ergodic) if for all  $t, h$  and sets  $A, B$ ,

$$\mathbb{P}((C_0, \dots, C_t) \in A, (C_h, \dots, C_{t+h}) \in B) \rightarrow \mathbb{P}((C_0, \dots, C_t) \in A)\mathbb{P}((C_0, \dots, c_t) \in B), \quad h \rightarrow \infty$$

The notion is intuitively, that events removed far in time from one another are independent. Importantly, they imply that various Laws of Large Numbers apply.

We now turn our attention to the *drift criterion* which establishes conditions under which LLNS and central limit theorems (CLTs) hold for time series. Let  $\{C_t\}_{t=0,1,\dots}$  be a Markov chain that satisfies the drift criterion. The first implication of satisfying said criterion is that the initial value,  $C_0$ , can be assigned a distribution such that  $X_t$  is stationary. The second implication is finiteness of certain moments for the stationary version. Moreover, variations of LLN and CLT can be applied. Let  $\{C_t\}_{t=0,1,\dots}$  denote an AR(1) process. Then, the distribution of  $C_t | (C_{t-1}, \dots, C_0)$ ,  $t \geq 1$  depends only on  $C_{t-1}$ , meaning,  $C_t | C_{t-1} \sim \mathcal{N}(C_{t-1}\rho, \sigma_\varepsilon^2)$ . As can be seen, the conditional distribution is Gaussian, which has some attractive properties.

We now state 2 assumptions based on those of [54] and [38].

**Assumption 2.2.** Assume that for  $\{C_t\}_{t=1,0,\dots}$  with  $C_t \in \mathbb{R}^p$  it holds that:

(i) The conditional distribution of  $C_t$  given  $(C_{t-1}, C_{t-2}, \dots, C_0)$  depends only on  $C_{t-1}$ , that is

$$C_t | C_{t-1}, C_{t-2}, \dots, C_0 \stackrel{d}{=} C_t | C_{t-1}.$$

(ii) The conditional distribution of  $C_t$  given  $C_{t-n}$ , for some  $n \geq 1$ , has a positive ( $n$ -step) conditional density  $f(y | x) > 0$ , which is continuous in both arguments.

Now, simply note that (i) in Assumption 2.2 implies that  $\{C_t\}_{t=0,1,\dots}$  is a Markov chain on  $\mathbb{R}^p$ , or sometimes called a *Markov chain on a general state space*.

**Example:** For the purpose of our analysis, consider the AR(1) process. As  $\varepsilon_t$  are i.i.d.  $\mathcal{N}(0, \sigma_\varepsilon^2)$

---

<sup>6</sup>The literature differs a lot on the notion on mixing; some even include different kind of mixing ( $\alpha, \beta$ -mixing).

and independent of  $(C_{t-1}, \dots, C_0)$ ,  $C_t$  conditional on  $(C_{t-1}, \dots, C_0)$  has density

$$f_{C_t|C_{t-1}}(c_t | c_{t-1}) = \frac{1}{\sqrt{2\pi\sigma_\varepsilon^2}} \exp\left(-\frac{(c_t - \rho c_{t-1})^2}{2\sigma_\varepsilon^2}\right),$$

which only depends on  $C_{t-1}$  and is well known to be positive and continuous in both arguments.

Next, define the so-called *drift function* that satisfies Assumption [Assumption 2.2](#). A drift function for some time series  $C_t$ , is some function  $\delta(C_t) \geq 1$  and which is not identically  $\infty$ . The role of a drift function is to measure the dynamical drift of  $X_t$  by studying the dynamics of the corresponding drift function  $\delta(C_t)$ . That is, we are interested in  $\mathbb{E}[\delta(C_t) | C_{t-m}]$  for some  $m \geq 1$ .

**Example:** Consider again the AR(1) process with drift function  $\delta(C_t) = 1 + C_t^2$  and  $m = 1$ . Then, using that  $C_{t-1}$  and  $\varepsilon_t$  are independent, we obtain

$$\begin{aligned} \mathbb{E}[\delta(C_t) | C_{t-1}] &= \mathbb{E}[1 + (\rho C_{t-1} + \varepsilon_t)^2 | C_{t-1}] \\ &= 1 + \rho^2 \mathbb{E}[C_{t-1}^2 | C_{t-1}] + 2\rho C_{t-1} \mathbb{E}[\varepsilon_t | C_{t-1}] + \mathbb{E}[\varepsilon_t^2 | C_{t-1}] \\ &= 1 + \rho^2 C_{t-1}^2 + 2\rho C_{t-1} \mathbb{E}[\varepsilon_t] + \mathbb{E}[\varepsilon_t^2] \\ &= 1 + \sigma^2 + \rho^2 C_{t-1}^2 \\ &= \rho^2 \delta(C_{t-1}) + c, \quad c = (1 - \rho^2 + \sigma_\varepsilon^2). \end{aligned}$$

Thus we obtain a process that mimics a AR(1) process in  $\delta(C_t)$ , apart from some constant  $c$ . In other words,

$$\delta(C_t) = \rho^2 \delta(C_{t-1}) + c + \eta_t,$$

with  $\eta_t = (\delta(C_t) - \mathbb{E}[\delta(C_t) | C_{t-1}])$  such that  $\mathbb{E}[\eta_t] = 0$ . As such, if  $\rho^2 < 1$ ,  $\delta(C_t)$  resembles a stationary AR(1) process. This leads us to the final assumption.

**Assumption 2.3.** Assume that  $\{C_t\}_{t=0,1,\dots}$ , with  $C_t \in \mathbb{R}^p$ , satisfies Assumption [Assumption 2.2](#). With drift function  $\delta$ ,  $\delta(C_t) \geq 1$ , assume that there exist positive constants  $M$ ,  $C$  and  $\varphi$  with  $\varphi < 1$ , such that for some  $m \geq 1$ ,

$$(i) \quad \mathbb{E}[\delta(C_{t+m}) | C_t = C] \leq \varphi \delta(C), \quad \text{for } \|C\| > M,$$

$$(ii) \quad \mathbb{E}[\delta(C_{t+m}) | C_t = C] \leq C < \infty, \quad \text{for } \|C\| \leq M.$$

**Example:** Consider again the AR(1). To obtain the desired properties for  $C_t$  by making restrictions on  $\rho$ , we apply the drift criterion with  $\delta(C_t) = 1 + C_t^2$ . Using our previous example, we saw that:

$$\mathbb{E}[\delta(C_t) | C_{t-1} = C] = \rho^2 \delta(C) + c, \quad c = (1 - \rho^2 + \sigma_\varepsilon^2).$$

Since

$$\lim_{|C| \rightarrow \infty} \frac{\mathbb{E}[\delta(C_t) | C_{t-1} = C]}{\delta(C)} = \rho^2,$$

we must require that  $|\rho| < 1$  for the existence of positive constants  $M, \varphi$  with  $\varphi < 1$ , such that  $\mathbb{E}[\delta(C_t) | C_{t-1} = C] \leq \varphi \delta(C)$  for  $|C| > M$ . By continuity of  $\delta(C)$  and  $c$ , it automatically follows that  $\mathbb{E}[\delta(C_t) | C_{t-1} = C] \leq C$  for some  $C > 0$  for  $|y| \leq M$ . In other words, if  $|\rho| < 1$ , the AR(1) process satisfies the drift criterion in [Assumption 2.3](#).

**Theorem 2.4.** *Assume that  $\{C_t\}_{t \geq 1}$  satisfies [Assumption 2.3](#) with drift function  $\delta$ . Then  $C_1$  can be given an initial distribution such that  $C_t$  initiated in  $X_1$  is stationary. With  $C_t$  denoting the stationary version, we have  $\mathbb{E}[\delta(C_t)] < \infty$ . Moreover,  $C_t$  is mixing in the sense that, for any initial value  $C_1$ , the LLN [Lemma 2.1](#) (seen below) applies.*

We turn our attention to distributional properties of the MLE estimators in the AR(1) process. By definition, the density of  $C_t$ , conditional on  $C_{t-1}$  is the Gaussian density with mean  $\rho C_{t-1}$  and variance  $\sigma_\varepsilon^2$ . Moreover, the joint density of  $\{C_t\}_{t=1}^T$  with the initial value  $C_0 = C$  fixed, factorizes as follows

$$f(C_T, C_{T-1}, \dots, C_1 | C_0) = \prod_{t=1}^T f(C_t | C_{t-1}) \quad (29)$$

Denote the likelihood function as  $\mathcal{L}(\rho, \sigma^2)$ . Then by the factorization in [Equation 29](#) of the joint density of  $C_1, \dots, C_T$  given  $C_0$ , gives the log-likelihood function

$$\begin{aligned} \ell(\rho, \sigma^2) &= \log \mathcal{L}(\rho, \sigma^2) = \log \left( \prod_{t=2}^T f(C_t | C_{t-1}) \right) \\ &= -\frac{T}{2} \underbrace{\log(2\pi)}_{*} - \frac{T}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{t=2}^T (C_t - \rho C_{t-1})^2. \end{aligned} \quad (30)$$

Note that the term  $*$  does not contain any parameters and does not matter for parameter estimation. However, we include it for likelihood calculations as AIC and BIC will be reported.

**Theorem 2.5.** *The MLEs of  $\rho$  and  $\sigma_\varepsilon$  for the AR(1) model are given by*

$$\begin{aligned} \hat{\rho} &= \frac{\frac{1}{T} \sum_{t=2}^T C_t C_{t-1}}{\frac{1}{T} \sum_{t=1}^T C_{t-1}^2} \\ \hat{\sigma}^2 &= \frac{1}{T} \sum_{t=1}^T (C_t - \hat{\rho} C_{t-1})^2 \end{aligned}$$

Ignoring a constant factor, the maximized likelihood function is given by

$$\mathcal{L}(\hat{\rho}, \hat{\sigma}^2) = (2\pi e \hat{\sigma}^2)^{-T/2}.$$

*Proof.* Differentiating Equation 30 with respect to the parameters  $(\rho, \sigma^2)$  gives us FOCs:

$$(1) : \sum_{t=2}^T (C_t - \rho C_{t-1}) C_{t-1} = 0, \quad (2) : \frac{1}{T} \sum_{t=2}^T (C_t - \rho C_{t-1})^2 = \sigma^2.$$

The first equality (1) instantly leads to the MLE of  $\hat{\rho}$  by multiplying the parenthesis and isolating  $\rho$ . Substituting  $\hat{\rho}$  into the second FOC (2) gives us that  $\hat{\sigma}^2$  is exactly the residual sum of squares

$$\hat{\sigma}^2 = \frac{1}{T} \sum_{t=1}^T (C_t - \hat{\rho} C_{t-1})^2.$$

Lastly, the second order derivatives evaluated at  $(\hat{\rho}, \hat{\sigma}^2)$  equal

$$\begin{aligned} \left. \frac{\partial^2}{\partial \rho^2} \log \mathcal{L} \right|_{(\hat{\rho}, \hat{\sigma}^2)} &= -\frac{T}{\hat{\sigma}^2} \left( \sum_{t=1}^T C_{t-1}^2 \right), \\ \left. \frac{\partial^2}{\partial (\sigma^2)^2} \log \mathcal{L} \right|_{(\hat{\rho}, \hat{\sigma}^2)} &= -\frac{T}{2\hat{\sigma}^4}, \\ \left. \frac{\partial^2}{\partial \sigma^2 \partial \rho} \log \mathcal{L} \right|_{(\hat{\rho}, \hat{\sigma}^2)} &= 0, \end{aligned}$$

implying that  $\mathcal{L}(\rho, \sigma^2)$  has maximum given by  $(2\pi)^{-T/2} e^{-T/2} (\hat{\sigma}^2)^{-T/2} = (2\pi e \hat{\sigma}^2)^{-T/2}$ .  $\square$

Next, we consider asymptotic properties of the MLEs.  ${}_0$  will denote the values of the parameters under which the probabilistic arguments are made, or in a more intuitive sense, the true-values. However, firstly, one of the infamous Law of Large Numbers (LLN) [43, p. 17] and Central Limit Theorem (CLT) [43, p. 20].

**Lemma 2.1.** *Assume that with  $C_t \in \mathbb{R}^p$ ,  $\{C_t\}_{t=0}^T$  is a geometrically ergodic markov chain with statioanry solution  $\{C_t^*\}$ . Assume furthermore that the function  $g : \mathbb{R}^{p(m+1)} \rightarrow \mathbb{R}$ ,  $m \geq 0$ , satisfies  $\mathbb{E}[g(C_t^*, C_{t-1}^*, \dots, C_{t-m}^*)] < \infty$ , then as  $T \rightarrow \infty$ ,*

$$\frac{1}{T} \sum_{t=1}^T g(C_t, C_{t-1}, \dots, X_{t-m}) \xrightarrow{\mathcal{P}} \mathbb{E}[g(C_t^*, C_{t-1}^*, \dots, C_{t-m}^*)].$$

**Lemma 2.2.** *For a given sequence  $\{C_t\}_{t \geq 1}$ , consider  $C_t = f(C_t, C_{t-1}, \dots, C_m)$ , with  $f$  continuous,*

with  $\mathbb{E}[Y_t | \mathcal{F}_{t-1}] = 0$ , where  $\mathcal{F}_t = (C_t, \dots, C_0)$ . If

$$\mathbf{I}: \quad \frac{1}{T} \sum_{t=1}^T \mathbb{E}[Y_t^2 | \mathcal{F}_{t-1}] \xrightarrow{\mathcal{P}} \sigma^2 > 0$$

and either **II** or **II'** hold,

$$\begin{aligned}\mathbf{II}: \quad & \frac{1}{T} \sum_{t=1}^T \mathbb{E}[Y_t^2 \mathbb{1}_{\{|Y_t| > \delta T^{1/2}\}}] \rightarrow 0, \\ \mathbf{II'}: \quad & \frac{1}{T} \sum_{t=1}^T \mathbb{E}[Y_t^2 \mathbb{1}_{\{|Y_t| > \delta T^{1/2}\}} | \mathcal{F}_{t-1}] \xrightarrow{\mathcal{P}} 0,\end{aligned}$$

for any  $\delta > 0$ , then  $\frac{1}{\sqrt{T}} \sum_{t=1}^T Y_t \xrightarrow{\mathcal{P}} \mathcal{N}(0, \sigma^2)$ .

**Theorem 2.6.** For  $|\rho_0| < 1$ , the MLEs of the AR(1) model are consistent,  $\hat{\rho} \xrightarrow{\mathcal{P}} \rho_0$  and  $\hat{\sigma} \xrightarrow{\mathcal{P}} \sigma_0$  as  $T \rightarrow \infty$ . Moreover,

$$\sqrt{T}(\hat{\rho} - \rho_0) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1 - \rho_0^2),$$

and a consistent estimator of the asymptotic variance is given by  $\hat{\sigma}^2 \left( \frac{1}{T} \sum_{t=1}^T C_{t-1}^2 \right)$ , such that

$$\sqrt{T}(\hat{\rho} - \rho_0) \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T C_{t-1}^2}{\hat{\sigma}^2}} \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1), \quad T \rightarrow \infty.$$

**Note:** The above distributional statements do not hold without the crucial condition  $|\rho_0| < 1$ .

*Proof.* Since  $|\rho_0| < 1$  we recall that  $C_t$  is geometrically ergodic and the LLN Lemma 2.1 applies. Consider  $\hat{\rho}$  as given by,

$$\hat{\rho} = \left( \frac{1}{T} \sum_{t=1}^T C_t C_{t-1} \right) \left( \frac{1}{T} \sum_{t=1}^T C_{t-1}^2 \right)^{-1}.$$

By the LLN applied to  $C_t C_{t-1}$  and  $C_{t-1}^2$  and as  $\mathbb{E}[C_t^*]^2 < \infty$ , it follows directly that

$$\hat{\rho} \xrightarrow{\mathcal{P}} \text{Cov}[C_t^*, C_{t-1}^*] / \mathbb{V}[C_t^*] = \rho_0.$$

Furthermore,

$$\begin{aligned}\hat{\sigma}^2 & \xrightarrow{\mathcal{P}} \mathbb{V}[C_t^*] - (\text{Cov}[C_t^*, C_{t-1}^*])^2 / \mathbb{V}[C_t^*] \\ & = \frac{\sigma_0^2}{1 - \rho_0^2} - \rho_0^2 \frac{\sigma_0^2}{1 - \rho_0^2} = \sigma_0^2,\end{aligned}$$

as claimed and it follows that

$$\hat{\sigma}^2 \left( \frac{1}{T} \sum_{t=1}^T C_{t-1}^2 \right) \xrightarrow{\mathcal{P}} 1 - \rho_0^2.$$

For the asymptotic distribution, we have that

$$\sqrt{T}(\hat{\rho} - \rho_0) = \frac{\frac{1}{\sqrt{T}} \sum_{t=1}^T \varepsilon_t C_{t-1}}{\frac{1}{T} \sum_{t=1}^T C_{t-1}^2},$$

where  $\frac{1}{T} \sum_{t=1}^T C_{t-1}^2 \xrightarrow{\mathcal{P}} \frac{\sigma_0^2}{1-\rho_0^2}$  by the LLN [Lemma 2.1](#). Define  $Y_t \equiv \varepsilon_t C_{t-1}$ . As  $\varepsilon_t = C_t - \rho_0 C_{t-1}$ ,  $Y_t$  is a martingale difference sequence with respect to  $\mathcal{F}_t$  and thus the CLT [Lemma 2.2](#) can be applied. Observat that

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[Y_t | \mathcal{F}_{t-1}] = \sigma_{\varepsilon,0}^2 \left( \frac{1}{T} \sum_{t=1}^T C_{t-1}^2 \right) \xrightarrow{\mathcal{P}} \frac{\sigma_{\varepsilon,0}^4}{1-\rho_0^2}.$$

Next let  $A = \{|X| > \delta\sqrt{T}\}$  with  $\delta > 0$ ,  $T > 0$  and assume  $\mathbb{E}[X^4] < \infty$ . Then by Cauchy–Schwarz and Markov’s inequality (see, e.g., [\[16, 4, 27\]](#))

$$\begin{aligned} \mathbb{E}[X^2 \mathbb{1}_{\{A\}}] &\leq \left( \mathbb{E}[X^4] \right)^{1/2} \left( \mathbb{E}[\mathbb{1}_{\{A\}}] \right)^{1/2} && \text{(by Cauchy–Schwarz)} \\ &= \left( \mathbb{E}[X^4] \right)^{1/2} \mathbb{P}(A)^{1/2} \\ &\leq \left( \mathbb{E}[X^4] \right)^{1/2} \left( \frac{\mathbb{E}[X^4]}{\delta^4 T^2} \right)^{1/2} && \text{(by Markov on } X^4\text{)} \\ &= \frac{\mathbb{E}[X^4]}{\delta^2 T}. \end{aligned}$$

Using the inequality yields

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[ Y_t^2 \mathbb{1}_{\{|Y_t| > \delta\sqrt{T}\}} \mid \mathcal{F}_{t-1} \right] &\leq \frac{1}{T^2 \delta^2} \sum_{t=1}^T \mathbb{E} [Y_t^4 \mid \mathcal{F}_{t-1}] \\ &= \frac{1}{T \delta^2} \left( \frac{1}{T} \sum_{t=1}^T C_{t-1}^4 \right) \mathbb{E} [\varepsilon_t^4] \xrightarrow{\mathcal{P}} 0. \end{aligned}$$

We conclude that

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T Y_t = \frac{1}{\sqrt{T}} \sum_{t=1}^T \varepsilon_t C_{t-1} \xrightarrow{\mathcal{D}} \mathcal{N} \left( 0, \frac{\sigma_{\varepsilon,0}^4}{1-\rho_0^2} \right).$$

Lastly, substituting in gives us the final resul for  $\hat{\rho}$ . □

**Unit roots** The underlying assumption of the before-mentioned analyses relied on the fact that  $|\rho| < 1$ . As stated in cite ?? (under theorem II.2.2), often met in the analysis of stock prices, it will be the case that  $\rho = 1$ . However, we shall shortly argue on why  $\rho = 1$ , i.e. the unit root case and cointegration, is not examined further in this thesis. When  $\rho = 1$ , it follows that the AR(1) process with  $C_0 = C$  fixed,  $\varepsilon_t$  i.i.d.  $\mathcal{N}(0, \sigma_\varepsilon^2)$ ,

$$C_t = C_{t-1} + \varepsilon_t = \sum_{i=1}^t \varepsilon_i + C. \quad (31)$$

In other words,  $C_t$  is the sum of a random walk  $\sum_{i=1}^t \varepsilon_i$  and the initial value  $C$ . When  $\rho = 1$ ,  $x_t$  is not stationary, not even asymptotically. This can easily be seen by using the fact that  $\varepsilon_t$  i.i.d.  $\mathcal{N}(0, \sigma_\varepsilon^2)$ , to see that the variance of  $C_t$  in Equation 31 is given by

$$\mathbb{V}[C_t] = \mathbb{V}\left[\sum_{i=1}^t \varepsilon_i + C\right] = t\sigma_\varepsilon^2,$$

which is increasing in  $t$ . Furthermore, note that

$$\Delta C_t = C_t - C_{t-1} = \sum_{i=1}^t \varepsilon_i + C - \sum_{i=1}^{t-1} \varepsilon_i + C = \varepsilon_t,$$

implying that the differenced process is stationary. Unit root analysis provides a framework to discriminate between the pair: The former is a non-stationary random walk case (the hypothesis of non-stationarity) and the latter a stationary case (the hypothesis of stationarity).

It is now apparent why the raw S&P 500 data was transformed to returns. However, note that we assumed that the state process is modelled through an autoregressive process of order 1 and not the asset prices.

A unit-root state behaves like a random walk, drifting without pullback. This destroys a stable baseline for “high/low” regimes and undermines interpretability. It also blurs identification (level shifts can be traded between the intercept and the state) and the state’s uncertainty grows with sample length, so there is no steady-state signal-to-noise. If the state feeds the mean or (especially) the variance, implied moments can blow up over time. For stable inference and meaningful regimes, we therefore impose  $|\rho| < 1$ . Furthermore, as the state-process is unobserved/hidden/latent, we have no means of testing for the hypothesis of non-stationarity for the state-process but it is possible for the asset price process.

### 2.3.2 Likelihood Formulation & Parameter Estimation

We consider the basic SSM in which the state process is univariate. Such a SSM is characterized by two processes (almost identical to that of the discrete-valued hidden Markov model):

1. A continuous-valued hidden Markov state process,  $\{C_t\}_{t \in \mathbb{N}}$ .
2. A observed process,  $\{X_t\}_{t \in \mathbb{N}}$ , whose realizations are assumed to be conditionally independent, given the states.

Formally, for a density function  $f$ , the assumptions can be formalized as:

$$f_{C_t|\mathbf{C}^{(t-1)}}(c_t | \mathbf{c}^{(t-1)}) = f(c_t | c_{t-1}), \quad t = 2, 3, \dots,$$

$$f_{X_t|\mathbf{X}^{(t-1)}, \mathbf{C}^{(t)}}(x_t | \mathbf{x}^{(t-1)}, \mathbf{c}^{(t)}) = f(s_t | c_t), \quad t \in \mathbb{N}.$$

The only difference in models between that of a HMM and a SSM is that the Markov process  $\{C_t\}$  is continuous-valued in the latter. However, As we will discretize the state space into a sufficiently large but finite number of states, we can evaluate an approximation of the likelihood of any given SSM, exactly like the discrete-valued HMM.

The discretization procedure is as follows: For some given SSM, we consider an essential range  $[b_0, b_m]$  of possible values of  $C_t$ . This range is then subdivided into  $m$  subintervals  $B_i = (b_{i-1}, b_i)$ ,  $i = 1, \dots, m$ . These subintervals need not be on a equidistant grid, however, for simplicity and computational ease, we assume they are equidistant. As such, they are all of the length  $h = (b_m - b_0)/m$ . Denote  $b_i^*$  a representative point in  $B_i$ , for example the midpoint. By making use of the SSM dependence structure and repeatedly approximating integrals  $\int_a^b f(c) dc$  by simple expressions of the form  $(b - a)f(c^*)$ , the likelihood of the observations  $\mathbf{x}^{(T)}$  can be approximated as follows

$$\begin{aligned} \mathcal{L}_T &= \underbrace{\int \dots \int}_{T-\text{integrals}} f_{\mathbf{X}^{(T)}, \mathbf{C}^{(T)}}(\mathbf{x}^{(T)}, \mathbf{c}^{(T)}) dc_T \dots dc_1 \\ &\stackrel{\dagger}{=} \int \dots \int f_{\mathbf{X}^{(T)}|\mathbf{C}^{(T)}}(\mathbf{x}^{(T)} | \mathbf{c}^{(T)}) f_{\mathbf{C}^{(T)}}(\mathbf{c}^{(T)}) dc_T \dots dc_1 \\ &\stackrel{\ddagger}{=} \int \dots \int f_{C_1}(c_1) f_{X_1|C_1}(x_1 | c_1) \prod_{t=2}^T f_{C_t|C_{t-1}}(c_t | c_{t-1}) f_{X_t|C_t}(x_t | c_t) dc_T \dots dc_1 \quad (32) \\ &\stackrel{\ddagger\ddagger}{\approx} \int_{b_0}^{b_m} \dots \int_{b_0}^{b_m} f_{C_1}(c_1) f_{X_1|C_1}(x_1 | c_1) \prod_{t=2}^T f_{C_t|C_{t-1}}(c_t | c_{t-1}) f_{X_t|C_t}(x_t | c_t) dc_T \dots dc_1 \\ &\stackrel{\ddagger\ddagger\ddagger}{\approx} h^T \sum_{i_1=1}^m \dots \sum_{i_T=1}^m f_{B_{i_1}}(b_{i_1}^*) f_{X_1|B_{i_1}}(x_1 | b_{i_1}^*) \prod_{t=2}^T f_{B_{i_t}|B_{i_{t-1}}}(b_{i_t}^* | b_{i_{t-1}}^*) f_{X_t|B_{i_t}}(x_t | b_{i_t}^*). \end{aligned}$$

$\dagger$  follows from definition of joint probability,  $\ddagger$  is simply rewriting into a product,  $\ddagger\ddagger$  is splitting the integrals into the essential range and  $\ddagger\ddagger\ddagger$  from the approximation, where the innermost

integral has been approximated as follows:

$$\int_{b_0}^{b_m} f_{C_T|C_{T-1}}(c_T | c_{T-1}) f_{X_T|C_T}(x_T | c_T) dc_T \approx h \sum_{i_T=1}^m f_{B_{i_T}|C_{T-1}}(b_{i_T}^* | c_{T-1}) f_{X_T|B_{i_T}}(x_T | b_{i_T}^*).$$

The terms appearing in the approximation in [Equation 32](#) are simple, however, the likelihood cannot be evaluated as is because of the extremely high number of summands ( $m^T$ ). However, the likelihood with a discrete state space, as the approximation is, yields a convenient form which allows us to employ our previously developed technique of using the forward algorithm to evaluate the likelihood<sup>7</sup>.

**Evaluation of the Approximate Likelihood** The discretization of the state space into some large number of intervals  $m$  corresponds to an approximation of the SSM by an  $m$ -state HMM. However, it is now possible to specify the components of this approximating HMM with ease. First, Consider the initial distribution of the state process. To obtain the exact expressions given in the last line of [Equation 32](#), we define the  $i$ 'th component of the  $m$ -dimensional vector  $\boldsymbol{\delta}$  to be  $\delta_i = hf(b_i^*)$ . then  $\delta_i$  is the approximate probability of the state process fallin in the interval  $B_i$  at time 1 (as it is the initial distribution). For example, assume that the state process is in its stationary distribution at the time of the first observation. Then  $f(b_i^*)$  is the density of the normal distribution evaluated at  $b_i^*$  with  $\mathbb{E}[C_t^*] = 0$  an  $\mathbb{V}[C_t^*] = \sigma_\varepsilon^2 / (1 - \rho^2)$ . In the exact same manner, define an  $m \times m$  t.p.m  $\boldsymbol{\Gamma} = (\gamma_{ij})_{i,j=1}^m$  by specifying  $\gamma_{ij} = hf(b_j^* | b_i^*)$ . The transition probabilities  $\gamma_{ij}$  are the approximate probability of the value of the state process falling into the intervals  $B_j$  at time  $t$  given that the process is in interval  $B_i$  at time  $t - 1$ . For the Gaussian AR(1) state process, the values of  $\gamma_{ij}$  is  $h$  times the density of the normal distribution with mean  $\rho b_i^*$  and variance  $\sigma_\varepsilon^2$  evaluated at  $b_j^*$ . Lastly, we define the component  $\boldsymbol{P}(x_t)$  to be the  $m \times m$  diagonal matrix with  $i$ th entry corresponding ot  $f(x_t | b_i^*)$ . This corresponds to an approximation of the conditional density of  $x_t$  given that the state process takes some value in the interval  $B_i$  at time  $t$ .

Assembling the components just defined, we can rewrite the multiple-sum expression for the apprximate likelihood given in [Equation 32](#) in the form of a matrix product

$$\begin{aligned} & h^T \sum_{i_1=1}^m \dots \sum_{i_T=1}^m f_{B_{i_1}}(b_{i_1}^*) f_{X_1|B_{i_1}}(x_1 | b_{i_1}^*) \prod_{t=2}^T f_{B_{i_t}|B_{i_{t-1}}}(b_{i_t}^* | b_{i_{t-1}}^*) f_{X_t|B_{i_t}}(x_t | b_{i_t}^*) \\ &= \boldsymbol{\delta} \boldsymbol{P}(x_1) \boldsymbol{\Gamma} \boldsymbol{P}(x_2) \boldsymbol{\Gamma} \boldsymbol{P}(x_3) \dots \boldsymbol{\Gamma} \boldsymbol{P}(x_{T-1}) \boldsymbol{\Gamma} \boldsymbol{P}(x_T) \mathbf{1}_N^\top. \end{aligned}$$

**Estimation Issues and Assessment** According to [59, p. 160], numerical maximization of the likelihood given in equation [Equation 32](#) is feasible even when the observation count and  $m$  is fairly large, which is what would be required for a relatively close approximation to the likelihood.

---

<sup>7</sup>This is not the only possible discretization. We choose a simple midpoint quadrature as it is computationally efficient.

In general, it seems to be that values around  $m = 50$  stabilize [59, 29]. Furthermore, as can easily be seen, the number of parameters does not depend on the magnitude of  $m$ . The entries of the approximate tpm  $\boldsymbol{\Gamma} \in \mathbb{R} \times \mathbb{R}$ , depends only on state process parameters of the SSM. The range  $[b_0, b_m]$  has to be chosen such that it is sufficiently large to cover the essential range of the state process. However, if it is chosen too large, the do not maintain sufficient fineness of the grid.

The other issues remain the same as for the  $N$ -state HMM as we are essentially fitting a  $m$ -state hidden Markov model; Local maxima, parameter constraints and especially, numerical under- and overflow. We use the exact same techniques as described for the HMM to deal with the latter issues of numerical under- and overflow.

As for the HMM, we extract the numerically estimated Hessian of the log-likelihood for the estimated parameters using the base R function `nlm`. Furthermore, we can use the Viterbi algorithm for decoding, pseudo-residuals and forecasts, exactly as for the HMM.

The BS-SSM $_{\beta}$  specified in (28) is parameterised by  $\boldsymbol{\zeta} = (\rho, \sigma_Z, \mu_0, \sigma, \beta_\mu, \beta_\sigma)$ , where  $\rho$  and  $\sigma_Z$  govern the latent AR(1) factor  $C_t$  and  $(\mu_0, \sigma, \beta_\mu, \beta_\sigma)$  determine the time-varying drift  $\mu_t$  and volatility  $\sigma_t$  of returns. This parametrisation is not globally identifiable. In particular, the latent factor and its loadings are only determined up to a scale and sign transformation. For any  $a > 0$  one can replace  $(\sigma_Z, \beta_\mu, \beta_\sigma) \mapsto (a \sigma_Z, \beta_\mu/a, \beta_\sigma/a)$ , while keeping  $(\rho, \mu_0, \sigma)$  fixed and obtain the same implied processes  $(\mu_t, \sigma_t)$  and hence the same law for the observed returns  $\{X_t\}$ . A further sign invariance arises by simultaneously mapping  $C_t \mapsto -C_t$  and  $(\beta_\mu, \beta_\sigma) \mapsto (-\beta_\mu, -\beta_\sigma)$ , which again leaves  $(\mu_t, \sigma_t)$  and the distribution of  $\{X_t\}$  unchanged. As a consequence, the individual parameters  $(\sigma_Z, \beta_\mu, \beta_\sigma)$  and the absolute scale and orientation of  $C_t$  are not uniquely identified from returns alone.

Crucially, this lack of full identifiability does not affect the quantities used in the empirical analysis. The maximised log-likelihood depends only on the distribution of  $\{X_t\}$  and is invariant along the above reparameterisation ridge, so AIC and BIC comparisons between the BSM, BS-HMM and BS-SSM $_{\beta}$  are well defined. Likewise, conditional forecast distributions  $F_{X_{t+h}|\mathcal{F}_t}$ , point forecasts (mixture means) and forecast intervals are functions of the implied drift and volatility paths  $(\mu_t, \sigma_t)$  and are therefore unchanged by any scale or sign transformation of  $(C_t, \sigma_Z, \beta_\mu, \beta_\sigma)$  that preserves  $(\mu_t, \sigma_t)$ . The same argument applies to pseudo-residuals, which are constructed from the conditional CDF of  $X_t$  under the fitted model and thus depend only on the forecast distribution of returns, not on a particular normalisation of the latent factor. All likelihood-based and distribution-based diagnostics reported here (AIC/BIC, pseudo-residuals, forecast-based MSE/RMSE/MAE) are therefore invariant to this non-identifiability and can be interpreted without qualification.

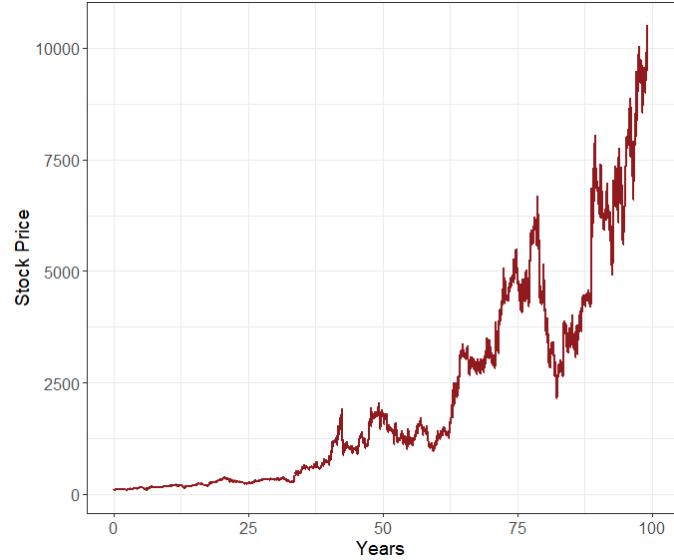
Some components of  $\boldsymbol{\zeta}$  remain identifiable in the usual sense. The AR(1) persistence  $\rho$ , the baseline drift  $\mu$  and the baseline volatility  $\sigma$  are uniquely pinned down by the dynamics of returns up to the usual economic decomposition of dividend depended drift. What is not uniquely identified is the internal decomposition of time-variation in  $(\mu_t, \sigma_t)$  into “how volatile the factor is”

$(\sigma_Z)$  versus “how strongly it loads into drift and volatility”  $(\beta_\mu, \beta_\sigma)$ . In this thesis,  $C_t$  is therefore interpreted only as a *relative* latent index of the return environment (high vs. low regimes over time), rather than as a factor with an intrinsically meaningful absolute scale. Statements about the timing and ordering of regimes (e.g. that periods with high  $C_t$  coincide with elevated volatility and lower conditional drift) are invariant and can be interpreted, whereas structural claims about the absolute magnitude of  $(\sigma_Z, \beta_\mu, \beta_\sigma)$  are not. Since the main focus is on model fit, in- and out-of-sample errors for returns and all of these are invariant to the latent-factor reparameterisation, the lack of full parameter identifiability is not a practical limitation for the empirical conclusions drawn here.

**Simulation** We first simulate the BS-SSM in Equation 27, that is, a Black–Scholes model with an AR(1) latent state  $C_t = \rho C_{t-1} + \varepsilon_t$  with innovations  $\varepsilon_t \sim \mathcal{N}(0, \sigma_\varepsilon^2)$  that multiplicatively scales both drift and volatility in the log-return formulation  $X_t$  (cf. Equation 4). The baseline parameters are  $\mu = 0.05$ ,  $\sigma = 0.15$ ,  $\rho = 0.98$  (high persistence),  $\sigma_\varepsilon = 0.10$ , with  $S_0 = 100$ ,  $n = 25000$  and daily observations  $\Delta = 1/252$  (approximately  $25000/252 \approx 99$  years). The simulated price path is shown in Figure 11 and estimated versus true parameters are reported in Table 3. We use `nlm` in R to maximise the discretised likelihood in Equation 32 for the BS-SSM.

We then consider the BS-SSM $_\beta$  specification in Equation 28, where the same AR(1) latent state enters the drift and log-volatility via linear factor loadings,  $\mu_t = \mu_0 + \beta_\mu C_t$  and  $\sigma_t = \sigma \exp(\beta_\sigma C_t)$ . For the simulation we choose  $\mu_0 = 0.05$ ,  $\beta_\mu = 0.20$ ,  $\sigma = 0.15$ ,  $\beta_\sigma = 0.60$ , together with  $\rho = 0.98$ ,  $\sigma_\varepsilon = 0.10$ ,  $S_0 = 100$ ,  $n = 25000$  and  $\Delta = 1/252$ . The corresponding price path is shown in Figure 12 and Table 4 collects true and estimated parameters.

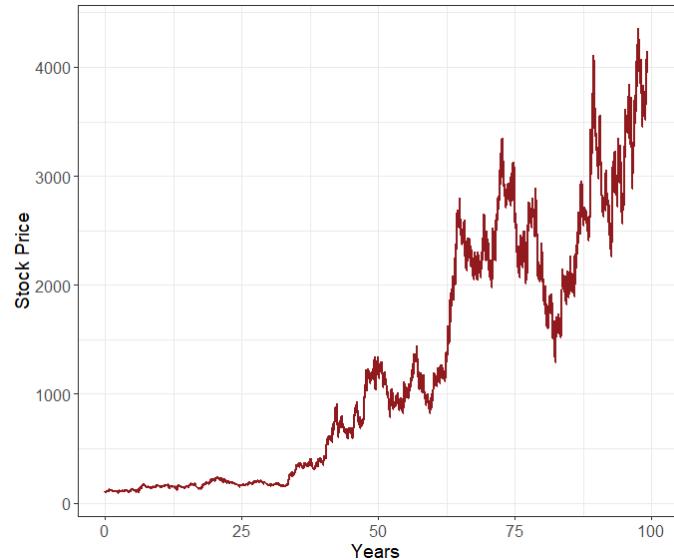
In both cases, the relative estimation errors for are small and do not raise concerns, bearing in mind the approximation error from the Euler discretisation and grid-based likelihood. The only stand-out parameter estimate is the loading factor  $\beta_\mu$  in the BS-SSM $_\beta$  that is somewhat less precisely estimated. This is consistent with expectation given that it primarily affects the conditional mean rather than the volatility, which we know to be prone to erroneous estimation.



**Figure 11:** BS-SSM simulated stock path.

Parameter	True Value	Estimated Value	Relative Error (%)
$\rho$	0.98000	0.98039	0.04
$\sigma_\varepsilon$	0.10000	0.09967	0.33
$\mu$	0.05000	0.05062	1.24
$\sigma$	0.15000	0.14653	2.31

**Table 3:** True vs. estimated parameters for the BS-SSM, including relative estimation errors.



**Figure 12:** BS-SSM $_\beta$  simulated stock path.

Parameter	True Value	Estimated Value	Relative Error (%)
$\rho$	0.98000	0.98082	0.08
$\sigma_\varepsilon$	0.10000	0.10006	0.06
$\mu_0$	0.05000	0.05062	1.24
$\sigma$	0.15000	0.14814	1.24
$\beta_\mu$	0.20000	0.22590	12.95
$\beta_\sigma$	0.60000	0.59251	1.25

**Table 4:** True vs. estimated parameters for the BS-SSM $_\beta$ , including relative estimation errors.

## 2.4 Model Selection Criteria & Assessment

### 2.4.1 Information Criteria: AIC & BIC

Two of the most popular approaches to model selection for HMMs will be used: The Akaike Information Criterion (AIC) and The Bayesian Information Criterion (BIC). These are supplementary methods to those discussed previously.

Assume that  $x_1, \dots, x_T$  were generated by the true data generating process,  $f$  and that one is interested in determining which model to choose among two different approximating families  $\{g_1 \in \mathcal{G}_1\}$  and  $\{g_2 \in \mathcal{G}_2\}$  under some criteria of being "the best". We thus need some operator to determine the lack of fit between the true data generating model and the fitted models,  $\Delta(f, \hat{g}_1)$  and  $\Delta(f, \hat{g}_2)$ . An immediate issue that arises is the lack of knowledge of  $f$ . As such, we can not determine from this discrepancy which model to select. However, we can use model selection criteria,  $\widehat{\mathbb{E}}_f[\Delta(f, \hat{g}_1)]$  and  $\widehat{\mathbb{E}}_f[\Delta(f, \hat{g}_2)]$ . These quantities bases selection on estimators of the expected discrepancies. The model selection criterion simplifies to the Akaike information criterion [59, p. 98] which, briefly stated, arises of the Kullback–Leibler discrepancy and conditions listed in [31, Appendix A]:

$$\text{AIC} = \underbrace{-2 \log \mathcal{L}_T}_{\text{measure of fit}} + \underbrace{2p}_{\text{penalty}}, \quad (33)$$

where  $\mathcal{L}$  is the log-likelihood of the fitted model and  $p$  denotes the number of parameters of the model<sup>8</sup>. It is immediately clear that increasing the number of parameters, by increasing the number of states or state-dependent parameters, will penalize the AIC. To compare model performances in terms of AIC, we follow [9, pp. 270–272] to some degree; Let  $\Delta i$  denote the difference in AIC between the best model (i.e. smallest AIC) and the one of comparison. The rule of thumb then states that we can assess the relative merits of models by:

- $\Delta i \leq 2 \Rightarrow$  Substantial support (evidence).
- $4 \leq \Delta i \leq 7 \Rightarrow$  Considerably less support (evidence).
- $\Delta i > 10 \Rightarrow$  Essentially no support (evidence).

Note, that [10] relaxed the rule of thumb and thus  $2 \leq \Delta i \leq 7$  have some support and should seldom be disregarded. However, this is not sufficient for model assessment as discussed in [Section 2.2.7](#).

Another approach to model selection is the Bayesian philosophy. The Bayesian philosophy to model selection differs slightly to the AIC approach. The Bayesian philosophy is to select the family which is estimated to be most likely to be true. Consistent with the Bayesian paradigm,

---

<sup>8</sup>see [Section 2.2.7](#) for number of parameter determination.

in the first step before considering observations at hand, one specifies the prior probabilities, that  $f$  stems from the approximating families  $\mathcal{G}_1, \mathcal{G}_2$ , namely,  $\mathbb{P}(f \in \mathcal{G}_1)$  and  $\mathbb{P}(f \in \mathcal{G}_2)$ . Secondly, one computes and compares the posterior probabilities that  $f$  belongs to the approximating families given the observations, namely,  $\mathbb{P}(f \in \mathcal{G}_1 | \mathbf{X}^{(T)} = \mathbf{x}^{(T)})$  and  $\mathbb{P}(f \in \mathcal{G}_2 | \mathbf{X}^{(T)} = \mathbf{x}^{(T)})$ . Briefly stated, under conditions seen in [57], the Bayesian information criterion arises [59, p. 98]:

$$\text{BIC} = \underbrace{-2 \log \mathcal{L}_T}_{\text{measure of fit}} + \underbrace{p \log T}_{\text{penalty}}, \quad (34)$$

where  $\mathcal{L}_T$  and  $p$  are as for the AIC and  $T$  is the number of observations, which is obviously not present whatsoever for the AIC. Compared to the AIC, the penalty term of the BIC has more weight for  $T > e^2$ , which holds in most practical applications. Thus, the BIC does, in general, favour models with fewer parameters than the AIC.

Summarizing, in both cases, the best model in the family is the one that minimizes these information criteria. Clearly, AIC does not depend directly on the sample size,  $T$ . Moreover, AIC presents the danger that it might overfit, whereas BIC presents the danger that it might underfit, simply in virtue of how each criterion penalize free parameters (see the under-braced penalty-terms in [Equation 33](#) and [Equation 34](#)).

#### 2.4.2 Pseudo-Residuals

Even after selecting what seems to be the best model according to some chosen criterion, it is still necessary to determine whether the model actually provides a good fit to the data. This requires tools that can assess the overall adequacy of the model and help identify potential outliers. In classical settings such as normal-theory regression, residuals are a well-known and widely used method for checking model fit. In this section, we introduce quantities called pseudo-residuals (also referred to as quantile residuals), which extend this idea to more general models and serve a similar purpose in the context of HMMs. We present two types of pseudo-residuals, both of which rely on the ability to compute likelihoods efficiently, something that HMMs naturally allow. The theory presented is based on that of [15, pp. 236–244] which they note is a special case of Cox–Snell residuals [13].

Each  $X_t$  has a distribution that depends on some latent state in the state space. As such, assessing outliers or model fit is non-trivial, since the conditional distribution of each  $X_t$  changes over time and depends on the hidden state sequence. A commonly used approach in HMMs for assessment of model fit is to transform the observations to a common scale using pseudo-residuals  $\{z_t\}_{t=1}^T$ , constructed via the probability integral transform [59, pp. 101–106]:

- I. Transform a observation  $x_t$  to  $u_t = F_{X_t}(x_t) \sim \mathcal{U}[0, 1]$ , where  $F_{X_t}$  is the CDF of  $X_t$ .
- II. Transform  $u_t$  to  $z_t = \Phi^{-1}(u_t) \sim \mathcal{N}(0, 1)$ , where  $\Phi$  is the standard normal CDF.

**III.** If the model is correctly specified, then the pseudo-residuals,  $z_t = \Phi^{-1}(F_{X_t}(x_t))$ , should be approximately independent and standard normally distributed. These can be evaluated using histograms and Q–Q plots.

We show the properties in **I.** and **II.** to be the case.

**Proposition 2.7.** *Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space and let  $X : \Omega \rightarrow \mathbb{R}$  be a real-valued r.v. with cumulative distribution function  $F_X$ . Suppose  $F_X$  is continuous and strictly increasing (hence invertible with inverse  $F_X^{-1}$ ). Define*

$$U := F_X(X) \quad \text{and} \quad Z := \Phi^{-1}(U),$$

where  $\Phi$  is the standard normal distribution function. Then  $U \sim \mathcal{U}[0, 1]$  and  $Z \sim \mathcal{N}(0, 1)$ .

*Proof.* First, for any  $u \in [0, 1]$ ,

$$\begin{aligned} \mathbb{P}(U \leq u) &= \mathbb{P}(F_X(X) \leq u) \\ &\stackrel{\dagger}{=} \mathbb{P}\left(X \leq F_X^{-1}(u)\right) \\ &= F_X(F_X^{-1}(u)) \\ &= u, \end{aligned}$$

where  $\dagger$  follows since  $F_X$  is strictly increasing. Next, for any  $z \in \mathbb{R}$ ,

$$\begin{aligned} F_Z(z) &= \mathbb{P}(Z \leq z) \\ &= \mathbb{P}(\Phi^{-1}(U) \leq z) \\ &= \mathbb{P}(U \leq \Phi(z)) \\ &= F_U(\Phi(z)) \\ &= \Phi(z). \end{aligned}$$

Therefore  $Z \sim \mathcal{N}(0, 1)$ . □

**Ordinary Pseudo-Residuals** The first approach examines each observation individually, identifying those that appear unusually extreme relative to the model and the remaining data in the series, indicating that they may differ in nature or origin. In practice, this involves computing a pseudo-residual  $\{z_t\}_{t=1}^T$  from the conditional distribution of  $X_t | \mathbf{X}^{(-t)}$ . For continuous observations the normal pseudo-residual is

$$z_t = \Phi^{-1} \left( \mathbb{P} \left( X_t \leq x_t | \mathbf{X}^{(-t)} = \mathbf{x}^{(-t)} \right) \right). \quad (35)$$

$z_t$  is approximately a realization of a standard normal r.v., if the model is correctly specified. Using results from [Section 2.2.6](#), specifically by integrating over [Equation 22](#) and [Equation 23](#), the ordinary pseudo-residuals can be calculated. However, we will use forecast pseudo-residuals for the assessment of models. This is because we aren't inherently interested in idiosyncratic outliers relative to the model. However, we are interested in the models predictive powers. As such, forecast pseudo-residuals are of advantageous use.

**Forecast Pseudo-Residuals** The second approach to detecting outliers focuses on identifying observations that appear unusually extreme when compared to what the model predicts based on all previous observations, rather than the entire data sequence. Here, the key quantity is the conditional distribution of  $X_t$  given  $\mathbf{X}^{(t-1)}$ . For continuous observations the pseudo-residual is

$$z_t = \Phi^{-1} \left( \underbrace{\mathbb{P} \left( X_t \leq x_t \mid \mathbf{X}^{(t-1)} = \mathbf{x}^{(t-1)} \right)}_{(*)} \right). \quad (36)$$

To compute the forecast pseudo-residuals, we evaluate the one-step-ahead forecast distribution of  $X_t$  given the observed history up to time  $t - 1$  by examining  $(*)$  in [Equation 36](#). Note that

$$\begin{aligned} \mathbb{P} \left( X_t \leq x_t \mid \mathbf{X}^{(t-1)} = \mathbf{x}^{(t-1)} \right) &= \sum_{j \in \mathcal{C}} \mathbb{P} \left( X_t \leq x_t, C_t = j \mid \mathbf{X}^{(t-1)} = \mathbf{x}^{(t-1)} \right) \\ &\stackrel{\dagger}{=} \sum_{j \in \mathcal{C}} \underbrace{\mathbb{P} \left( X_t \leq x_t \mid \mathbf{X}^{(t-1)} = \mathbf{x}^{(t-1)}, C_t = j \right)}_{:= F_{X_t,j}(x_t)} \underbrace{\mathbb{P} \left( C_t = j \mid \mathbf{X}^{(t-1)} = \mathbf{x}^{(t-1)} \right)}_{:= \psi_t(j)} \\ &\stackrel{\ddagger}{=} \sum_{j \in \mathcal{C}} \psi_t(j) F_{X_t,j}(x_t), \end{aligned}$$

where  $\dagger$  follows from the Law of Total Probability and  $\ddagger$  from the HMM conditional-independence assumption. In short, in the HMM setting, this forecast distribution is a mixture of state-dependent conditional distributions with  $\psi_t(j) = \mathbb{P} \left( C_t = j \mid \mathbf{X}^{(t-1)} = \mathbf{x}^{(t-1)} \right)$  the one-step-ahead predicted state probabilities and for  $q, p \in \mathbb{R}$

$$F_{X_t,j}(x_t) = \Phi \left( \frac{x_t - m_j}{s_j} \right),$$

with  $m_i := (\mu_i - \frac{1}{2}\sigma_i^2)\Delta$  and  $s_i^2 = \sigma_i^2\Delta$ . The predicted state probabilities are obtained by propagating them forward using the transition probabilities and the normalized vector of forward

variables

$$\begin{aligned}
\psi_t(j) &= \mathbb{P}\left(C_t = j \mid \mathbf{X}^{(t-1)} = \mathbf{x}^{(t-1)}\right) \\
&= \sum_{i \in \mathcal{C}} \underbrace{\mathbb{P}(C_t = j \mid C_{t-1} = i)}_{:= \gamma_{i,j}} \underbrace{\mathbb{P}(C_{t-1} = i \mid \mathbf{X}^{(t-1)} = \mathbf{x}^{(t-1)})}_{:= \phi_{t-1}(j)} \\
&= \sum_{i \in \mathcal{C}} \gamma_{ij} \phi_{t-1}(i) \\
&\Rightarrow \\
\psi_t &= \boldsymbol{\phi}_{t-1} \boldsymbol{\Gamma}.
\end{aligned}$$

For numerical stability we use exponentiated normalized forward probability vectors in the code implementation. Specifically, in the BS-HMM, the state-dependent cumulative distributions  $F_{X_{t,j}}(x_t)$  are governed by the Gaussian distribution

$$X_t \mid \{C_t = i\} \sim \mathcal{N}\left(\left(\mu_i - \frac{1}{2}\sigma_i^2\right)\Delta, \sigma_i^2\Delta\right).$$

The pseudo-residuals are then given by

$$z_t = \Phi^{-1} \left( \sum_{j \in \mathcal{S}} \psi_t(j) \cdot F_{X_{t,j}}(x_t) \right), \quad t = 2, \dots, T.$$

For the first residual  $z_1$ , the one-step-ahead state probabilities  $\psi_1(j)$  cannot be computed from previous forward probabilities, since there is no observation before  $X_1 = x_1$ . A convenient circumvention for this complication is to approximate them using the stationary distribution  $\boldsymbol{\delta}$  which represents the long-run state probabilities of the Markov chain. Thus, the first residual is computed as

$$z_1 = \Phi^{-1} \left( \sum_{j \in \mathcal{S}} \delta_j \cdot F_{X_{t,j}}(x_t) \right).$$

If the model is correctly specified, the pseudo-residuals  $\{z_t\}_{t=1}^T$  should be approximately independent and standard normally distributed.

### 3 Data (II of II)

Throughout the empirical analysis we work with logarithmic returns rather than price levels. Let  $S_t$  denote the closing level of the S&P 500 at trading day  $t$ . The one-day log return over some fixed time horizon  $t = 1, 2, \dots, T$  is

$$X_t = \log S_t - \log S_{t-1} = \log\left(\frac{S_t}{S_{t-1}}\right) \quad (37)$$

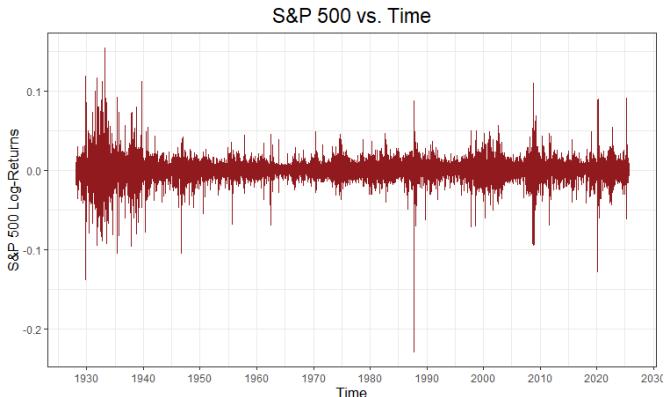
This transformation is standard in financial econometrics for several reasons.

First, price levels for broad equity indices are well known to be non-stationary and to exhibit unit-root behavior over long samples, while their first differences, i.e. returns, are much closer to weak stationarity. Stationarity (or approximate stationarity) is a prerequisite for a wide class of likelihood-based time-series methods (ARMA, GARCH and state-space models). Working in returns mitigates spurious regression phenomena associated with trending levels and yields series with stable mean near zero and finite variance, albeit with conditional heteroskedasticity. However, we again note, we never state that we model the raw S&P 500 data as an autoregressive process but rather log-normal as we use the BSM.

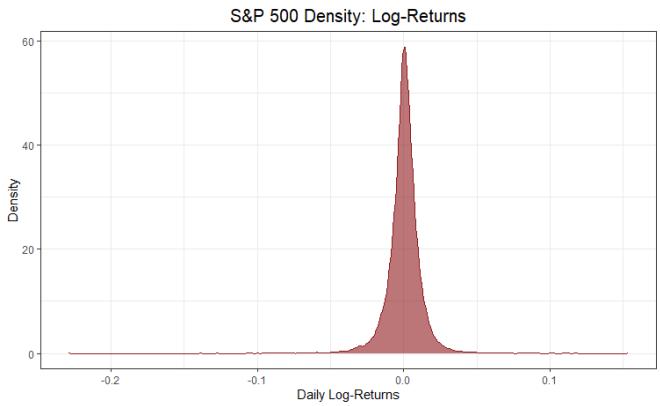
Second, log returns aggregate additively across time. For any horizon  $h \in \mathbb{N}$ ,

$$X_{t:t+h} = \log\left(\frac{S_{t+h}}{S_t}\right) = \sum_{j=1}^h X_{t+j}. \quad (38)$$

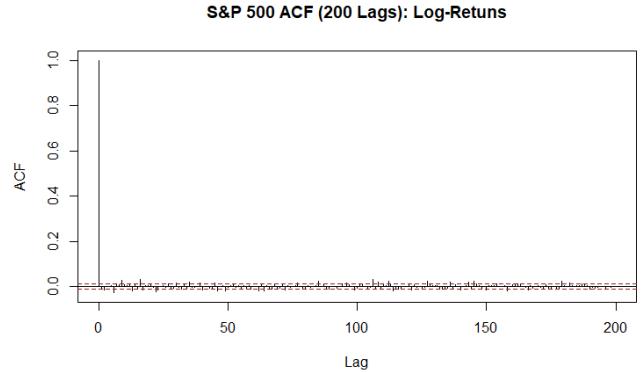
Additivity greatly simplifies multi-horizon likelihoods, forecasting and decomposition of long-horizon performance into daily contributions. By contrast, simple (arithmetic) returns compound multiplicatively. For high-frequency horizons where  $|X_t| = (S_t - S_{t-1})/S_{t-1}$  is small,  $\log(1 + X_t) \approx X_t$ , so log and simple returns coincide to first order while preserving exact additivity in Equation 38.



**Figure 13:** S&P 500 index time series of log-returns.



**Figure 14:** Kernel density of S&P 500 log-returns.



**Figure 15:** Autocorrelation of S&P 500 log-returns (200 lags).

Third, log returns are scale-free and invariant to rescaling of the numeraire. If prices are multiplied by any constant  $c > 0$  (e.g., index rebasings), the difference of logs in [Equation 37](#) is unchanged. This facilitates cross-asset and cross-period comparisons and stabilizes variance relative to the price level.

Finally, using log returns improves numerical stability in estimation. Likelihood functions built on price levels inherit the explosive scale and collinearity of  $S_t$ , whereas those in differences avoid ill-conditioning and reduce sensitivity to arbitrary index base values.

To examine the data after the transformation we examine the same descriptive statistics and plots. The time series plot can be seen in [Figure 13](#). As claimed, the data are now scale-free and returns do not compound multiplicatively as with the non-transformed data.

[Figure 14](#) shows the density of the log-returns. It is apparent that the data are unimodal, symmetric and approximately normal with low variance and approximately zero mean. Indeed, the first quartile is  $-0.0045519$ , the median is  $0.0004940$  and the third quartile is  $0.005458$ . The largest observation is  $0.1536613$  and the smallest is  $-0.2289973$ .

Including 200 lags (200 closing trading days), the S&P 500 log-returns ACF is shown in [Figure 15](#). The autocorrelation is notably reduced by the simple transformation.

Because of the issues related to raw prices and the advantages of using the transformed data, we use log-returns for the remainder of the thesis, but we do report MLEs on both raw and transformed data for comparison.

**Estimating a State-Wise Dividend Yield  $\hat{q}_i$**  As detailed in [Section 1](#), we construct a daily log-dividend series

$$q_t^{(\log)} = \log \frac{T_t}{T_{t-1}} - \log \frac{P_t}{P_{t-1}},$$

from total-return  $T_t$  and price  $P_t$ . Because the most likely state sequence is not generally ordered in consecutive blocks (states  $i$  and  $j$  may interleave arbitrarily), we cannot form state-wise averages from single contiguous subsamples. Instead, we use the Viterbi algorithm to obtain the most probable state sequence  $(i_1^*, \dots, i_T^*)$ . For each state  $i$ , define the index set

$$\mathcal{T}_i := \{t \in \{2, \dots, T\} : i_t^* = i\}, \quad N_i := |\mathcal{T}_i|.$$

Conditional on state  $i$ , we model  $\{q_t^{(\log)} : t \in \mathcal{T}_i\}$  as a weakly dependent time series with constant mean

$$m_{q,i} := \mathbb{E} \left[ q_t^{(\log)} \mid i_t^* = i \right],$$

and write the intercept-only regression

$$q_t^{(\log)} = \alpha_i + u_{i,t}, \quad t \in \mathcal{T}_i,$$

where  $u_{i,t}$  is a zero-mean disturbance capturing short-run deviations of  $q_t^{(\log)}$  around the state- $i$  mean, with  $\mathbb{E}[u_{i,t} \mid i_t^* = i] = 0$ . The ordinary least squares estimator of  $\alpha_i$  is the sample mean

$$\hat{\alpha}_i := \bar{q}_i^{(\log)} = \frac{1}{N_i} \sum_{t \in \mathcal{T}_i} q_t^{(\log)},$$

which we take as the estimator of the daily state-wise mean log dividend yield  $m_{q,i}$ .

For reporting, we convert the daily mean in each state to an annualised continuous dividend yield by the linear rescaling

$$\hat{q}_i := 252 \hat{\alpha}_i.$$

For each state we then report the pair governing price dynamics and the total-return counterpart:

$$\hat{\mu}_{\text{cap},i}, \hat{\sigma}_i \quad \text{and} \quad \hat{\mu}_{\text{tot},i} = \hat{\mu}_{\text{cap},i} + \hat{q}_i,$$

but inference below is based on the capital-gains drifts  $\hat{\mu}_{\text{cap},i}$ .

As a benchmark, the constant-parameter Black–Scholes specification uses a single global dividend yield  $q$ , obtained by repeating the same construction on the full sample  $\mathcal{T}$  without conditioning on a particular state. Writing

$$q_t^{(\log)} = \alpha + u_t, \quad t \in \mathcal{T},$$

the global daily mean  $\hat{\alpha} := |\mathcal{T}|^{-1} \sum_{t \in \mathcal{T}} q_t^{(\log)}$  is estimated by OLS and the corresponding continuous annualised dividend yield is

$$\hat{q} := 252 \hat{\alpha}.$$

To quantify the estimation uncertainty of the dividend-yield estimators, we treat  $\{q_t^{(\log)}\}$  as a strictly stationary and ergodic time series with finite second moment and absolutely summable autocovariances

$$\Psi_h := \text{Cov}\left(q_t^{(\log)}, q_{t+h}^{(\log)}\right), \quad \sum_{h=-\infty}^{\infty} |\Psi_h| < \infty.$$

Let  $\mathcal{I}$  denote a generic index set of observation times (either the full sample  $\mathcal{T}$  for the global estimator, or one of the state-specific sets  $\mathcal{T}_i$ ) and write  $n_{\mathcal{I}} := |\mathcal{I}|$ . The corresponding sample mean is

$$\hat{\alpha}_{\mathcal{I}} := \frac{1}{n_{\mathcal{I}}} \sum_{t \in \mathcal{I}} q_t^{(\log)},$$

with population mean  $m_{q,\mathcal{I}} := \mathbb{E}\left[q_t^{(\log)}\right]$  under the relevant conditioning. Thus the point estimator is always the sample mean over the relevant index set  $\mathcal{I}$ .

Define the long-run variance

$$\Omega_{q,\mathcal{I}} := \sum_{h=-\infty}^{\infty} \Psi_{h,\mathcal{I}} = \Psi_{0,\mathcal{I}} + 2 \sum_{h=1}^{\infty} \Psi_{h,\mathcal{I}},$$

where  $\Psi_{h,\mathcal{I}}$  denotes the autocovariance at lag  $h$  for the process underlying the sample indexed by  $\mathcal{I}$  (for example, the state- $i$  process when  $\mathcal{I} = \mathcal{T}_i$ ). Under the above weak-dependence conditions (see, e.g., [2, 25]), a central limit theorem yields

$$\sqrt{n_{\mathcal{I}}}(\widehat{\alpha}_{\mathcal{I}} - m_{q,\mathcal{I}}) \xrightarrow{d} \mathcal{N}(0, \Omega_{q,\mathcal{I}}),$$

so that the asymptotic variance of  $\widehat{\alpha}_{\mathcal{I}}$  is  $\Omega_{q,\mathcal{I}}/n_{\mathcal{I}}$ .

In practice,  $\Omega_{q,\mathcal{I}}$  is unknown and we replace it by a Newey-West HAC estimator. Let

$$\hat{u}_t := q_t^{(\log)} - \widehat{\alpha}_{\mathcal{I}}, \quad t \in \mathcal{I},$$

denote the regression residuals and let  $\hat{\Psi}_{h,\mathcal{I}}$  be their sample autocovariance at lag  $h$ . For kernel weights  $\{w_h\}_{h=0}^{H_{n_{\mathcal{I}}}}$  (e.g., Bartlett) with a bandwidth  $H_{n_{\mathcal{I}}}$  satisfying  $H_{n_{\mathcal{I}}} \rightarrow \infty$  and  $H_{n_{\mathcal{I}}}/n_{\mathcal{I}} \rightarrow 0$ , the HAC estimator of the long-run variance is

$$\widehat{\Omega}_{q,\mathcal{I}} = \hat{\Psi}_{0,\mathcal{I}} + 2 \sum_{h=1}^{H_{n_{\mathcal{I}}}} w_h \hat{\Psi}_{h,\mathcal{I}},$$

and a consistent estimator of the asymptotic variance of  $\widehat{\alpha}_{\mathcal{I}}$  is

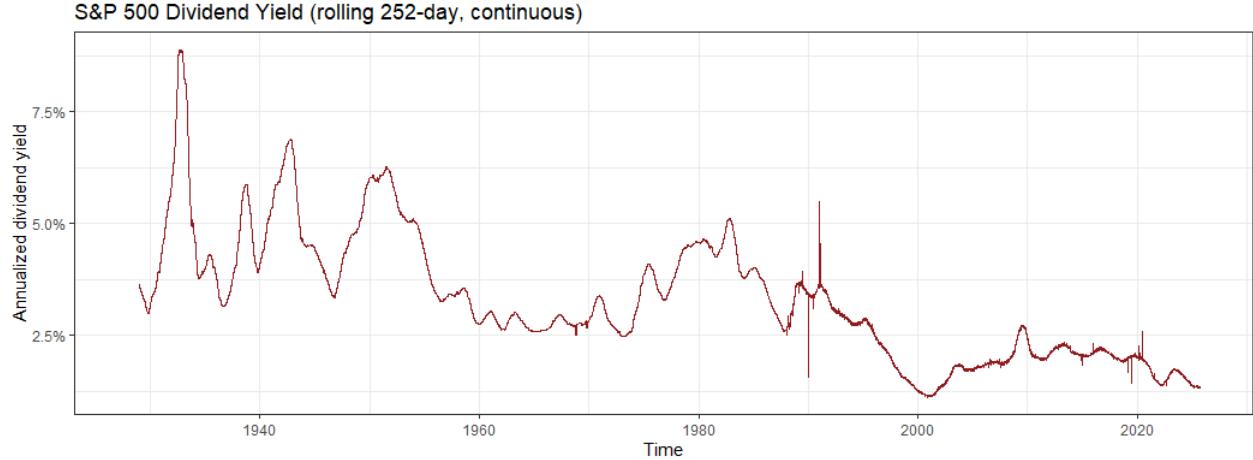
$$\widehat{\mathbb{V}}(\widehat{\alpha}_{\mathcal{I}}) = \frac{\widehat{\Omega}_{q,\mathcal{I}}}{n_{\mathcal{I}}}, \quad \widehat{\text{SE}}(\widehat{\alpha}_{\mathcal{I}}) = \sqrt{\frac{\widehat{\Omega}_{q,\mathcal{I}}}{n_{\mathcal{I}}}}.$$

Finally, the continuous annualised dividend yields and their standard errors are obtained by linear rescaling,

$$\widehat{q}_{\mathcal{I}} := 252 \widehat{\alpha}_{\mathcal{I}}, \quad \widehat{\text{SE}}(\widehat{q}_{\mathcal{I}}) := 252 \widehat{\text{SE}}(\widehat{\alpha}_{\mathcal{I}}).$$

The global estimator corresponds to  $\mathcal{I} = \mathcal{T}$ , while the state-wise estimators  $(\widehat{q}_i, \widehat{\text{SE}}(\widehat{q}_i))$  arise by taking  $\mathcal{I} = \mathcal{T}_i$ .

For the constant-dividend estimator  $\widehat{q}$  we also compute, as a benchmark, a naive i.i.d. standard error  $\widehat{\text{SE}}_{\text{iid}}(\widehat{q})$  based on the usual sample-variance formula that treats  $\{q_t^{(\log)}\}$  as independent. All reported confidence intervals in tables, however, are based on the HAC standard errors  $\widehat{\text{SE}}(\widehat{q}_{\mathcal{I}})$ .



**Figure 16:** Annualised dividend yield vs. time for the S&P 500, constructed from Shiller's monthly price–dividend data in the pre-TR period and from TR–PR differencing in the post-1988 period.

We visualize the estimated dividend yields in Figure 16. Over the full sample (1927–2025), the estimated average continuous annualised dividend yield is, to two significant digits, 3.3% under both construction methods. Using the HAC / Newey–West standard error, the standard error of the constant-dividend estimator  $\hat{q}$  is approximately  $9.2 \times 10^{-4}$ . The maximum of the estimated daily dividend yield is 8.9% on 2 November 1932, while the minimum is 1.1% on 27 November 2000. Dividend yields have trended down over the sample, consistent with a gradual shift from cash dividends towards share repurchases in U.S. equity markets.

## 4 Empirical Data Application

### 4.1 Model Selection & Assessment

In what follows, let  $p$  denote the number of estimated parameters of a given model.

**The Black-Scholes Model** The model criteria for the BSM is seen in [Table 5](#) and the residuals in [Figure A.3.1](#).

Black-Scholes Model (BSM)			
Model	$p$	AIC	BIC
BSM	2	-139353.758	-139337.662

*Table 5:* AIC and BIC for the standard BSM.

It is quite evident that the BSM yields extremely heavy tails. Heavy tails imply that large positive and negative returns occur with a much higher frequency than predicted by the Gaussian assumption in the BSM, so the model severely underestimates tail risk. Specifically, we notice from the time-series plot that the overly large residuals stem from [\[20\]](#):

- The Wall Street crash of 1929 & recession of 1937-1938.
- The early 1980s recession & Black Monday 1987.
- The Dot-com Bubble late 1990s & early 2000s.
- The 2008 Financial Crisis.

The autocorrelation for the returns is mostly unchanged and the model does still largely over- and underestimate returns outside of economical crises. These findings suggest that the BSM does not capture extreme economical environments adequately.

**Black-Scholes Hidden Markov Models** It is evident from [Table 6](#) that the best performing model in terms of the model information criteria, AIC and BIC, was the 5-state BS-HMM where we allow for state-dependency of the variables  $\sigma$  and  $\mu$ . Secondly is the 5-state BS-HMM where  $\sigma$  is allowed state-dependency but  $\mu$  is state-independent. Thirdly and last up for consideration, is the 4-state BS-HMM where we allow for state-dependency of the variables  $\sigma$  and  $\mu$ . In [Figure A.3.2](#) we examine the state-dependent density plots to examine for any unreasonable fittings that would point towards overfitting issues. Every model can be seen in [Section A.4](#).

Black-Scholes Hidden Markov Model (BS-HMM)			
Model	$p$	AIC	BIC
2-state BS-HMM ( $\mu$ )	5	-139348.0	-139308.0
3-state BS-HMM ( $\mu$ )	10	-139338.0	-139257.0
4-state BS-HMM ( $\mu$ )	17	-139324.0	-139187.0
5-state BS-HMM ( $\mu$ )	26	-139306.0	-139097.0
2-state BS-HMM ( $\sigma$ )	5	-150725.0	-150685.0
3-state BS-HMM ( $\sigma$ )	10	-152591.0	-152510.0
4-state BS-HMM ( $\sigma$ )	17	-153156.0	-153019.0
5-state BS-HMM ( $\sigma$ )	26	-153266.0	-153057.0
2-state BS-HMM ( $\sigma$ & $\mu$ )	6	-150746.0	-150698.0
3-state BS-HMM ( $\sigma$ & $\mu$ )	12	-152613.0	-152517.0
4-state BS-HMM ( $\sigma$ & $\mu$ )	20	-153199.0	-153038.0
5-state BS-HMM ( $\sigma$ & $\mu$ )	30	<b>-153342.0</b>	<b>-153101.0</b>

**Table 6:** AIC and BIC for fitted HMMs by state-dependent parameter family. The globally lowest (best) AIC and BIC are shown in bold.

Firstly, we examine state-dependent densities seen in Figure A.3.2 for the best performing model in terms of information criteria, 5-state BS-HMM where we allow for state-dependency of the parameters  $\sigma$  and  $\mu$ . We notice that states 2 and 3 are strikingly similar to state 2 for the 4-state BS-HMM where we allow for state-dependency of the parameters  $\sigma$  and  $\mu$ , which could suggest a case of overfitting. Furthermore, by examining the proportion of time spent in each state, it is almost exactly the sum of states 2 and 3 in the 4-state model that gives the time spent in state 2 for the 5-state model. However, it is advantageous for us to examine parameter estimates as well, which are seen in Table A.4.12.

As is now quite evident, the volatility parameter of states 2 and 3 for the 5-state BS-HMM with  $\mu$  and  $\sigma$  state-dependent, are almost identical for the first 3 digits at  $\sigma_2 \approx \sigma_3 \approx 0.11$ . This estimate is reasonable but does raise concern by the similarity. We turn our attention to  $\mu_2$  and  $\mu_3$ . Immediately we see that the estimates for  $\mu$  differ by quite a large margin, which is not a concern. However, both estimates are unreasonable. Especially,  $\mu_3 \approx -1.39$  emphasizes that the model for consideration provides unreasonable parameter estimates by overfitting to the data at hand. We therefore move our attention to the 4-state BS-HMM with  $\mu$  and  $\sigma$  state-dependent. Parameter estimates are seen in Table A.4.11.

The parameter estimates are reasonable, with no immediate outliers other than some large volatility parameters that could be problematic. We will present these in the next section and why they allow for a reasonable market interpretation.

**Black-Scholes Continuous State-Space Models** As described in Section 2.3.2, there is no exact result in choosing  $b_{\max}$  and  $m$  but only general experimental knowledge. As such, we firstly calculate the negative log-likelihood using a range and combination of  $b_{\max} \in \{0.5, 1, 2, 3, 4\}$  and

$m \in \{20, 40, 70, 100, 200\}$ . This is done to find the model yielding the smallest negative log-likelihood that also yields reasonable parameter estimates. The results are seen in [Table 7](#) and [Table 8](#). Models marked with a “—” yield parameter estimates that are unreasonable and suggest the model is too coarse.

$b_{\max}$	$m$				
	20	40	70	100	200
0.5	—	—	-74572.223296	-74568.761073	-74566.365696
1	—	-76247.016297	-76243.344459	-76242.462436	-76241.832368
2	—	-76635.707635	-76635.587101	-76635.558157	-76635.537345
3	—	—	<b>-76637.766845*</b>	-76637.766841*	-76637.766841*
4	—	—	—	-76637.766841*	-76637.766841*

**Table 7:** BS-SSM negative log-likelihoods rounded to 6 decimals; the overall minimum is in bold. Dashes indicate runs that didn’t converge or were too coarse. A \* marks values with identical 10 first digits (i.e., -76637.76684x), indicating they are essentially tied with the best model.

$b_{\max}$	$m$				
	20	40	70	100	200
0.5	—	—	—	-76709.211244	-76709.211244
1	—	—	-76709.211207	-76709.211205	-76709.211205
2	—	—	-76709.211239	-76709.211239	-76709.211239
3	—	—	—	-76709.211243	-76709.211243
4	—	—	—	<b>-76709.211773*</b>	-76709.211243

**Table 8:** BS-SSM $_{\beta}$  negative log-likelihoods rounded to 6 decimals; the overall minimum is in bold. Dashes indicate runs that didn’t converge or were too coarse.

It is evident the best BS-SSM model is  $m = 70$  and  $b_{\max} = 3$  and the best BS-SSM $_{\beta}$  is  $m = 100$  and  $b_{\max} = 4$ . Robustness checks across models show that the objective has converged, so the comparison is not sensitive to the precise grid choice. AIC & BIC comparison for the best BS-SSMs is seen in [Table 9](#).

Black-Scholes State-Space Models (BS-SSM)			
Model	$p$	AIC	BIC
BS-SSM	4	-153267.5	-153235.1
BS-SSM $_{\beta}$	6	<b>-153406.4</b>	<b>-153357.8</b>

**Table 9:** AIC and BIC for the BS-SSM and BS-SSM $_{\beta}$ . The globally lowest (best) AIC and BIC are shown in bold.

The BS-SSM $_{\beta}$  is clearly superior in terms of model criteria. Proceeding to the model assessment we examine pseudo-residuals which can be seen in [Figure A.3.7](#) and [Figure A.3.8](#). We start with

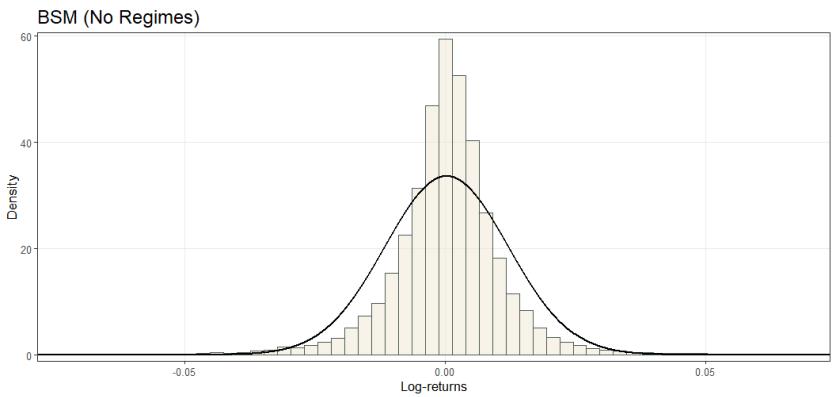
the former. The QQ-plot shows that the residuals are heavy-tailed and negatively skewed, with a particularly pronounced left tail relative to the assumed normal distribution. This implies that large negative residuals occur more frequently and are more extreme than the model predicts, so it systematically understates downside moves in returns and violates the normality assumption for the innovations. This is mostly similar to the BSM pseudo-residuals. Furthermore, the large residuals are present on the exact same days as the BSM suggesting a lack of fit on extreme changes in economical environments that deviates from the norm. However, a key difference is seen in the size of residuals. The BS-SSMs predicted returns are systematically too high, whereas they are largely symmetrical for the BSM. Pseudo-residuals of magnitude  $|5| >$  are not a rare occurrence. The large pseudo-residuals imply that one should continue with the model with extreme uncertainty and that conclusions drawn are most likely not adequate. We will discuss why the BS-SSM intuitively is not great for price modelling in [Section 5](#).

The  $\text{BS-SSM}_\beta$  does have a heavy left tail but the results are now magnitudes better than for the BS-SSM. Pseudo-residuals are mostly of adequate order with some outliers during the before-mentioned extreme periods. As such, we choose the  $\text{BS-SSM}_\beta$  for presentation.

## 4.2 Model Presentation

**Black-Scholes Model** Firstly, a presentation of the simple BSM will be given with parameter estimates seen in [Table 10](#) and density plot overlaid to the distribution of daily log-returns, using said parameters, in [Figure 17](#) the BSM specification fitted to daily log returns on the S&P 500 price index over our baseline sample period (see [Section 1](#) and [Section 3](#)). The likelihood is constructed under the assumption that log returns are conditionally i.i.d. Gaussian with constant drift and volatility and that the dividend yield enters as a constant proportional payout rate. All parameters are reported in annualised units and the standard errors and 95% confidence intervals are obtained from the inverse Hessian of the minimised log-likelihood.

The point estimate of the dividend yield is  $\hat{q} = 0.03297$  with a very tight 95% confidence interval constructed by the HAC SEs [0.03205, 0.03389], corresponding to an annual dividend yield of roughly 3.3%. The capital-gain drift of the price index is estimated as  $\hat{\mu}_{\text{cap}} = 0.07454$ , with a 95% confidence interval [0.03599, 0.11310] that lies entirely above zero. This implies a statistically and economically significant positive expected excess return on the index. Adding the dividend component, the implied total-return drift is  $\hat{\mu}_{\text{tot}} = 0.10751$ , corresponding to an expected annual total return of about 10.8%.



**Figure 17:** Histogram plot of daily log-returns and density of the BSM.

Parameter	Estimate (95% CI)	Std. Error
$\hat{q}$	0.03297 (0.03205, 0.03389)	0.0009171
$\hat{\mu}_{\text{cap}}$	0.07454 (0.03599, 0.11310)	0.01967
$\hat{\mu}_{\text{tot}}$	0.10751 (—, —)	—
$\hat{\sigma}$	0.18830 (0.18660, 0.19000)	0.0008761

**Table 10:** BSM parameter estimates with 95% CIs based on the inverse Hessian of the minimized log-likelihood.

The volatility estimate is  $\hat{\sigma} = 0.18830$  with a narrow 95% confidence interval [0.18660, 0.19000], indicating that the unconditional return volatility is tightly pinned down by the data. Interpreted in annual terms, the model implies an S&P 500 volatility of approximately 18.8%. Overall, the BSM benchmark thus corresponds to a relatively stable annual dividend yield around 3%, an expected total return around 11% and a volatility just below 20%. We use these constant-parameter estimates as a reference point for assessing the added flexibility and empirical fit of the regime-switching specifications developed in the following analysis.

**Black-Scholes Hidden Markov Model** We now present the 4-state BS–HMM with state-dependent drift and volatility  $(\mu_i, \sigma_i)$ . To aid interpretation, the main parameter estimates are reported in Table A.4.11 and, for readability, in Table 11, while the remaining parameters are collected in Appendix A.4. The MLEs indicate that the four latent states correspond to qualitatively distinct market phases, which can be broadly interpreted as two “bull” (expansion) states and two “bear” (contraction) states of varying intensity. States 1 and 2 are associated with positive expected price growth (capital gains), whereas states 3 and 4 exhibit negative expected price growth. The magnitude of the annualised drift and volatility in each state differs markedly, spanning a spectrum from a very tranquil bull market to a severe crisis regime. Dividend yields  $q_i$  vary relatively little across the first three states and rise noticeably only in the most adverse regime.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.2346 (0.1990, 0.2703)	0.01820
$\hat{\mu}_{cap,2}$	0.1062 (0.0621, 0.1502)	0.02246
$\hat{\mu}_{cap,3}$	-0.1031 (-0.2193, 0.0132)	0.05930
$\hat{\mu}_{cap,4}$	-0.3818 (-0.9260, 0.1623)	0.27763
$\hat{q}_1$	0.0306 (0.0286, 0.0327)	0.00104
$\hat{q}_2$	0.0336 (0.0316, 0.0355)	0.00099
$\hat{q}_3$	0.0306 (0.0279, 0.0333)	0.00139
$\hat{q}_4$	0.0498 (0.0413, 0.0583)	0.00432
$\hat{\mu}_{tot,1}$	0.2653 (—, —)	—
$\hat{\mu}_{tot,2}$	0.1398 (—, —)	—
$\hat{\mu}_{tot,3}$	-0.0725 (—, —)	—
$\hat{\mu}_{tot,4}$	-0.3320 (—, —)	—
$\hat{\sigma}_1$	0.0707 (0.0680, 0.0734)	0.00139
$\hat{\sigma}_2$	0.1270 (0.1229, 0.1311)	0.00210
$\hat{\sigma}_3$	0.2301 (0.2205, 0.2396)	0.00487
$\hat{\sigma}_4$	0.5651 (0.5329, 0.5974)	0.01646
$\hat{\Gamma}$	$\begin{pmatrix} 0.9649 (0.0044) & 0.0341 (0.0045) & 0.0002 (0.0006) & 0.0007 (0.0005) \\ 0.0207 (0.0030) & 0.9683 (0.0033) & 0.0110 (0.0015) & 0.0000 (0.0000) \\ 0.0000 (0.0000) & 0.0259 (0.0033) & 0.9632 (0.0039) & 0.0109 (0.0021) \\ 0.0000 (0.0000) & 0.0000 (0.0000) & 0.0517 (0.0094) & 0.9483 (0.0094) \end{pmatrix}$	
$\hat{\delta}$	$(0.2766 (0.0254), 0.4675 (0.0233), 0.2080 (0.0220), 0.0479 (0.0099))$	

**Table 11:** 4-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  and volatility  $\sigma_i$ . SEs are marked in parentheses for the HMM parameters for readability and space.

State 1 is characterised by an exceptionally high annualised capital-gains drift of  $\hat{\mu}_{cap,1} = 0.2346$ , implying roughly a 23.5% expected annual price appreciation. Including the state-specific dividend yield  $\hat{q}_1 \approx 0.0306$  (about 3.1% per year), the total-return drift increases to  $\hat{\mu}_{tot,1} = 0.2653$ , corresponding to an expected annual total return of about 26.5%. This is by far the most optimistic state. At the same time, the volatility is very low,  $\hat{\sigma}_1 = 0.0707$  (around 7.1% annually), with a

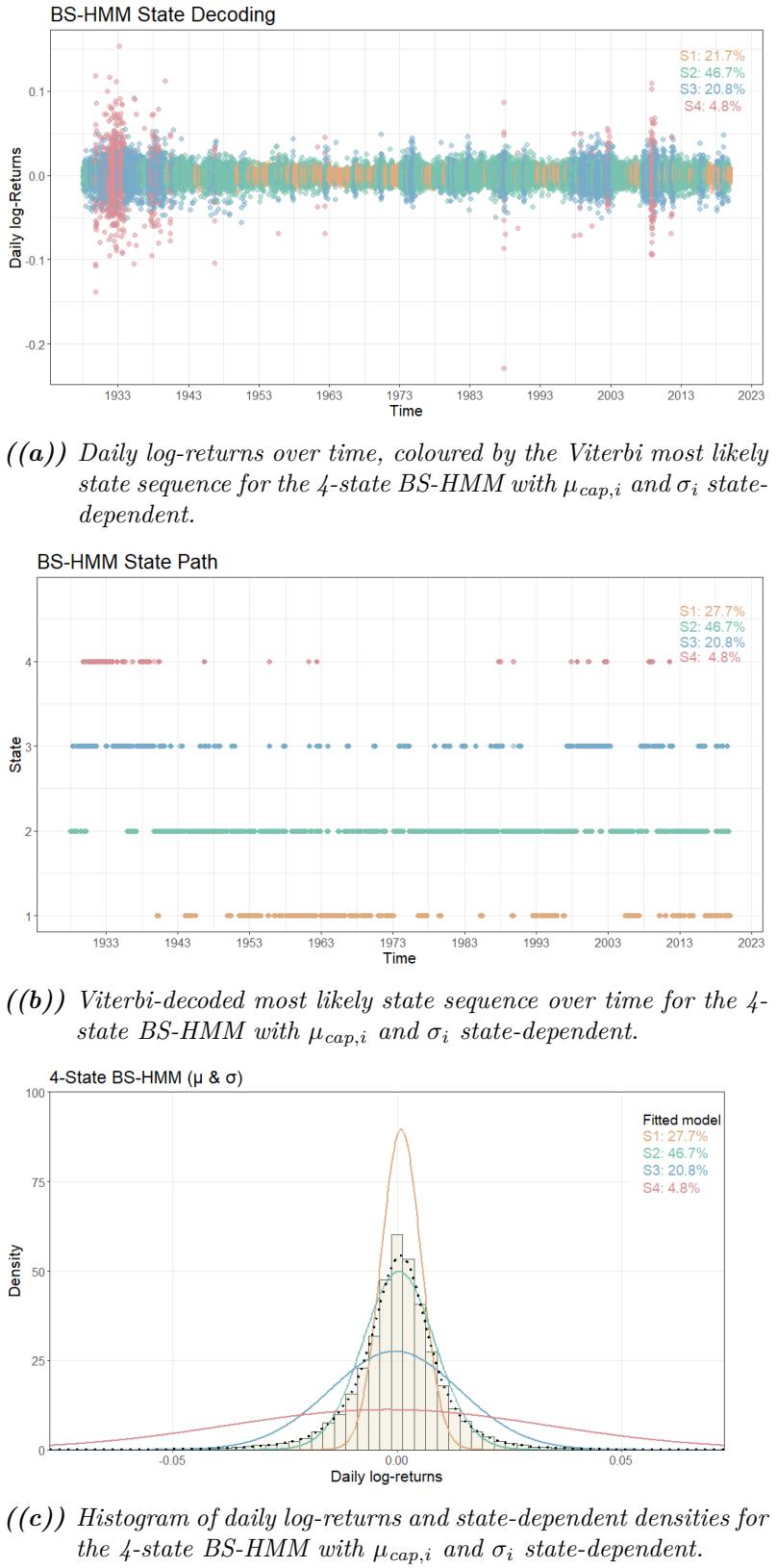
narrow confidence interval [0.0680, 0.0734]. The 95% confidence interval for  $\hat{\mu}_{\text{cap},1}$ , [0.1990, 0.2703], lies well above zero, indicating statistically and economically significant positive excess returns. State 1 can thus be interpreted as a “high-growth, low-volatility” bull state, capturing boom-like periods of sustained, relatively tranquil expansion.

State 2 also has a clearly bullish profile but is more moderate than State 1. The estimated capital-gains drift  $\hat{\mu}_{\text{cap},2} = 0.1062$  corresponds to an annual price appreciation of about 10.6%, with a 95% confidence interval [0.0621, 0.1502] that lies strictly above zero. Combined with the dividend yield  $\hat{q}_2 \approx 0.0336$  (roughly 3.4% per year), this yields a total-return drift of  $\hat{\mu}_{\text{tot},2} = 0.1398$ , i.e. an expected annual total return of about 14%. Volatility is higher than in State 1 but still moderate,  $\hat{\sigma}_2 = 0.1270$  (around 12.7% annually) and below the constant-volatility BSM benchmark  $\hat{\sigma} \approx 0.1883$  (see [Table 10](#)). State 2 therefore reflects a more conventional expansionary or “normal bull” regime with solid positive returns and moderate risk.

State 3 marks a transition to bearish conditions. The capital-gains drift  $\hat{\mu}_{\text{cap},3} = -0.1031$  implies an annual price decline of about 10.3% and even after adding the dividend yield  $\hat{q}_3 \approx 0.0306$  (around 3.1%) the total-return drift remains slightly negative at  $\hat{\mu}_{\text{tot},3} = -0.0725$  (about -7.3% annually). The confidence interval for  $\hat{\mu}_{\text{cap},3}$ , [-0.2193, 0.0132], narrowly straddles zero, indicating a moderate but not overwhelmingly strong downturn. In contrast, volatility increases substantially to  $\hat{\sigma}_3 = 0.2301$  (about 23.0% annually), clearly above the constant-parameter BSM estimate and indicative of elevated market turbulence. State 3 can thus be viewed as a “mild bear” or correction state with near-zero to moderately negative expected returns and high volatility.

State 4 represents the most pessimistic and volatile conditions. The capital-gains drift  $\hat{\mu}_{\text{cap},4} = -0.3818$  implies an annual price decline of about 38.2% and even with the relatively high dividend yield  $\hat{q}_4 \approx 0.0498$  (about 5.0%) the total-return drift remains very negative at  $\hat{\mu}_{\text{tot},4} = -0.3320$  (around -33.2% per year). This corresponds to a deep bear or crisis state with extremely poor expected performance. The confidence interval for  $\hat{\mu}_{\text{cap},4}$ , [-0.9260, 0.1623], is wide, reflecting both the rarity and severity of such episodes. Volatility is extreme,  $\hat{\sigma}_4 = 0.5651$  (around 56.5% annually), with a confidence interval [0.5329, 0.5974] that lies far above the volatility levels in the other states and in the BSM benchmark. State 4 is therefore naturally interpreted as a crisis regime with crash-like dynamics and pronounced market stress.

We observe in [Figure 18\(a\)](#) and [Figure 18\(b\)](#) that the HMM successfully identifies major episodes of market turbulence, including the Wall Street Crash of 1929 and the 1937–1938 recession, the early-1980s recession, Black Monday in 1987, the late-1990s and early-2000s Dot-com Bubble and the 2008 Financial Crisis. In contrast, the constant-parameter BSM produces extremely large or small residuals in these periods. Episodes of the most turbulent bear market conditions (state 4) are relatively rare, whereas the more moderate bear regime (state 3) occurs more frequently. The Markov chain visits the pessimistic and highly volatile state 4 roughly 4.8% of the time and the mild bear state 3 about 20.8%, numbers that are closely aligned with the stationary distribution  $\hat{\delta} = (0.2766, 0.4675, 0.2080, 0.0479)$  implied by the estimated transition matrix. The transition probability matrix  $\hat{\Gamma}$  reinforces the interpretation of an ordered set of regimes. The diagonal elements are all close to one, with  $\hat{\gamma}_{11} = 0.9649$ ,  $\hat{\gamma}_{22} = 0.9683$ ,  $\hat{\gamma}_{33} = 0.9632$  and  $\hat{\gamma}_{44} = 0.9483$ , implying expected regime durations on the order of one to one-and-a-half months in all four states, with crisis episodes somewhat shorter-lived. Off-diagonal entries exhibit a nearly tridiagonal structure, that is, transitions occur predominantly between neighbouring regimes in terms of severity (from 1 to 2, from 2 to 3 and from 3 to 2 or 4), while direct jumps between the most extreme states are essentially absent.



**Figure 18:** BS-HMM presentation.

In particular, conditional on leaving the crisis state, the chain moves to state 3 with probability essentially equal to one, whereas transitions from state 4 directly to the bull states 1 or 2 are ruled out by the estimates. Economically, this means that the economy does not jump from an extremely adverse crisis regime to a normal or boom state, but passes through an intermediate, still bearish but less extreme regime. This gradual adjustment pattern is clearly visible in [Figure 18\(b\)](#), where crises in state 4 are typically preceded and followed by spells in state 3.

The model also appears to correctly identify bull states. The extremely favourable bull market conditions associated with state 1 are prominent during well-known expansionary periods, such as the post-war boom of the 1950s–1960s and much of the 1982–2000 secular bull market, typically linked to disinflation, deregulation and technology-led growth. As expected, the more moderate bull state 2, with moderate volatility and expected returns, is the most frequently occupied regime: the Markov chain spends approximately 46.7% of the time in state 2, fairly evenly spread across the sample, as seen from [Figure 18\(b\)](#).

Taken together, the four states provide a coherent ordering of market conditions. States 1 and 2 represent expansionary bull-market phases with positive expected returns and relatively low volatility, with State 1 capturing particularly strong and stable booms and State 2 corresponding to more typical growth periods. States 3 and 4 capture contractionary bear-market phases: State 3 is a mild bear state with near-zero to moderately negative returns and high volatility, while State 4 corresponds to a deep crisis state with very negative expected returns and extraordinary volatility. Volatility increases monotonically as one moves from the most optimistic to the most pessimistic state, in line with the empirical observation that volatility is higher in downturns. Dividend yields are relatively stable at around 3–3.4% in the first three states, but rise to nearly 5% in the crisis regime, consistent with sharply depressed equity prices in severe downturns.

Relative to the constant-parameter BSM benchmark, which imposes a single drift and volatility ( $\hat{\mu}_{\text{cap}}, \hat{\sigma}$ ) for the entire sample, the BS–HMM offers a much richer and more realistic description of the time-varying risk–return trade-off. Instead of a single “average” regime, the model allows the S&P 500 to switch between distinct bull and bear states with markedly different expected returns and risk levels, including a tranquil high-growth state and a turbulent crisis state. This state dependence is empirically important: periods of strong performance are associated with systematically lower volatility, while severe downturns are accompanied by sharply higher volatility and only partly compensated by higher dividend yields. The BS–HMM therefore provides a natural benchmark for the more elaborate regime-switching specifications developed in the subsequent analysis and highlights the limitations of the constant-parameter BSM in the presence of pronounced time-variation in both expected returns and risk.

**Black-Scholes Continuous State-Space Model** We now consider the BS-SSM $_{\beta}$  from Equation 28.

The parameter estimates are reported in Table 12. The latent factor  $C_t$  is highly persistent, with  $\hat{\rho} = 0.9827$  (95% CI [0.9794, 0.9861]), implying a half-life of shocks on the order of 40 trading days. Combined with an innovation volatility  $\hat{\sigma}_{\varepsilon} = 0.0769$ , this yields a slowly varying but economically meaningful state variable that captures medium-run shifts in market conditions.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_{\varepsilon}$	0.0769 (—, —)	—
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017532
$\hat{\mu}_{\text{tot}}$	0.1238 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_{\mu}$	-0.2459 (—, —)	—
$\hat{\beta}_{\sigma}$	1.2824 (—, —)	—

**Table 12:** BS-SSM $_{\beta}$  using an  $m = 100$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

The Hessian-based covariance matrix for the working-scale parameters is somewhat ill-conditioned, especially in the directions corresponding to  $(\sigma_{\varepsilon}, \beta_{\mu}, \beta_{\sigma})$ . This reflects the fact that, as discussed in the identifiability remark on the BS-SSM $_{\beta}$  specification, the parametrisation is only weakly identified along a scale and sign transformation of the latent factor and its loadings. In our numerical implementation this shows up as slightly negative variance estimates for  $\beta_{\mu}$  and  $\beta_{\sigma}$  on the natural scale, which after truncation at zero lead to vanishing standard errors for these coefficients in Table 12. These entries should therefore not be read as evidence of almost no estimation error, but rather as an indication that the Hessian-based uncertainty quantification for  $(\sigma_{\varepsilon}, \beta_{\mu}, \beta_{\sigma})$  is unreliable because the likelihood is nearly flat along a reparameterisation ridge. By contrast, the baseline level parameters  $(\rho, \mu_{\text{cap}}, q, \sigma)$  are well-behaved and can be interpreted in the usual way and all objects that depend on the implied drift and volatility paths  $(\mu_t, \sigma_t)$ —such as pseudo-residuals, forecast distributions and the out-of-sample error measures reported below—are invariant to this lack of full identifiability.

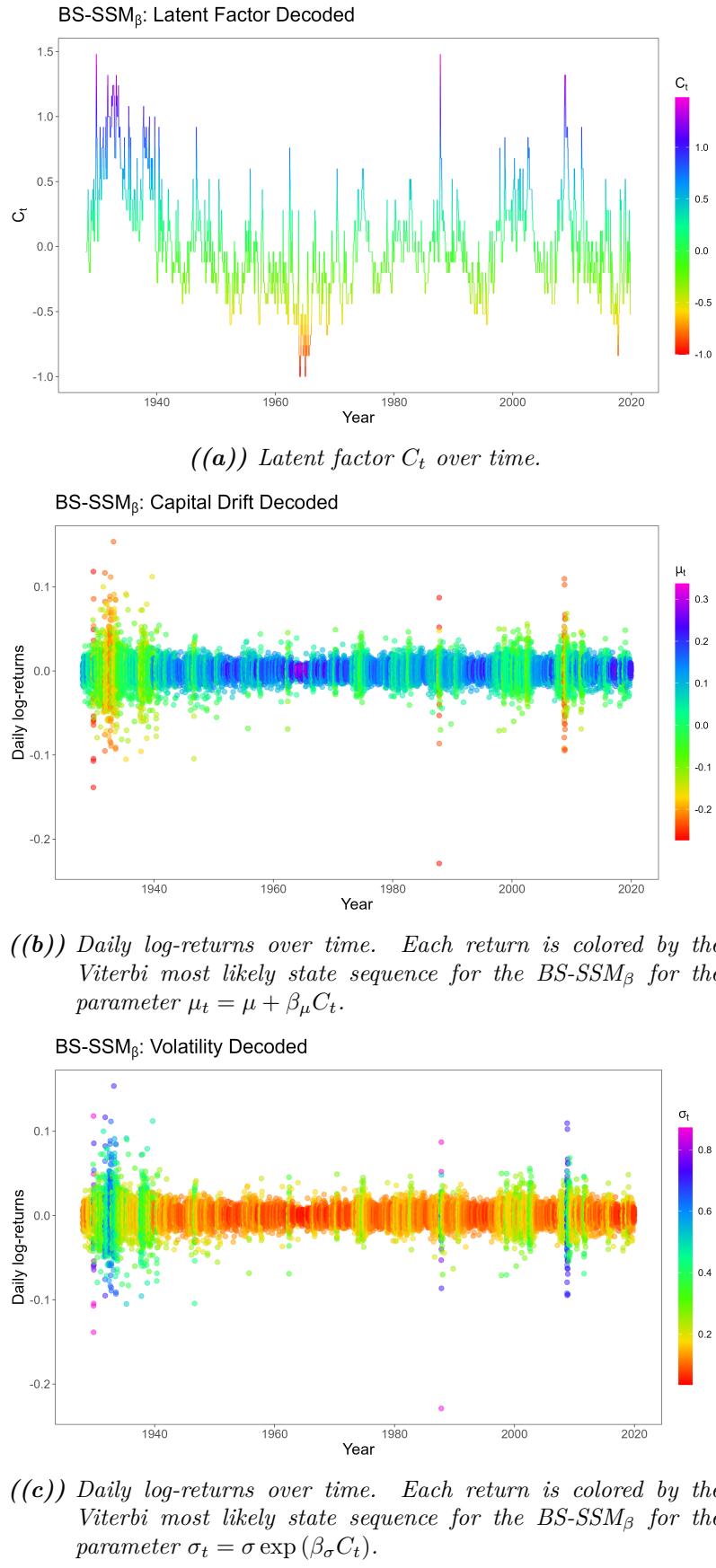
The level parameters have a natural interpretation. The baseline annualised capital-gains drift at  $C_t = 0$  is  $\hat{\mu}_{\text{cap}} = \hat{\mu}_0 = 0.0909$ , corresponding to roughly 9.1% expected annual price appreciation, while the constant dividend yield is estimated as  $\hat{q} \approx 0.0330$  (about 3.3% per year). Together, these imply a baseline total-return drift of  $\hat{\mu}_{\text{tot}} = 0.1238$ , i.e. an expected annual total return of about 12.4% when the latent factor is neutral. The unconditional volatility level is  $\hat{\sigma} = 0.1306$  (around 13.1% annually), which is substantially lower than the constant volatility in the BSM benchmark, reflecting that extreme movements are now captured via the time-varying  $\sigma_t$  rather than by inflating a single constant parameter.

The factor loadings  $\hat{\beta}_{\mu}$  and  $\hat{\beta}_{\sigma}$  govern how the latent state translates into the conditional risk–return trade-off. The estimate  $\hat{\beta}_{\mu} = -0.2459$  implies that, for the chosen orientation of the factor, positive realisations of  $C_t$  reduce the capital-gains drift, while negative values increase it. Con-

versely, the positive loading  $\widehat{\beta}_\sigma = 1.2824$  implies that volatility rises sharply when  $C_t$  is high and falls when  $C_t$  is low. Since  $C_t$  is a latent index measured in arbitrary units, the absolute numerical values of  $(\widehat{\sigma}_\varepsilon, \widehat{\beta}_\mu, \widehat{\beta}_\sigma)$  are not of primary interest; what matters is the relative pattern that periods with elevated  $C_t$  are associated with lower conditional drift and higher conditional volatility. Given the estimated distribution of  $C_t$ , this combination means that periods with elevated volatility tend to coincide with depressed expected returns, whereas tranquil, low-volatility periods are associated with higher expected returns. In other words, the model induces a smooth, continuous analogue of the bull and bear regimes seen in the BS–HMM: instead of discrete jumps between four states, the market moves gradually along a persistent latent factor that jointly modulates drift and volatility.

From Figure 19(a) we see that the latent factor  $C_t$  rises sharply during major crisis episodes and is depressed in tranquil expansionary periods. With  $\beta_\mu < 0$  and  $\beta_\sigma > 0$ , these movements feed directly into the conditional risk–return profile, since  $\mu_t = \mu + \beta_\mu C_t$  decreases when  $C_t$  is high and increases when  $C_t$  is low. As shown in Figure 19(b), the most adverse configuration occurs on 28 October 1929, at the onset of the Wall Street Crash, where the model implies an annualised capital-gains drift of about  $\mu_t \approx -0.27$ , whereas during the calm mid-1960s expansion the drift peaks around  $\mu_t \approx 0.34$  on 6 February 1964.

The volatility channel moves in the opposite direction. By construction  $\sigma_t = \sigma \exp(\beta_\sigma C_t)$  is always positive and increases with  $C_t$  and Figure 19(c) shows that  $\sigma_t$  reaches roughly 0.87 on 28 October 1929 and falls to about 0.04 on 6 February 1964 when conditions are exceptionally tranquil. Taken together, Figure 19(a) and Figure 19(c) demonstrate that the BS-SSM $_\beta$  generates a smooth, continuous analogue of the bull and bear regimes in the BS-HMM, with crisis periods characterized by high volatility and strongly depressed expected returns and calm periods by low volatility and high expected returns. In this way the continuous state-space specification distills the discrete regime structure of the BS-HMM into a single latent factor that tracks the gradual build-up and subsequent unwinding of market stress.



**Figure 19:** BS-SSM $_\beta$  diagnostics.

Relative to the constant-parameter BSM and the discrete-state BS–HMM, the BS-SSM $_{\beta}$  offers a parsimonious yet flexible description of time-variation in both expected returns and risk. A single persistent factor is sufficient to generate episodes that resemble the “high-growth, low-volatility” and “crisis” regimes of the HMM, but the transitions between these configurations are gradual rather than abrupt. This continuous state-space representation fits the data slightly better than the simpler BS-SSM without factor loadings, as reflected in the AIC/BIC comparison, while retaining a clear economic interpretation in terms of a slowly moving latent business-cycle that drives both the level and the variability of equity returns.

### 4.3 Forecast

State-space models often achieve an excellent in-sample fit to historical data yet deliver markedly weaker out-of-sample forecasting performance. This pattern is frequently noted in finance, where such models can retrospectively explain past regime changes or volatility dynamics but often fail to outperform simpler benchmarks on new data [18, 7]. For instance, [18] show that a Markov-switching exchange rate model fits the sample well but does not generate better forecasts than a random walk. More recent evidence likewise suggests that even flexible HMM-based specifications with multiple states, despite high in-sample likelihood, tend to yield only modest gains in out-of-sample accuracy [24]. Similar discrepancies appear outside finance, consistent with the broader tendency of complex models to fit idiosyncratic in-sample variation without improving predictive generalization [33]. Common explanations include overfitting to transient patterns, additional uncertainty from imperfectly inferred latent states and the intrinsic difficulty of anticipating abrupt structural change. Consequently, despite their sophistication and strong in-sample performance, HMMs and related state-space models often provide limited forecasting improvements out-of-sample, motivating ongoing work on more robust forecasting strategies [7].

To assess the point forecasting performance of the three BS-type specifications, we examine the Mean Squared Error (MSE), the Root Mean Square Error (RMSE) and the Mean Absolute Error (MAE) in [Table 13](#).

Model	MSE		RMSE		MAE	
	MSE	Rel. err. [%]	RMSE	Rel. err. [%]	MAE	Rel. err. [%]
BSM	$10^{-4} \times 1.82$	—	0.0135	—	0.00876	0.0251
BS-HMM	$10^{-4} \times 1.82$	0.162	0.0135	0.0809	0.00876	0.0154
BS-SSM $_{\beta}$	$10^{-4} \times 1.82$	0.120	0.0135	0.0598	0.00876	—

**Table 13:** Out-of-sample one-step-ahead prediction errors at three significant digits. Errors are in daily log-return units. Relative error is the percentage increase over the best-performing model for each metric.

The results in [Table 13](#) indicate that, in terms of one-step-ahead point forecasting, the three

Black–Scholes type specifications perform almost indistinguishably on the test sample. The constant-parameter BSM attains the lowest MSE and RMSE, with values of  $10^{-4} \times 1.82$  and 0.0135, respectively, and thus serves as the benchmark for the relative error measures. The BS–HMM and BS- $\text{SSM}_\beta$  are only marginally worse in squared-error terms: their MSEs differ from the BSM by about 0.2% and their RMSEs by less than 0.1%, i.e. at the fourth decimal place in daily log-return units. For MAE, the BS- $\text{SSM}_\beta$  attains the smallest error (0.00876), with the BSM and BS–HMM exhibiting relative increases of approximately 0.03% and 0.02%, respectively. In absolute terms, these differences correspond to changes in average forecast error of order  $10^{-6}$ . Taken together, [Table 13](#) shows that neither the discrete state-space structure of the BS–HMM nor the continuous latent factor of the BS- $\text{SSM}_\beta$  yields a materially better one-step-ahead forecast than the simplest BSM benchmark. This aligns with the empirical literature cited above, which finds that more flexible state-space models rarely deliver large out-of-sample gains relative to parsimonious alternatives when evaluated on short-horizon point forecasts.

[Table 14](#) complements this comparison by summarising the full one-day-ahead forecast distributions at several horizons and forecast origins. For the BSM, the forecast mode, median and mean are identical across all horizons and origins, reflecting the time-homogeneous Gaussian structure of the model: conditional on the estimated parameters, the one-day-ahead return distribution is the same whether the forecast is for one day, one week or several years into the future. The nominal 90% forecast interval is therefore also invariant, at approximately  $(-1.93 \times 10^{-2}, 1.97 \times 10^{-2})$ , and can be interpreted as a horizon-independent benchmark for daily return risk under the constant-parameter specification.

Year   Horizon	Mode	Median	Mean	90% interval
<b>BSM</b>				
2020   1 day	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$(-1.93 \times 10^{-2}, 1.97 \times 10^{-2})$
2020   1 week	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$(-1.93 \times 10^{-2}, 1.97 \times 10^{-2})$
2020   1 month	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$(-1.93 \times 10^{-2}, 1.97 \times 10^{-2})$
2020   3 months	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$(-1.93 \times 10^{-2}, 1.97 \times 10^{-2})$
2020   1 year	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$(-1.93 \times 10^{-2}, 1.97 \times 10^{-2})$
2022   3 years	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$(-1.93 \times 10^{-2}, 1.97 \times 10^{-2})$
2025   full horizon ( $\approx 5.66$ yrs)	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$2.25 \times 10^{-4}$	$(-1.93 \times 10^{-2}, 1.97 \times 10^{-2})$
<b>BS-HMM</b>				
2020   1 day	$8.99 \times 10^{-4}$	$8.54 \times 10^{-4}$	$8.55 \times 10^{-4}$	$(-7.23 \times 10^{-3}, 8.87 \times 10^{-3})$
2020   1 week	$8.97 \times 10^{-4}$	$8.46 \times 10^{-4}$	$7.87 \times 10^{-4}$	$(-8.14 \times 10^{-3}, 9.59 \times 10^{-3})$
2020   1 month	$8.56 \times 10^{-4}$	$7.28 \times 10^{-4}$	$5.92 \times 10^{-4}$	$(-1.11 \times 10^{-2}, 1.18 \times 10^{-2})$
2020   3 months	$8.03 \times 10^{-4}$	$5.87 \times 10^{-4}$	$3.58 \times 10^{-4}$	$(-1.47 \times 10^{-2}, 1.46 \times 10^{-2})$
2020   1 year	$7.80 \times 10^{-4}$	$5.19 \times 10^{-4}$	$2.28 \times 10^{-4}$	$(-1.69 \times 10^{-2}, 1.65 \times 10^{-2})$
2022   3 years	$7.79 \times 10^{-4}$	$5.18 \times 10^{-4}$	$2.27 \times 10^{-4}$	$(-1.70 \times 10^{-2}, 1.65 \times 10^{-2})$
2025   full horizon ( $\approx 5.66$ yrs)	$7.79 \times 10^{-4}$	$5.18 \times 10^{-4}$	$2.27 \times 10^{-4}$	$(-1.70 \times 10^{-2}, 1.65 \times 10^{-2})$
<b>BS-<math>\text{SSM}_\beta</math></b>				
2020   1 day	$9.89 \times 10^{-4}$	$8.91 \times 10^{-4}$	$8.30 \times 10^{-4}$	$(-6.92 \times 10^{-3}, 8.40 \times 10^{-3})$
2020   1 week	$1.02 \times 10^{-3}$	$8.73 \times 10^{-4}$	$7.96 \times 10^{-4}$	$(-7.51 \times 10^{-3}, 8.86 \times 10^{-3})$
2020   1 month	$1.10 \times 10^{-3}$	$8.21 \times 10^{-4}$	$6.80 \times 10^{-4}$	$(-9.59 \times 10^{-3}, 1.05 \times 10^{-2})$
2020   3 months	$1.11 \times 10^{-3}$	$6.98 \times 10^{-4}$	$4.87 \times 10^{-4}$	$(-1.32 \times 10^{-2}, 1.36 \times 10^{-2})$
2020   1 year	$9.89 \times 10^{-4}$	$5.49 \times 10^{-4}$	$3.08 \times 10^{-4}$	$(-1.68 \times 10^{-2}, 1.67 \times 10^{-2})$
2022   3 years	$9.83 \times 10^{-4}$	$5.43 \times 10^{-4}$	$3.01 \times 10^{-4}$	$(-1.69 \times 10^{-2}, 1.69 \times 10^{-2})$
2025   full horizon ( $\approx 5.66$ yrs)	$9.83 \times 10^{-4}$	$5.43 \times 10^{-4}$	$3.01 \times 10^{-4}$	$(-1.69 \times 10^{-2}, 1.69 \times 10^{-2})$

**Table 14:** Multi-horizon forecast summaries. Horizons are measured from the end of the estimation sample (2019-12-31). Nominal 90% forecast intervals are given in parentheses.

By contrast, the BS–HMM and BS–SSM $_{\beta}$  generate horizon-dependent forecast distributions through their latent state dynamics. At very short horizons (1 day and 1 week from 2019-12-31), both models produce substantially higher conditional means and noticeably tighter 90% intervals than the BSM. For example, one day ahead the BSM forecast mean is  $2.25 \times 10^{-4}$ , whereas the BS–HMM and BS–SSM $_{\beta}$  means are  $8.55 \times 10^{-4}$  and  $8.30 \times 10^{-4}$ , roughly four times larger. At the same time, the BS–HMM 90% interval ( $-7.23 \times 10^{-3}$ ,  $8.87 \times 10^{-3}$ ) and the BS–SSM $_{\beta}$  interval ( $-6.92 \times 10^{-3}$ ,  $8.40 \times 10^{-3}$ ) are only about 40% as wide as the BSM band ( $-1.93 \times 10^{-2}$ ,  $1.97 \times 10^{-2}$ ). This reflects the fact that, at the end of the estimation sample, the filtered state distributions in both state-space models put high probability on relatively benign configurations, leading to a more concentrated short-horizon conditional distribution of returns than under the unconditional BSM benchmark.

As the horizon lengthens to one month and three months, the forecast means in the BS–HMM and BS–SSM $_{\beta}$  decline toward smaller positive values, while the 90% intervals widen monotonically. For instance, at a one-year horizon (still forecasting a single daily return but from a start date one year after 2019-12-31), the BS–HMM mean has fallen to  $2.28 \times 10^{-4}$  and the BS–SSM $_{\beta}$  mean to  $3.08 \times 10^{-4}$ , both close to the BSM benchmark of  $2.25 \times 10^{-4}$ . The corresponding 90% intervals, ( $-1.69 \times 10^{-2}$ ,  $1.65 \times 10^{-2}$ ) for the BS–HMM and ( $-1.68 \times 10^{-2}$ ,  $1.67 \times 10^{-2}$ ) for the BS–SSM $_{\beta}$ , are now much closer in width to the BSM band, all three lying in the range of approximately  $\pm 1.6\text{--}2.0\%$  per day.

By the time the horizon reaches three years and the full remaining sample ( $\approx 5.66$  years), the rows for the BS–HMM and BS–SSM $_{\beta}$  have essentially stabilised: the forecast means and 90% intervals are numerically almost identical at three years and at the full horizon. This indicates that the predictive densities have converged to their stationary counterparts, in line with the theoretical result in [Equation 24](#), which states that the  $h$ -step-ahead forecast distribution converges to the marginal stationary density as  $h \rightarrow \infty$ . The speed of this convergence is noteworthy: already at horizons of one to three years, the additional state-dependent information in the BS–HMM and BS–SSM $_{\beta}$  is largely washed out from the perspective of a one-day-ahead forecast.

The chosen set of horizons in [Table 14](#)—one day, one week, one month, three months, one year, three years and the full out-of-sample period—is intended to span the range from very short-term trading and risk management horizons to medium- and long-term investment horizons. The forecast tables therefore provide a compact view of how each model translates the information at the end of 2019 into short-run and long-run predictive distributions. Overall, the evidence suggests that while the BS–HMM and BS–SSM $_{\beta}$  offer richer state-dependent forecast densities and more realistic short-horizon uncertainty quantification, their incremental gains in point forecasting accuracy relative to the simple BSM are modest, consistent with the broader findings in the forecasting literature.

## 5 Discussion

SSM

- Makes sense it is bad as economic regimes and shocks often happen abruptly
- another continuous state process
- another discretization
- another interval other than simple quadrature
- could be great for gradual regime switches that are obvious (heat prices as seasonal effects happen gradually for example)
- $e^{C_t}$  will always be positive and it seems to be  $\sigma$  that dictates the most in the likelihood which leads to  $\mu$  being positive and never negative although we see  $C_t$  can be negative.

## 6 Conclusion

The present research demonstrates that the Black-Scholes Model (BSM), despite its theoretical elegance and computational efficiency, fails to adequately describe the empirical reality of the S&P 500 index over the long run. The fundamental assumption of constant drift and volatility proves to be a distortion rather than a mere simplification, masking the distinct and regime-dependent nature of financial risk that has defined the US equity market over the last century.

The in-sample evidence overwhelmingly indicates that the market does not behave as a geometric Brownian motion with stable variance, but rather as a complex system that cycles through identifiable states of economic sentiment. The estimated 4-state Black-Scholes Hidden Markov Model (BS-HMM) successfully decodes this history and assigns rigorous quantitative parameters to the narratives of market cycles. The model reveals a distinct “Euphoria” state (State 1) characterized by aggressive annualized growth of approximately 23.5% paired with a remarkably low volatility of 7.1%. This stands in contrast to the “Stability” state (State 2), which represents the baseline of the US economy with sustainable 10.6% growth and moderate 12.7% volatility.

Crucially, the analysis confirms that a market crash is not simply a negative draw from a normal distribution but a fundamentally different statistical regime. As observed in the state decoding ([Figure 18\(a\)](#)), the identified “Crisis” state (State 4) successfully isolates major episodes of turbulence—including the 1929 Crash, the 1937–1938 recession, Black Monday (1987), and the 2008 Financial Crisis—into a regime characterized by extreme volatility exceeding 56.5% and deep negative drift of approximately -38.2%. This extreme state is relatively rare, visited roughly 4.8% of the time, whereas the “Correction” state (State 3)—acting as a transitional buffer with -10.3% drift and 23.0% volatility—occurs more frequently, about 20.8% of the time. These frequencies align closely with the stationary distribution  $\hat{\delta} = (0.2766, 0.4675, 0.2080, 0.0479)$ , confirming that while the market spends the majority of its time in growth regimes, specific bear market phases are statistically distinct and persistent phenomena.

The structural dynamics of these regimes are further illuminated by the transition probability matrix  $\hat{\Gamma}$ . The diagonal elements are all close to unity ( $\hat{\gamma}_{11} \approx 0.96, \dots, \hat{\gamma}_{44} \approx 0.95$ ), implying expected regime durations on the order of one to one-and-a-half months. This debunks the notion of fleeting volatility spikes; once the market enters a crisis or correction, it tends to remain there for a meaningful period. Furthermore, the matrix exhibits a tridiagonal structure, indicating that transitions occur predominantly between neighboring regimes (e.g., from Stability to Correction, rather than Stability to Crisis). This discrete segmentation is corroborated by the continuous Black-Scholes State-Space Model ( $BS - SSM_\beta$ ), which identifies a persistent latent stress factor ( $\hat{\rho} \approx 0.98$ ) that acts as a continuous index driving drift negatively and volatility positively.

However, the out-of-sample results serve as a sobering check on the limits of quantitative forecasting. The analysis confirms the existence of a forecasting paradox where a superior description of historical data does not guarantee improved prediction of future one-step-ahead returns. The

complexity inherent in regime-switching models introduces parameter uncertainty and lag, which effectively neutralizes their predictive advantage for short-term point forecasts. In this specific context, the static Black-Scholes Model remains a resilient benchmark for simple directional prediction because its rigidity makes it robust to the noise inherent in daily returns.

Therefore, the primary value of the regime-switching extensions lies not in their ability to act as a crystal ball for the next day's closing price, but in their capacity to function as a sophisticated risk radar. By quantifying the probability of transitioning into a high-volatility state, models like the BS-HMM and the continuous BS-SSM provide risk managers with horizon-dependent density forecasts that the static model cannot offer. This capability allows for the stress-testing of portfolios against specific crisis regimes rather than generic standard deviations, offering a dynamic view of tail risk that adapts to the prevailing economic environment.

In the final analysis, this thesis recommends the adoption of regime-switching models for robust risk management and derivative pricing rather than speculative trading. While predicting the exact moment a storm will break remains impossible, the framework developed here allows market participants to quantify the changing probability of adverse conditions. In financial markets, distinguishing between a low-volatility correction and a high-volatility crisis is the vital difference between solvency and ruin.

# Bibliography

## REFERENCES

- [1] Yacine Aït-Sahalia and Jean Jacod. “Testing for Jumps in a Discretely Observed Process.” In: *Annals of Statistics* 37.1 (2009), pp. 184–222.
- [2] Theodore W. Anderson. *The Statistical Analysis of Time Series*. New York: John Wiley & Sons, 1971.
- [3] Ole E. Barndorff-Nielsen and Neil Shephard. “Power and Bipower Variation with Stochastic Volatility and Jumps.” In: *Journal of Financial Econometrics* 2.1 (2004), pp. 1–37.
- [4] Patrick Billingsley. *Probability and Measure*. 3rd ed. Wiley, 1995.
- [5] Tomas Björk. *Arbitrage theory in continuous time*. 4th ed. Oxford university press, 2020.
- [6] Fischer Black and Myron Scholes. “The pricing of options and corporate liabilities.” In: *Journal of political economy* 81.3 (1973), pp. 637–654.
- [7] Tom Boot and Andreas Pick. “Optimal forecasts from Markov switching models.” In: *Journal of Business & Economic Statistics* 36.4 (2018), pp. 628–642.
- [8] Mark Broadie and Özgür Kaya. “Exact simulation of stochastic volatility and other affine jump diffusion processes.” In: *Operations research* 54.2 (2006), pp. 217–231.
- [9] Kenneth P Burnham and David R Anderson. “Multimodel inference: understanding AIC and BIC in model selection.” In: *Sociological methods & research* 33.2 (2004), pp. 261–304.
- [10] Kenneth P Burnham, David R Anderson, and Kathryn P Huyvaert. “AIC model selection and multimodel inference in behavioral ecology: some background, observations, and comparisons.” In: *Behavioral ecology and sociobiology* 65 (2011), pp. 23–35.
- [11] Olivier Cappé, Eric Moulines, and Tobias Rydén. *Inference in hidden Markov models*. Springer, 2005.
- [12] Kim Christensen, Roel CA Oomen, and Mark Podolskij. “Fact or friction: Jumps at ultra high frequency.” In: *Journal of Financial Economics* 114.3 (2014), pp. 576–599.
- [13] D. R. Cox and E. J. Snell. “A General Definition of Residuals (with discussion).” In: *Journal of the Royal Statistical Society: Series B (Methodological)* 30 (1968), pp. 248–275.
- [14] J. E. Dennis Jr. and R. B. Schnabel. *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Englewood Cliffs, NJ: Prentice-Hall, 1983.
- [15] Peter K. Dunn and Gordon K. Smyth. “Randomized Quantile Residuals.” In: *Journal of Computational and Graphical Statistics* 5 (1996), pp. 236–244.
- [16] Rick Durrett. *Probability: Theory and Examples*. 5th ed. Cambridge University Press, 2019.

- [17] Sean R Eddy. *Biological sequence analysis Probabilistic models of proteins and nucleic acids*. 1998.
- [18] Charles Engel. “Can the Markov Switching Model Forecast Exchange Rates?” In: *Journal of International Economics* 36.1 (1994), pp. 151–165.
- [19] Lorella Fatone et al. “The Use of Statistical Tests to Calibrate the Black-Scholes Asset Dynamics Model Applied to Pricing Options with Uncertain Volatility.” In: *Journal of Probability and Statistics* 2012 (2012), Article ID 931609, 20 pages. DOI: [10.1155/2012/931609](https://doi.org/10.1155/2012/931609).
- [20] Federal Reserve Bank of St. Louis. *Review (Federal Reserve Bank of St. Louis)*. <https://fraser.stlouisfed.org/title/820>. Publication series; issues from the 2000s accessed via FRASER. 1917-2025. (Visited on 12/04/2025).
- [21] William Feller. *An introduction to probability theory and its applications, Volume 2*. Vol. 2. John Wiley & Sons, 1991.
- [22] Jr. Forney G. David. “The Viterbi Algorithm.” In: *Proceedings of the IEEE* 61.3 (1973).
- [23] Sylvia Frühwirth-Schnatter. *Finite Mixture and Markov Switching Models*. Springer Series in Statistics. New York: Springer, 2006.
- [24] Arun Gopalakrishnan and Eric Bradlow. *Hidden Markov Model Backcasting Versus Forecasting Performance*. Tech. rep. Working paper, available at SSRN: <https://ssrn.com/abstract=5623200>. The Wharton School, 2025.
- [25] James D. Hamilton. *Time Series Analysis*. Princeton, NJ: Princeton University Press, 1994.
- [26] Zoé van Havre et al. “Overfitting hidden Markov models with an unknown number of states.” In: *arXiv preprint arXiv:1602.02466* (2016).
- [27] Achim Klenke. *Probability Theory: A Comprehensive Course*. 3rd ed. Springer, 2020.
- [28] Peter E. Kloeden and Eckhard Platen. *Numerical Solution of Stochastic Differential Equations*. Vol. 23. Applications of Mathematics. New York, NY: Springer, 1992.
- [29] Roland Langrock, Iain L MacDonald, and Walter Zucchini. “Some nonstandard stochastic volatility models and their estimation using structured hidden Markov models.” In: *Journal of Empirical Finance* 19.1 (2012), pp. 147–161.
- [30] Brian G Leroux and Martin L Puterman. “Maximum-penalized-likelihood estimation for independent and Markov-dependent mixture models.” In: *Biometrics* (1992), pp. 545–558.
- [31] H Linhart and W Zucchini. “Model Selection, Wiley.” In: *New York* (1986).
- [32] Roger Lord, Remmert Koekkoek, and Dick Van Dijk. “A comparison of biased simulation schemes for stochastic volatility models.” In: *Quantitative Finance* 10.2 (2010), pp. 177–194.

- [33] Spyros Makridakis and Michèle Hibon. “The M3-Competition: Results, Conclusions and Implications.” In: *International Journal of Forecasting* 16.4 (2000), pp. 451–476.
- [34] Brett T. McClintock and Théo Michelot. “momentuHMM: R package for generalized hidden Markov models of animal movement.” In: *Methods in Ecology and Evolution* 9.6 (2018), pp. 1518–1530. DOI: [10.1111/2041-210X.12995](https://doi.org/10.1111/2041-210X.12995). URL: <https://doi.org/10.1111/2041-210X.12995>.
- [35] Geoffrey J McLachlan and David Peel. *Finite mixture models*. John Wiley & Sons, 2000.
- [36] Robert C Merton. “An intertemporal capital asset pricing model.” In: *Econometrica: Journal of the Econometric Society* (1973), pp. 867–887.
- [37] Robert C. Merton. “Option Pricing When Underlying Stock Returns Are Discontinuous.” In: *Journal of Financial Economics* 3.1-2 (1976), pp. 125–144. DOI: [10.1016/0304-405X\(76\)90022-2](https://doi.org/10.1016/0304-405X(76)90022-2).
- [38] Sean P Meyn and Richard L Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012.
- [39] Théo Michelot, Roland Langrock, and Toby Patterson. “moveHMM: An R package for the analysis of animal movement data.” In: *Computer software* (2019).
- [40] John F Monahan. *Numerical methods of statistics*. Cambridge University Press, 2011.
- [41] Manh Cuong Ngô, Mads Peter Heide-Jørgensen, and Susanne Ditlevsen. “Understanding narwhal diving behaviour using Hidden Markov Models with dependent state distributions and long range dependence.” In: *PLoS computational biology* 15.3 (2019), e1006425.
- [42] Jennifer Pohle et al. “Selecting the number of states in hidden Markov models: pragmatic solutions illustrated using animal movement.” In: *Journal of Agricultural, Biological and Environmental Statistics* 22 (2017), pp. 270–293.
- [43] Anders Rahbek and Rasmus Søndergaard Pedersen. *Lecture Notes for NMAK24011U Financial Econometric Time Series Analysis (FinMetrics)*. Course lecture notes, Department of Mathematical Sciences. Aug. 2024.
- [44] Cyrus A. Ramezani and Yong Zeng. “Maximum Likelihood Estimation of the Double Exponential Jump-Diffusion Process.” In: *Annals of Finance* 3.4 (2007), pp. 487–507.
- [45] F. W. Scholz. “Maximum likelihood estimation.” In: *Encyclopedia of Statistical Sciences*. Ed. by Samuel Kotz et al. 2nd. Hoboken, NJ: Wiley, 2006, pp. 4629–4639.
- [46] Robert J. Shiller. *Data Appendix: U.S. Stock Market (notes and documentation)*. Documentation of sources and splicing for price, dividends, and earnings. 2025. URL: <https://www.econ.yale.edu/~shiller/data/chapt26.html> (visited on 11/09/2025).

- [47] Robert J. Shiller. *Online Data: U.S. Stock Market Prices, Dividends, Earnings, and CPI*. Monthly S&P price & dividend series starting in 1871. 2025. URL: <https://www.econ.yale.edu/~shiller/data.htm> (visited on 11/09/2025).
- [48] S&P Dow Jones Indices. *FAQ: S&P 500 Dividend Points Index*. Explains the Dividend Points index and annual reset. 2023. URL: <https://www.spglobal.com/spdji/en/documents/additional-material/faq-sp-500-dividend-points-index.pdf> (visited on 11/09/2025).
- [49] S&P Dow Jones Indices. *Index Mathematics Methodology*. URL: <https://www.spglobal.com/spdji/en/methodology/article/index-mathematics-methodology/> (visited on 11/09/2025).
- [50] S&P Dow Jones Indices. *Index Mathematics Methodology*. Defines price, total return, and net total return; dividends are reinvested on ex-date. 2024. URL: <https://www.spglobal.com/spdji/zh/documents/methodologies/methodology-index-math.pdf> (visited on 11/09/2025).
- [51] S&P Dow Jones Indices. *Index Mathematics Methodology*. 2025. URL: <https://www.spglobal.com/spdji/zh/documents/methodologies/methodology-index-math.pdf> (visited on 11/09/2025).
- [52] S&P Dow Jones Indices. *S&P U.S. Indices Methodology*. 2025. URL: <https://www.spglobal.com/spdji/en/documents/methodologies/methodology-sp-us-indices.pdf> (visited on 11/09/2025).
- [53] S&P Dow Jones Indices. *S&P U.S. Indices Methodology*. 2025. URL: <https://www.spglobal.com/spdji/en/documents/methodologies/methodology-sp-us-indices.pdf> (visited on 11/09/2025).
- [54] Dag Tjøstheim. “Non-linear time series and Markov chains.” In: *Advances in applied probability* 22.3 (1990), pp. 587–611.
- [55] Ingmar Visser and Maarten Speekenbrink. “depmixS4: An R Package for Hidden Markov Models.” In: *Journal of Statistical Software* 36.7 (2010), pp. 1–21. DOI: [10.18637/jss.v036.i07](https://doi.org/10.18637/jss.v036.i07). URL: <https://www.jstatsoft.org/v36/i07/>.
- [56] Andrew J. Viterbi. “Error Bounds for Convolutional Codes and an Asymptotically Optimum Decoding Algorithm.” In: *IEEE Transactions on Information Theory* 13.2 (1967).
- [57] Larry Wasserman. “Bayesian model selection and model averaging.” In: *Journal of mathematical psychology* 44.1 (2000), pp. 92–107.
- [58] Yahoo Finance. *S&P 500 (^GSPC) — Quote and Historical Data*. Data source. 2025. URL: <https://finance.yahoo.com/quote/%5EGSPC/> (visited on 11/09/2025).

- [59] Walter Zucchini and Iain L MacDonald. *Hidden Markov models for time series: an introduction using R*. Chapman and Hall/CRC, 2009.

# Appendix

## A.1 Code

Code used for this paper (and some for the keen reader) is available on this [hyper link to a GitHub repository dedicated to this paper](#).

Short descriptions of each (source) file in the repository is given below with a corresponding mark C or D depending on if the file contains code or data, respectively. Furthermore, we write what type each file is by R, Python or Excel.

- (C, R)

## A.2 Derivations & Proofs

**HMM and SSM** We start by defining two key concepts to assist in most of the proofs: a directed acyclic graph and a parent of a r.v.<sup>9</sup>.

**Definition A.2.1.** Let  $G = (V, E)$  be a directed graph, where  $V$  is a finite set of vertices (or nodes) and  $E \subseteq V \times V$  is a set of directed edges, where an edge  $(u, v) \in E$  indicates a directed link from  $u$  to  $v$ . The graph  $G$  is called a directed acyclic graph (DAG) if and only if it contains no directed cycles; that is, there do not exist distinct vertices  $v_1, v_2, \dots, v_k \in V$  such that  $(v_i, v_{i+1}) \in E$  for all  $i = 1, \dots, k - 1$ , and  $(v_k, v_1) \in E$ . Equivalently,  $G$  is acyclic if there exists a topological ordering of the vertices  $v_1, v_2, \dots, v_n$  such that  $(u, v) \in E \implies$  the index of  $u$  is less than that of  $v$ .

**Definition A.2.2.** Let  $G = (V, E)$  be a directed acyclic graph (DAG), where  $V = \{V_1, \dots, V_n\}$  denotes the set of vertices (r.v.'s) and  $E \subseteq V \times V$  denotes the set of directed edges. For a node  $V_i \in V$ , the parent set of  $V_i$  is defined as

$$\text{pa}(V_i) := \{V_j \in V : (V_j, V_i) \in E\}.$$

That is,  $V_j$  is said to be a parent of  $V_i$  if and only if there exists a directed edge from  $V_j$  into  $V_i$  in the graph.

The driving tool for any of the proofs is the following factorization for the joint distribution of the set of r.v.'s  $V_i$   $i \in \{1, \dots, N\}$  in a directed acyclic graph

$$f_{\mathbf{V}^{(N)}}(\mathbf{v}^{(N)}) = \prod_{i=1}^N f_{V_i|\text{pa}(V_i)}(v_i \mid \text{pa}(v_i)), \quad (39)$$

where  $\text{pa}(V_i)$  denotes all the parents of  $V_i$  in the set  $\{V_1, V_2, \dots, V_N\}$ . For example, consider our usual hidden Markov model setup such as that in Figure 7. The only parent of  $X_k$  is  $C_k$  and for  $k = 2, 3, \dots$  the only parent of  $C_k$  is  $C_{k-1}$  (obviously,  $C_1$  has no parent). As an example, the joint distribution of  $\mathbf{X}^{(t)}$  and  $\mathbf{C}^{(t)}$  is therefore given by

$$f_{\mathbf{X}^{(t)}, C_t}(\mathbf{x}^{(t)}, c) = \mathbb{P}(C_1) \prod_{k=2}^t \mathbb{P}(C_k \mid C_{k-1}) \prod_{k=1}^t f_{X_k, C_k}(x_k, c_k). \quad (40)$$

**Lemma A.2.1.** For  $t \in \mathbb{Z}^+$  and histories  $\mathbf{X}^{(t)}$  and  $\mathbf{C}^{(t)}$  we have that

$$f_{\mathbf{X}^{(t+1)}, C_t, C_{t+1}}(\mathbf{x}^{(t+1)}, c_t, c_{t+1}) = f_{\mathbf{X}^{(t)}, C_t}(\mathbf{x}^{(t+1)}, c_t) \mathbb{P}(C_{t+1} \mid C_t) f_{X_{t+1}, C_{t+1}}(x_{t+1}, c_{t+1})$$

---

<sup>9</sup>Proceeding in the appendix, we use the notation, that a draw of a r.v. (say  $C$ ) is simply denoted by the lower caps of said r.v. (say  $c$ ) and not necessarily the usual state values  $i, j \in \mathcal{C}$ .

*Proof.* By Equation 40 and the analogous expression for the expression  $f_{\mathbf{X}^{(t+1)}, \mathbf{C}^{(t+1)}}(\mathbf{x}^{(t+1)}, \mathbf{c}^{(t+1)})$  imply that

$$f_{\mathbf{X}^{(t+1)}, \mathbf{C}^{(t+1)}}(\mathbf{x}^{(t+1)}, \mathbf{c}^{(t+1)}) = \mathbb{P}(C_{t+1} | C_t) f_{\mathbf{X}^{(t)}, \mathbf{C}^{(t)}}(\mathbf{x}^{(t)}, \mathbf{c}^{(t)}) f_{X_{t+1}|C_{t+1}}(x_{t+1}, c_{t+1})$$

Summing over  $\mathbf{C}^{(t-1)}$  yields the desired result.  $\square$

**Lemma A.2.2.** For  $t = 1, 2, \dots, T-1$ ,

$$f_{\mathbf{X}_{t+1}^T | C_{t+1}}(\mathbf{x}_{t+1}^T | c_{t+1}) = f_{X_{t+1}|C_{t+1}}(x_{t+1} | c_{t+1}) f_{\mathbf{X}_{t+2}}(\mathbf{x}_{t+2})$$

*Proof.* The result follows by

$$\begin{aligned} f_{\mathbf{X}_{t+1}^T, \mathbf{C}_{t+1}^T}(\mathbf{X}_{t+1}^T, \mathbf{C}_{t+1}^T) &= f_{X_{t+1}|C_{t+1}}(x_{t+1} | c_{t+1}) \left( \mathbb{P}(C_{t+1}) \prod_{k=t+2}^T \mathbb{P}(C_k | C_{k-1}) \prod_{k=t+2}^T f_{X_k|C_k}(x_k | c_k) \right) \\ &= f_{X_{t+1}|C_{t+1}}(x_{t+1} | c_{t+1}) f_{\mathbf{X}_{t+2}^T, \mathbf{C}_{t+1}^T}(\mathbf{x}_{t+2}^T, \mathbf{c}_{t+1}^T) \end{aligned}$$

and then summing over  $\mathbf{C}_{t+2}^T$  and dividing by  $\mathbb{P}(C_{t+1})$ .  $\square$

**Lemma A.2.3.** For  $t = 1, 2, \dots, T-1$ ,

$$f_{\mathbf{X}_{t+1}^T | C_{t+1}}(\mathbf{x}_{t+1}^T | c_{t+1}) = f_{\mathbf{X}_{t+1}^T | C_{t+1}, C_t}(\mathbf{x}_{t+1}^T | c_t, c_{t+1}). \quad (\dagger)$$

*Proof.* Simply rewrite the RHS of  $(\dagger)$  to

$$\frac{1}{\mathbb{P}(C_t, C_{t+1})} \sum_{\mathbf{C}_{t+2}^T} f_{\mathbf{X}_{t+1}^T, \mathbf{C}_t^T}(\mathbf{x}_{t+1}^T, \mathbf{c}_t^T),$$

which by Equation 39 reduces to

$$\sum_{\mathbf{C}_{t+2}^T} \prod_{k=t+2}^T \mathbb{P}(C_k | C_{k-1}) \prod_{k=t+1}^T f_{X_k|C_k}(x_k | c_k).$$

The LHS of  $(\dagger)$  is

$$\frac{1}{\mathbb{P}(C_t)} \sum_{\mathbf{C}_{t+2}^T} f_{\mathbf{X}_{t+1}^T, \mathbf{C}_t^T}(\mathbf{x}_{t+1}^T, \mathbf{c}_t^T) = \sum_{\mathbf{C}_{t+2}^T} \prod_{k=t+2}^T \mathbb{P}(C_k | C_{k-1}) \prod_{k=t+1}^T f_{X_k|C_k}(x_k | c_k),$$

which show that both sides reduce to the same expression.  $\square$

**Lemma A.2.4.** For  $r = 1, \dots, T$  and  $i_r \in \{1, \dots, m\}$ , the vectors defined by

$$\alpha_1(i_1) \equiv h(i_1), \quad \alpha_{r+1}(i_{r+1}) \equiv \sum_{i_r=1}^m \alpha_r(i_r) f_{r+1}(i_r, i_{r+1})$$

satisfy

$$\alpha_r(i_r) = \sum_{i_1=1}^m \cdots \sum_{i_{r-1}=1}^m h(i_1) \prod_{t=2}^r f_t(i_{t-1}, i_t). \quad (\dagger)$$

*Proof.* We proceed by induction on  $r$ .

*Base case* ( $r = 1$ ). The product over an empty index set equals 1, so

$$\alpha_1(i_1) = h(i_1) = \sum_{\text{empty}} h(i_1) \cdot 1,$$

which is  $(\dagger)$  for  $r = 1$ .

*Inductive step.* Assume  $(\dagger)$  holds for some  $r \in \{1, \dots, T-1\}$ . Then

$$\begin{aligned} \alpha_{r+1}(i_{r+1}) &= \sum_{i_r=1}^m \alpha_r(i_r) f_{r+1}(i_r, i_{r+1}) \\ &= \sum_{i_r=1}^m \left[ \sum_{i_1=1}^m \cdots \sum_{i_{r-1}=1}^m h(i_1) \prod_{t=2}^r f_t(i_{t-1}, i_t) \right] f_{r+1}(i_r, i_{r+1}) \\ &= \sum_{i_1=1}^m \cdots \sum_{i_r=1}^m h(i_1) \prod_{t=2}^{r+1} f_t(i_{t-1}, i_t), \end{aligned}$$

which is  $(\dagger)$  with  $r$  replaced by  $r+1$ . By induction, the claim holds for all  $r$  and in particular for  $r = T$ :

$$\alpha_T(i_T) = \sum_{i_1=1}^m \cdots \sum_{i_{T-1}=1}^m h(i_1) \prod_{t=2}^T f_t(i_{t-1}, i_t).$$

□

**Lemma A.2.5.** Let  $\mathbf{F}_t$  be the  $m \times m$  matrix with  $(i, j)$  entry  $f_t(i, j)$  and let  $\mathbf{1}_N$  denote the  $m \times 1$  vector of ones. Then

$$\mathbb{S} = \sum_{i_1=1}^m \cdots \sum_{i_T=1}^m h(i_1) \prod_{t=2}^T f_t(i_{t-1}, i_t) = \sum_{i_T=1}^m \alpha_T(i_T) = \boldsymbol{\alpha}_T \mathbf{1}_N = \boldsymbol{\alpha}_1 \mathbf{F}_2 \mathbf{F}_3 \cdots \mathbf{F}_T \mathbf{1}_N.$$

*Proof.* The first equality is the definition of  $\mathbb{S}$ . The second follows from Lemma A.2.4 with  $r = T$ , which shows  $\alpha_T(i_T)$  is the sum over all indices except  $i_T$ . Summing over  $i_T$  gives  $\mathbb{S} = \sum_{i_T} \alpha_T(i_T) = \boldsymbol{\alpha}_T \mathbf{1}_N$ . Finally, by construction  $\boldsymbol{\alpha}_{r+1} = \boldsymbol{\alpha}_r \mathbf{F}_{r+1}$ , hence  $\boldsymbol{\alpha}_T = \boldsymbol{\alpha}_1 \mathbf{F}_2 \cdots \mathbf{F}_T$ , yielding the displayed

formula.  $\square$

**AR(1) process** We prove a Lemma to help us rewrite the AR(1) process in a more convenient form under certain conditions.

**Lemma A.2.6.** *With  $\rho \in \mathbb{R}$  and  $\rho \neq 1$ , then*

$$1 + \rho + \rho^2 + \dots + \rho^n = \sum_{i=0}^n \rho^i = (1 - \rho^{n+1}) / (1 - \rho).$$

If moreover  $|\rho| < 1$ ,  $\rho^n \rightarrow 0$  as  $n \rightarrow \infty$  and

$$\sum_{i=0}^{\infty} \rho^i = 1 / (1 - \rho)$$

*Proof.* Let  $S_n := \sum_{i=0}^n \rho^i$  with  $\rho \in \mathbb{R}$  and  $\rho \neq 1$ . Then

$$(1 - \rho)S_n = \sum_{i=0}^n \rho^i - \sum_{i=0}^n \rho^{i+1} = (1 + \rho + \rho^2 + \dots + \rho^n) - (\rho + \rho^2 + \dots + \rho^{n+1}) = 1 - \rho^{n+1}.$$

Hence

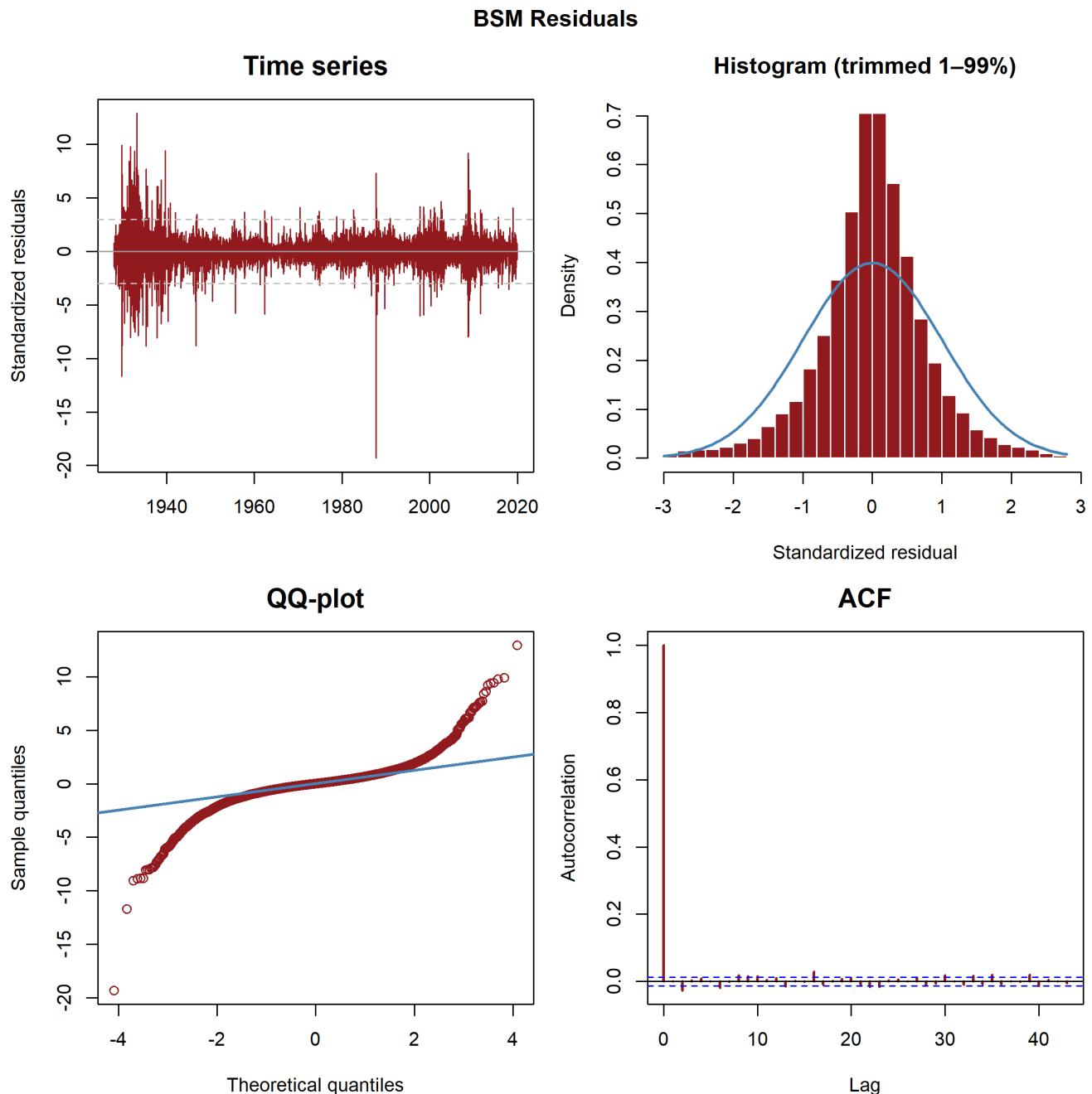
$$S_n = \frac{1 - \rho^{n+1}}{1 - \rho},$$

which proves the finite-sum identity. If moreover  $|\rho| < 1$ , then  $|\rho|^n \rightarrow 0$  as  $n \rightarrow \infty$ . Taking limits in the identity above yields

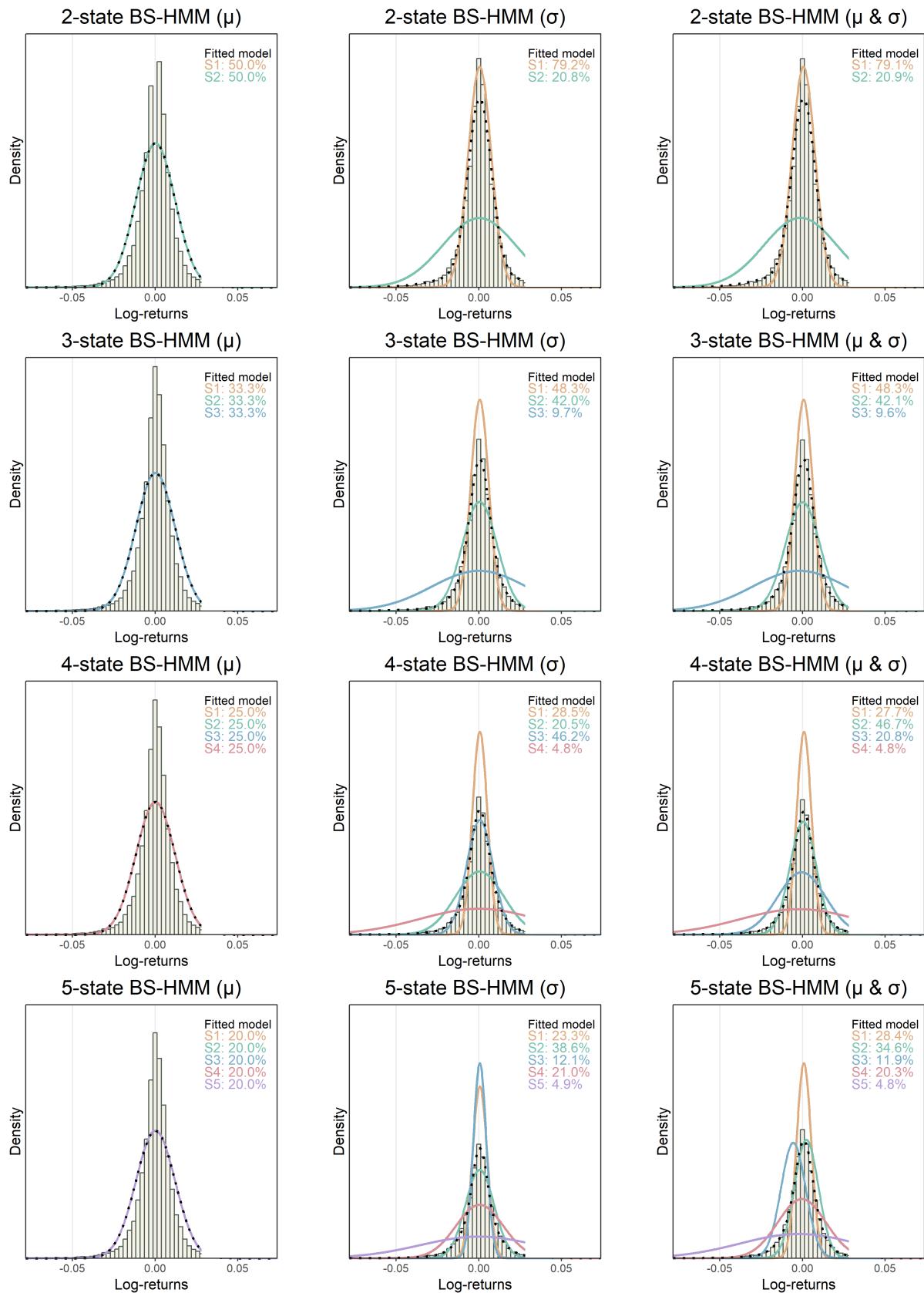
$$\sum_{i=0}^{\infty} \rho^i = \lim_{n \rightarrow \infty} S_n = \lim_{n \rightarrow \infty} \frac{1 - \rho^{n+1}}{1 - \rho} = \frac{1 - 0}{1 - \rho} = \frac{1}{1 - \rho}.$$

$\square$

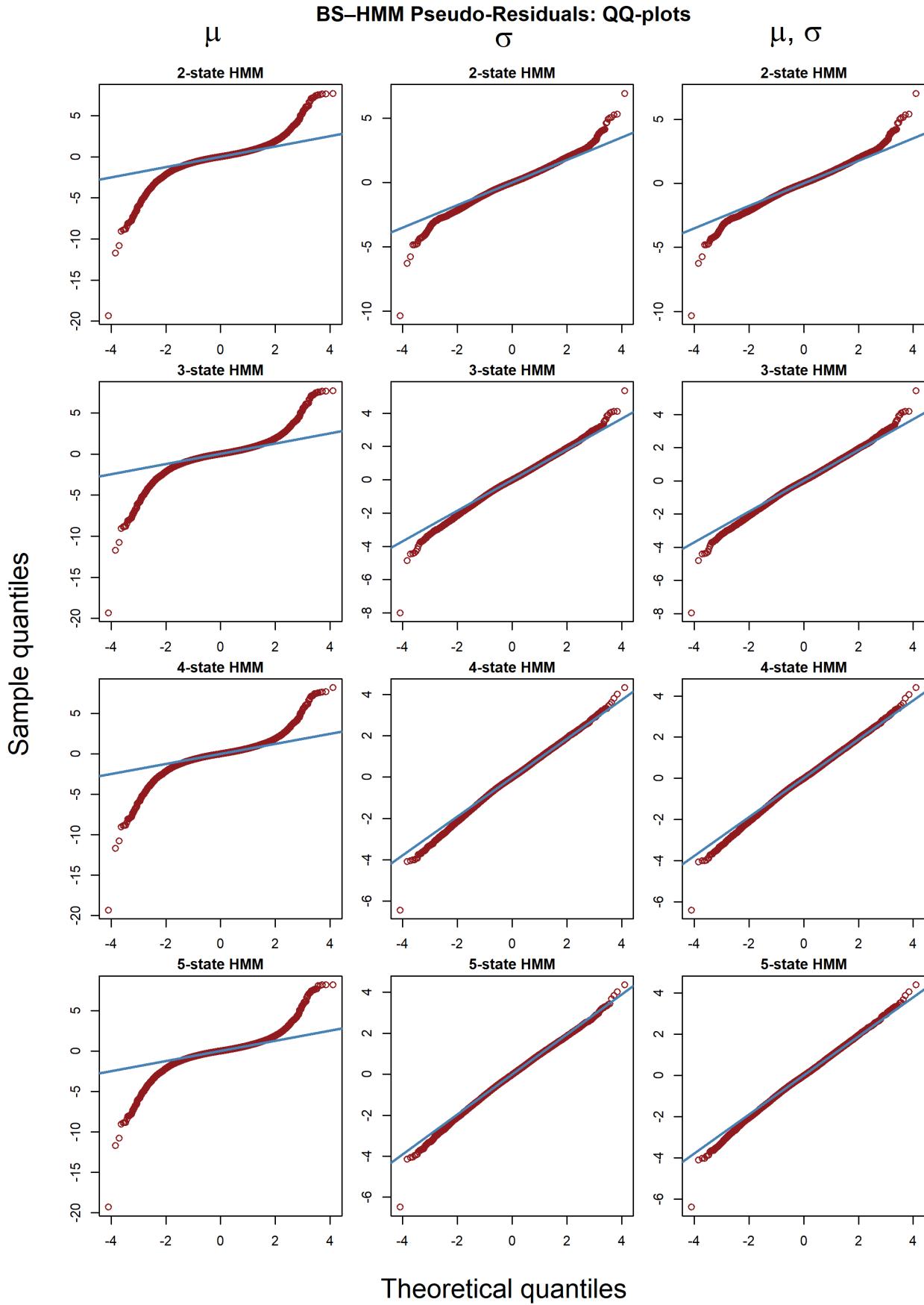
### A.3 Figures



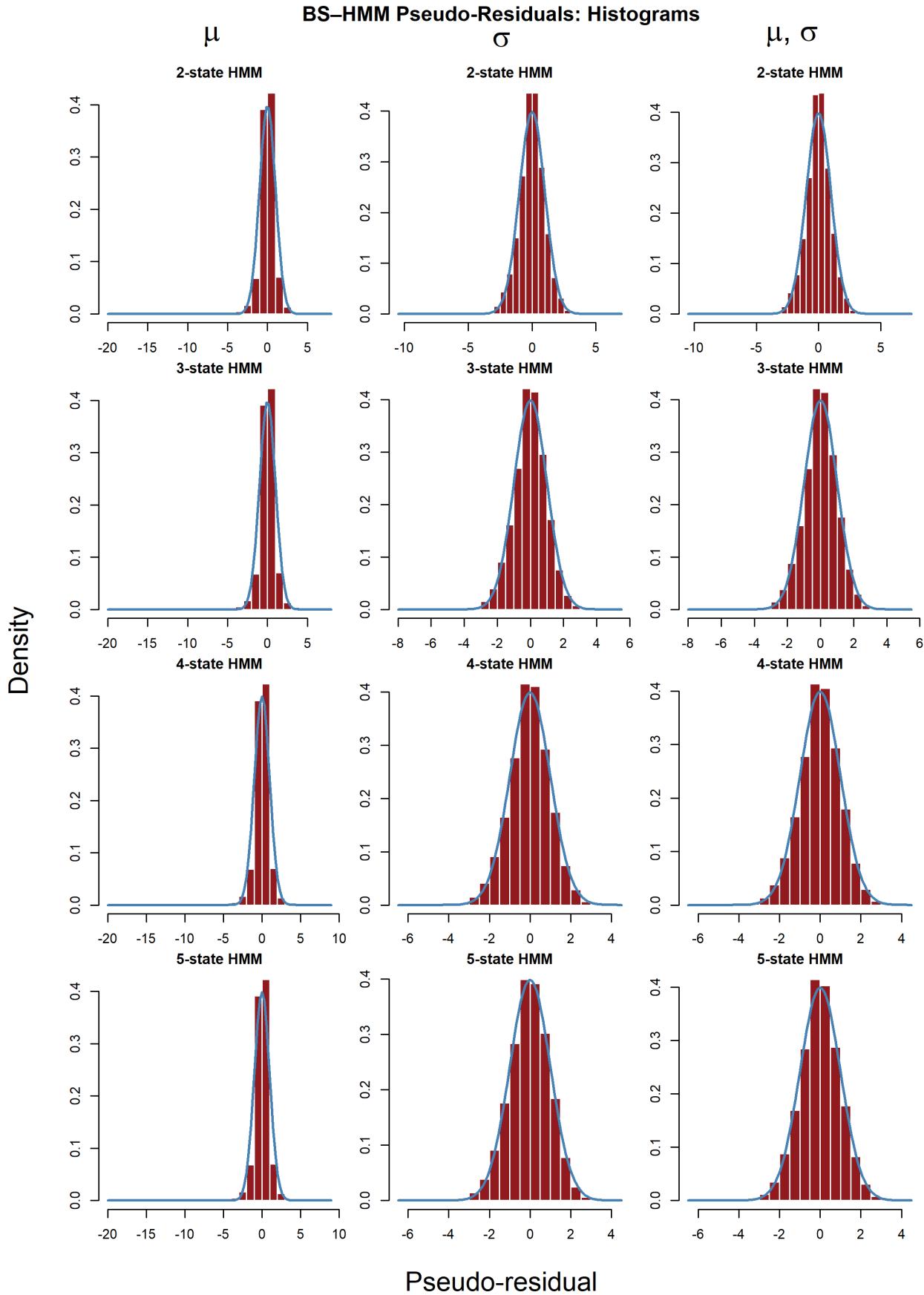
**Figure A.3.1:** Standardized residuals for the BSM. The histogram is trimmed for the purpose of inspection.



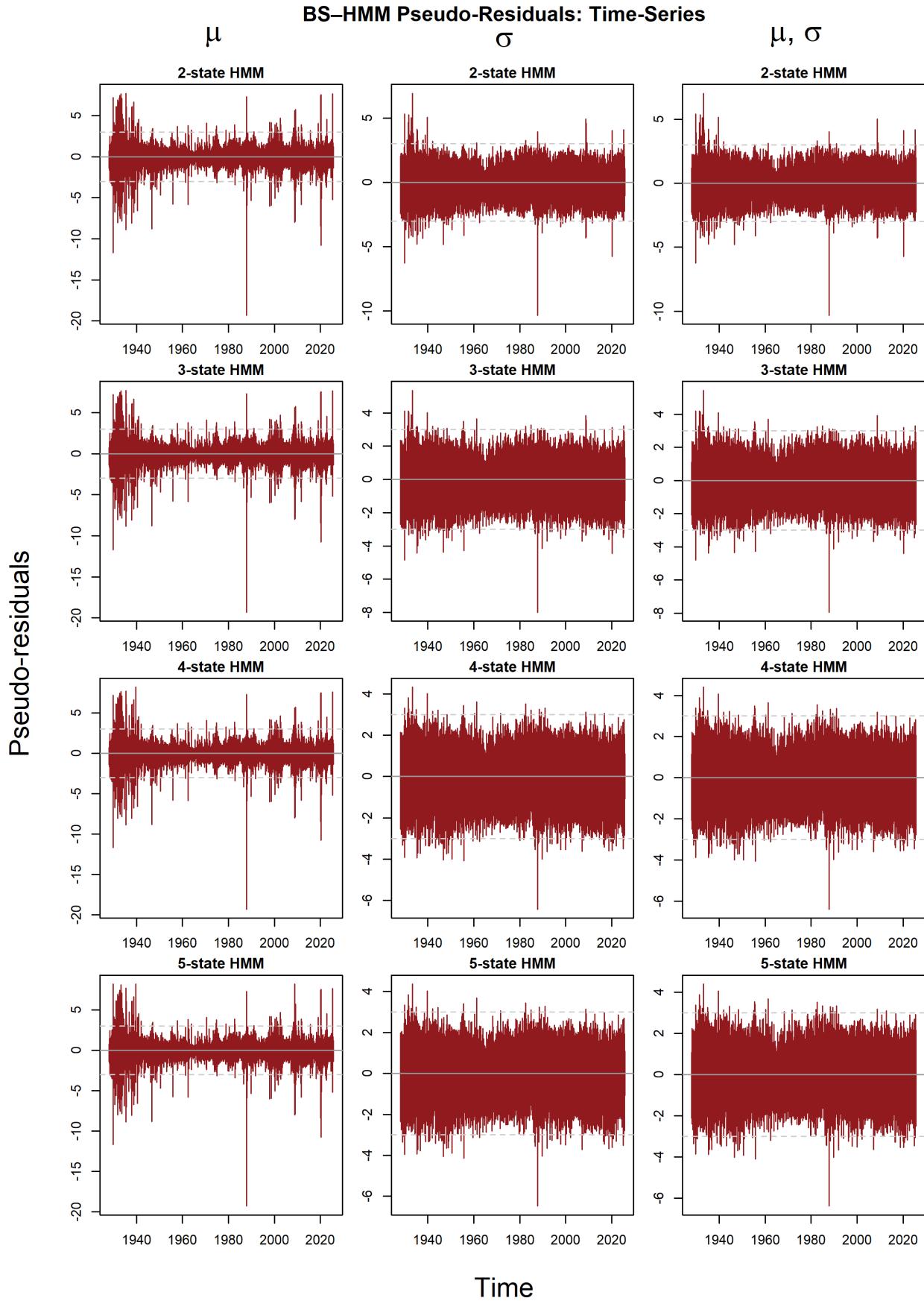
**Figure A.3.2:** State-dependent density plots for the BS-HMMs



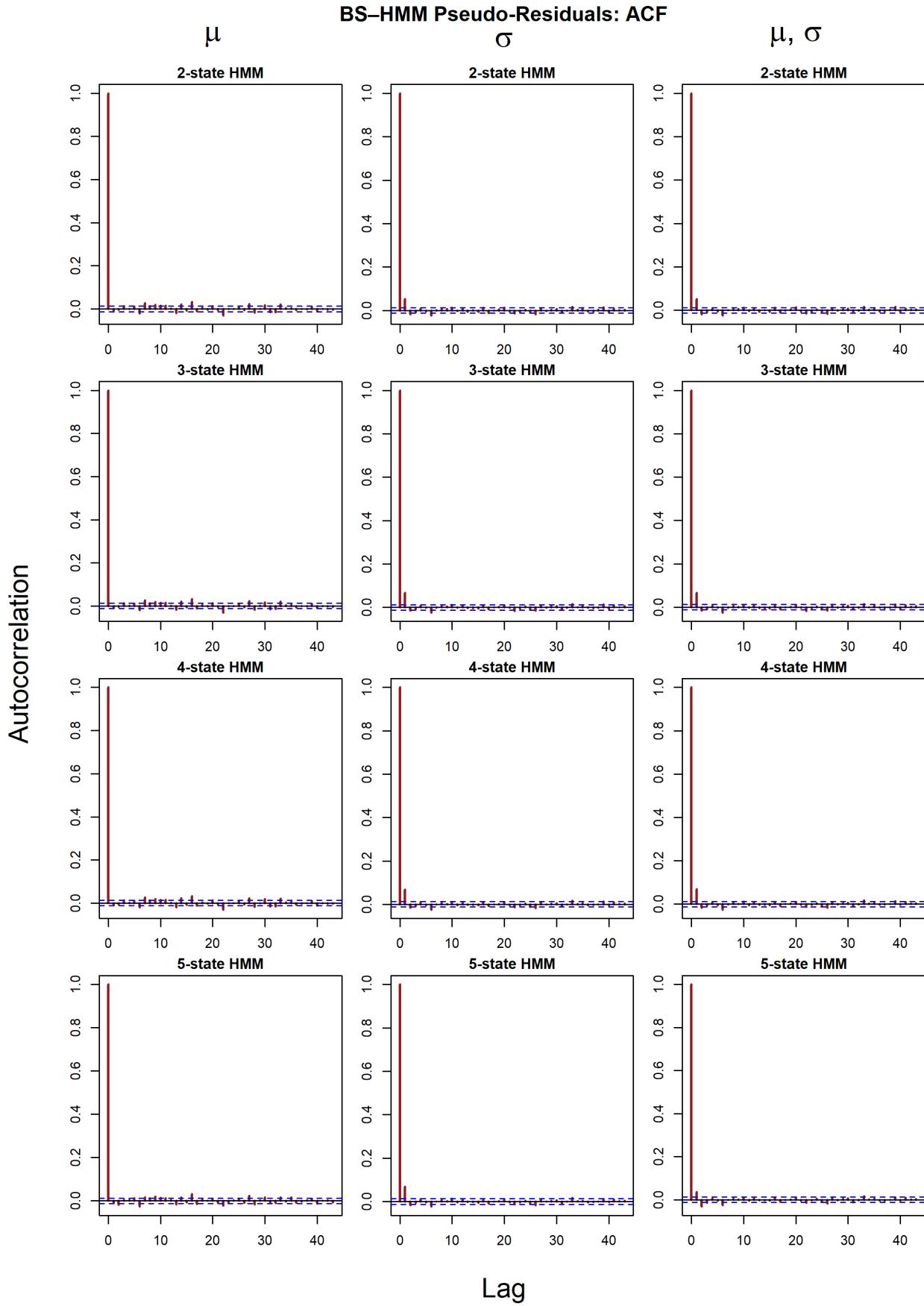
**Figure A.3.3:** QQ-plot for the BS-HMMs.



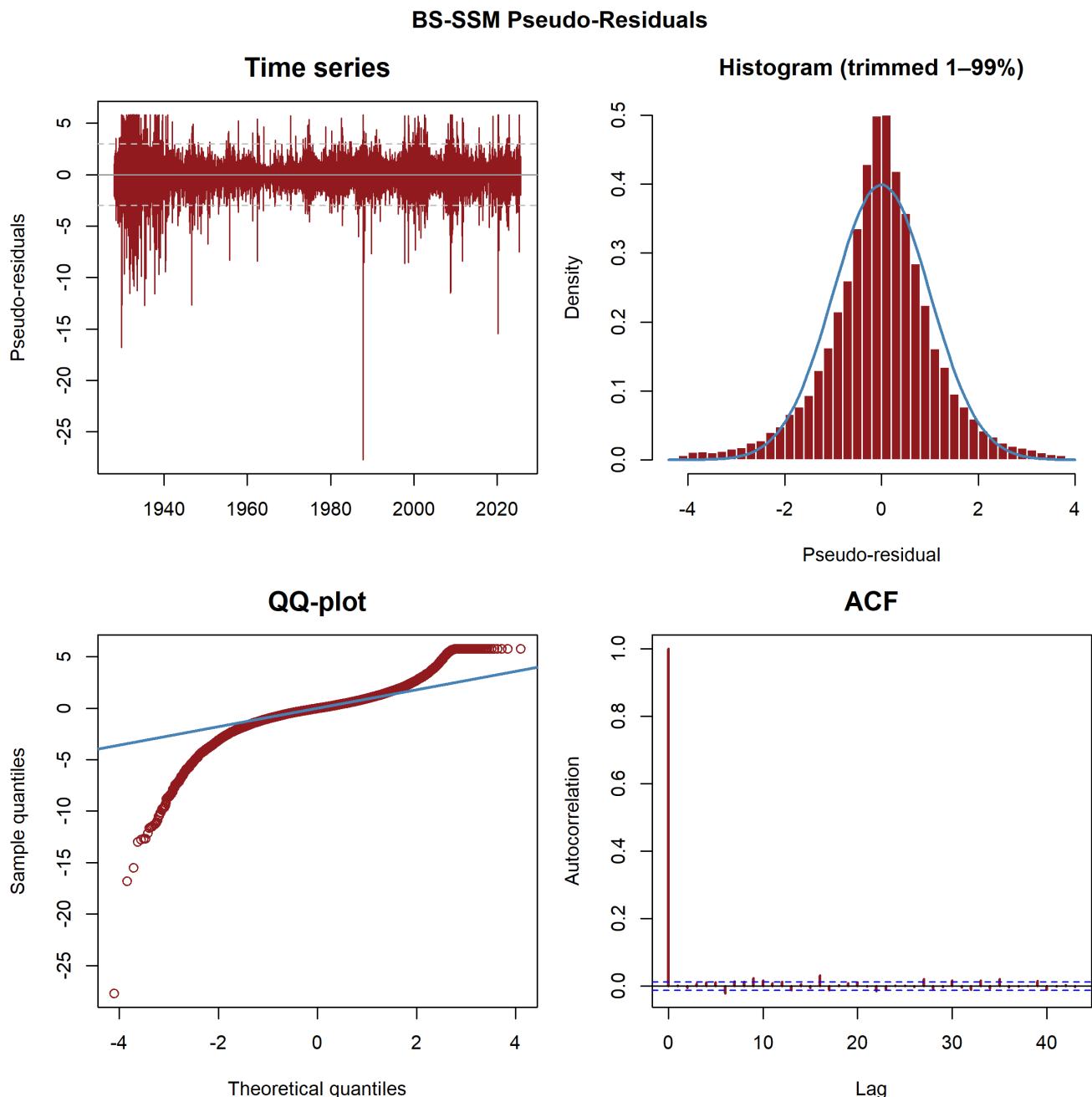
**Figure A.3.4:** Histograms for the BS-HMMs. We trimmed the residuals for inspection



**Figure A.3.5:** Histograms for the BS-HMMs. We trimmed the residuals for inspection

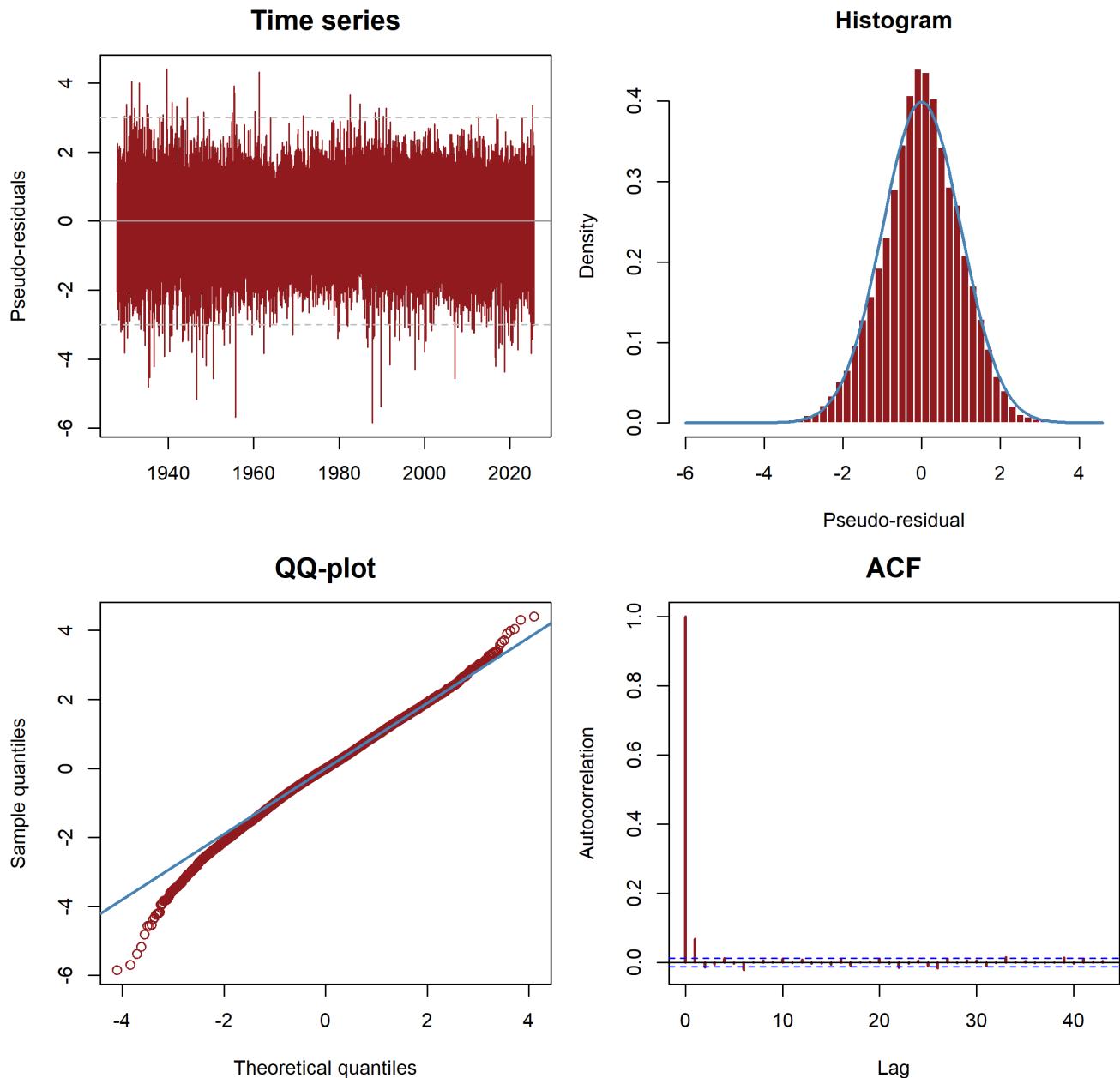


*Figure A.3.6:* ACF plot for the BS-HMMs.



**Figure A.3.7:** Pseudo-residuals for the BS-SSM. The histogram is trimmed for the purpose of inspection.

BS – SSM <sub>$\beta$</sub>  Pseudo-Residuals



**Figure A.3.8:** Pseudo-residuals for the BS-SSM <sub>$\beta$</sub> . The histogram is trimmed for the purpose of inspection.

## A.4 Tables

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.0746 (-0.1947, 0.3439)	0.13741
$\hat{\mu}_{cap,2}$	0.0745 (-0.1943, 0.3433)	0.13715
$\hat{q}_1$	NA (—, —)	—
$\hat{q}_2$	0.0330 (0.0312, 0.0348)	0.00092
$\hat{\mu}_{tot,1}$	NA (—, —)	—
$\hat{\mu}_{tot,2}$	0.1074 (—, —)	—
$\hat{\sigma}$	0.1883 (0.1866, 0.1900)	0.00088
$\widehat{\boldsymbol{\Gamma}} = \begin{pmatrix} 0.8807 (0.6190) & 0.1193 (0.6190) \\ 0.1191 (0.5486) & 0.8809 (0.5486) \end{pmatrix}$		
$\widehat{\boldsymbol{\delta}} = (0.4997 (1.3079), 0.5003 (1.3079))$		

**Table A.4.1:** 2-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  with common volatility. SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.0723 (-0.3200, 0.4646)	0.20014
$\hat{\mu}_{cap,2}$	0.0783 (-0.3000, 0.4566)	0.19300
$\hat{\mu}_{cap,3}$	0.0730 (-0.3157, 0.4617)	0.19830
$\hat{q}_1$	NA (—, —)	—
$\hat{q}_2$	NA (—, —)	—
$\hat{q}_3$	0.0330 (0.0312, 0.0348)	0.00092
$\hat{\mu}_{tot,1}$	NA (—, —)	—
$\hat{\mu}_{tot,2}$	NA (—, —)	—
$\hat{\mu}_{tot,3}$	0.1060 (—, —)	—
$\hat{\sigma}$	0.1883 (0.1866, 0.1900)	0.00088
$\widehat{\boldsymbol{\Gamma}} = \begin{pmatrix} 0.7869 (0.1183) & 0.1065 (0.3642) & — (—) \\ — (—) & — (—) & 0.1065 (0.2975) \\ — (—) & — (—) & 0.7870 (0.1004) \end{pmatrix}$		
$\widehat{\boldsymbol{\delta}} = (— (—), — (—), — (—))$		

**Table A.4.2:** 3-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  with common volatility. SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.0716 (-0.4523, 0.5955)	0.26728
$\hat{\mu}_{cap,2}$	0.0770 (-0.4153, 0.5694)	0.25120
$\hat{\mu}_{cap,3}$	0.0770 (-0.4231, 0.5771)	0.25516
$\hat{\mu}_{cap,4}$	0.0725 (-0.4402, 0.5853)	0.26162
$\hat{q}_1$	NA (—, —)	—
$\hat{q}_2$	NA (—, —)	—
$\hat{q}_3$	NA (—, —)	—
$\hat{q}_4$	0.0330 (0.0312, 0.0348)	0.00092
$\hat{\mu}_{tot,1}$	NA (—, —)	—
$\hat{\mu}_{tot,2}$	NA (—, —)	—
$\hat{\mu}_{tot,3}$	NA (—, —)	—
$\hat{\mu}_{tot,4}$	0.1055 (—, —)	—
$\hat{\sigma}$	0.1883 (0.1866, 0.1900)	0.00088
<hr/>		
$\hat{\Gamma} =$	$\begin{pmatrix} — (—) & 0.0963 (0.1594) & 0.0963 (0.3092) & 0.0963 (0.1551) \\ — (—) & — (—) & — (—) & — (—) \\ 0.0962 (0.3498) & 0.0962 (0.2067) & — (—) & 0.0963 (0.3889) \\ 0.0962 (0.3526) & 0.0963 (0.1787) & 0.0963 (0.1636) & — (—) \end{pmatrix}$	
$\hat{\delta} =$	$(— (—), 0.2500 (0.1264), — (—), 0.2501 (0.1970))$	

**Table A.4.3:** 4-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  with common volatility. SEs are marked in parenthesis for the HMM parameters for readability and space..

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.0714 (-0.5945, 0.7373)	0.33974
$\hat{\mu}_{cap,2}$	0.0758 (-0.5424, 0.6941)	0.31541
$\hat{\mu}_{cap,3}$	0.0773 (-0.5290, 0.6836)	0.30934
$\hat{\mu}_{cap,4}$	0.0760 (-0.5487, 0.7007)	0.31873
$\hat{\mu}_{cap,5}$	0.0722 (-0.5680, 0.7125)	0.32667
$\hat{q}_1$	NA (—, —)	—
$\hat{q}_2$	NA (—, —)	—
$\hat{q}_3$	NA (—, —)	—
$\hat{q}_4$	0.0330 (0.0312, 0.0348)	0.00092
$\hat{q}_5$	NA (—, —)	—
$\hat{\mu}_{tot,1}$	NA (—, —)	—
$\hat{\mu}_{tot,2}$	NA (—, —)	—
$\hat{\mu}_{tot,3}$	NA (—, —)	—
$\hat{\mu}_{tot,4}$	0.1089 (—, —)	—
$\hat{\mu}_{tot,5}$	NA (—, —)	—
$\hat{\sigma}$	0.1883 (0.1866, 0.1900)	0.00088

$$\widehat{\boldsymbol{\Gamma}} = \begin{pmatrix} —(—) & —(—) & —(—) & —(—) & —(—) \\ —(—) & —(—) & —(—) & —(—) & —(—) \\ 0.0878 (0.1984) & —(—) & —(—) & —(—) & —(—) \\ 0.0878 (0.0336) & —(—) & 0.0878 (0.2075) & 0.6488 (0.2634) & 0.0878 (0.1643) \\ 0.0878 (0.2050) & —(—) & —(—) & —(—) & —(—) \end{pmatrix}$$

$$\widehat{\boldsymbol{\delta}} = (0.2000 (0.1697), —(—), 0.2000 (0.2024), —(—), —(—))$$

**Table A.4.4:** 5-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  with common volatility.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap}$	0.1443 (0.1182, 0.1703)	0.01329
$\hat{q}_1$	0.0323 (0.0307, 0.0339)	0.00083
$\hat{q}_2$	0.0356 (0.0318, 0.0394)	0.00193
$\hat{\mu}_{tot,1}$	0.1765 (—, —)	—
$\hat{\mu}_{tot,2}$	0.1798 (—, —)	—
$\hat{\sigma}_1$	0.1105 (0.1089, 0.1122)	0.00083
$\hat{\sigma}_2$	0.3517 (0.3425, 0.3609)	0.00470

$$\widehat{\boldsymbol{\Gamma}} = \begin{pmatrix} 0.9888 (0.0011) & 0.0112 (0.0011) \\ 0.0427 (0.0041) & 0.9573 (0.0041) \end{pmatrix}$$

$$\widehat{\boldsymbol{\delta}} = (0.7915 (0.0167), 0.2085 (0.0167))$$

**Table A.4.5:** 2-state BS-HMM with state-dependent volatility  $\sigma_i$  with common drift. SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{\text{cap}}$	0.1456 (0.1215, 0.1698)	0.01231
$\hat{q}_1$	0.0315 (0.0297, 0.0333)	0.00091
$\hat{q}_2$	0.0326 (0.0305, 0.0348)	0.00109
$\hat{q}_3$	0.0421 (0.0365, 0.0478)	0.00288
$\hat{\mu}_{\text{tot},1}$	0.1771 (—, —)	—
$\hat{\mu}_{\text{tot},2}$	0.1783 (—, —)	—
$\hat{\mu}_{\text{tot},3}$	0.1878 (—, —)	—
$\hat{\sigma}_1$	0.0865 (0.0843, 0.0886)	0.00108
$\hat{\sigma}_2$	0.1675 (0.1629, 0.1721)	0.00237
$\hat{\sigma}_3$	0.4554 (0.4388, 0.4721)	0.00850
$\widehat{\boldsymbol{\Gamma}}$	$\begin{pmatrix} 0.9821 (0.0019) & 0.0174 (0.0020) & 0.0005 (0.0005) \\ 0.0206 (0.0022) & 0.9708 (0.0026) & 0.0086 (0.0013) \\ 0.0000 (0.0000) & 0.0399 (0.0054) & 0.9601 (0.0054) \end{pmatrix}$	
$\widehat{\boldsymbol{\delta}}$	$(0.4833 (0.0281), 0.4201 (0.0236), 0.0966 (0.0140))$	

**Table A.4.6:** 3-state BS-HMM with state-dependent volatility  $\sigma_i$  with common drift. SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{\text{cap}}$	0.1618 (0.1382, 0.1854)	0.01204
$\hat{q}_1$	0.0309 (0.0289, 0.0330)	0.00105
$\hat{q}_2$	0.0308 (0.0280, 0.0336)	0.00141
$\hat{q}_3$	0.0334 (0.0315, 0.0353)	0.00099
$\hat{q}_4$	0.0500 (0.0415, 0.0585)	0.00433
$\hat{\mu}_{\text{tot},1}$	0.1927 (—, —)	—
$\hat{\mu}_{\text{tot},2}$	0.1926 (—, —)	—
$\hat{\mu}_{\text{tot},3}$	0.1952 (—, —)	—
$\hat{\mu}_{\text{tot},4}$	0.2118 (—, —)	—
$\hat{\sigma}_1$	0.0721 (0.0692, 0.0750)	0.00148
$\hat{\sigma}_2$	0.2317 (0.2219, 0.2415)	0.00500
$\hat{\sigma}_3$	0.1278 (0.1235, 0.1322)	0.00221
$\hat{\sigma}_4$	0.5677 (0.5351, 0.6002)	0.01660
$\widehat{\boldsymbol{\Gamma}}$	$\begin{pmatrix} 0.9688 (0.0040) & 0.0000 (0.0000) & 0.0305 (0.0040) & 0.0007 (0.0004) \\ 0.0000 (0.0000) & 0.9630 (0.0039) & 0.0260 (0.0034) & 0.0110 (0.0021) \\ 0.0193 (0.0028) & 0.0110 (0.0015) & 0.9696 (0.0031) & — (—) \\ 0.0000 (0.0000) & 0.0521 (0.0094) & 0.0000 (0.0000) & 0.9479 (0.0094) \end{pmatrix}$	
$\widehat{\boldsymbol{\delta}}$	$(0.2850 (0.0270), 0.2052 (0.0221), 0.4622 (0.0238), 0.0477 (0.0098))$	

**Table A.4.7:** 4-state BS-HMM with state-dependent volatility  $\sigma_i$  with common drift. SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error		
$\hat{\mu}_{cap}$	0.1578 (0.1346, 0.1810)	0.01184		
$\hat{q}_1$	0.0317 (0.0298, 0.0335)	0.00095		
$\hat{q}_2$	0.0305 (0.0271, 0.0339)	0.00175		
$\hat{q}_3$	0.0302 (0.0248, 0.0356)	0.00277		
$\hat{q}_4$	0.0330 (0.0308, 0.0351)	0.00111		
$\hat{q}_5$	0.0493 (0.0410, 0.0576)	0.00425		
$\hat{\mu}_{tot,1}$	0.1895 (—, —)	—		
$\hat{\mu}_{tot,2}$	0.1883 (—, —)	—		
$\hat{\mu}_{tot,3}$	0.1880 (—, —)	—		
$\hat{\mu}_{tot,4}$	0.1908 (—, —)	—		
$\hat{\mu}_{tot,5}$	0.2071 (—, —)	—		
$\hat{\sigma}_1$	0.0707 (0.0678, 0.0736)	0.00147		
$\hat{\sigma}_2$	0.1375 (0.1314, 0.1435)	0.00309		
$\hat{\sigma}_3$	0.0623 (0.0509, 0.0737)	0.00582		
$\hat{\sigma}_4$	0.2274 (0.2175, 0.2374)	0.00506		
$\hat{\sigma}_5$	0.5623 (0.5301, 0.5945)	0.01642		
$\hat{\Gamma}$	$\begin{pmatrix} 0.9757 (0.0035) & 0.0000 (0.0000) & 0.0238 (0.0035) & 0.0000 (0.0000) & 0.0005 (0.0006) \\ 0.0147 (0.0024) & 0.7156 (0.0508) & 0.2580 (0.0500) & 0.0111 (0.0017) & 0.0006 (0.0006) \\ — (—) & 0.9066 (0.0748) & 0.0933 (0.0748) & 0.0000 (0.0000) & — (—) \\ 0.0000 (0.0000) & — (—) & 0.0222 (0.0030) & 0.9672 (0.0036) & 0.0106 (0.0021) \\ 0.0000 (0.0000) & 0.0000 (0.0000) & 0.0000 (0.0000) & 0.0524 (0.0092) & 0.9476 (0.0092) \end{pmatrix}$			
$\hat{\delta}$	(0.2330 (0.0279), 0.3864 (0.0270), 0.1212 (0.0217), 0.2102 (0.0237), 0.0493 (0.0101))			

**Table A.4.8:** 5-state BS-HMM with state-dependent volatility  $\sigma_i$  with common drift. SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.1587 (0.1319, 0.1854)	0.01363
$\hat{\mu}_{cap,2}$	-0.2431 (-0.4030, -0.0832)	0.08158
$\hat{q}_1$	0.0323 (0.0307, 0.0339)	0.00083
$\hat{q}_2$	0.0356 (0.0318, 0.0394)	0.00193
$\hat{\mu}_{tot,1}$	0.1909 (—, —)	—
$\hat{\mu}_{tot,2}$	-0.2075 (—, —)	—
$\hat{\sigma}_1$	0.1104 (0.1088, 0.1120)	0.00083
$\hat{\sigma}_2$	0.3502 (0.3410, 0.3593)	0.00466
$\hat{\Gamma}$	$\begin{pmatrix} 0.9887 (0.0011) & 0.0113 (0.0011) \\ 0.0428 (0.0041) & 0.9572 (0.0041) \end{pmatrix}$	
$\hat{\delta}$	(0.7908 (0.0167), 0.2092 (0.0167))	

**Table A.4.9:** 2-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  and volatility  $\sigma_i$ . SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.1824 (0.1533, 0.2116)	0.01486
$\hat{\mu}_{cap,2}$	0.0425 (-0.0155, 0.1004)	0.02956
$\hat{\mu}_{cap,3}$	-0.3245 (-0.6299, -0.0190)	0.15583
$\hat{q}_1$	0.0315 (0.0298, 0.0333)	0.00091
$\hat{q}_2$	0.0326 (0.0304, 0.0347)	0.00109
$\hat{q}_3$	0.0421 (0.0364, 0.0478)	0.00289
$\hat{\mu}_{tot,1}$	0.2140 (—, —)	—
$\hat{\mu}_{tot,2}$	0.0750 (—, —)	—
$\hat{\mu}_{tot,3}$	-0.2824 (—, —)	—
$\hat{\sigma}_1$	0.0862 (0.0840, 0.0883)	0.00110
$\hat{\sigma}_2$	0.1675 (0.1628, 0.1722)	0.00239
$\hat{\sigma}_3$	0.4543 (0.4376, 0.4710)	0.00850
$\widehat{\Gamma}$	$\begin{pmatrix} 0.9812 (0.0021) & 0.0183 (0.0021) & 0.0005 (0.0004) \\ 0.0216 (0.0024) & 0.9699 (0.0027) & 0.0086 (0.0013) \\ — (—) & 0.0398 (0.0054) & 0.9602 (0.0054) \end{pmatrix}$	
$\widehat{\delta}$	$(0.4827 (0.0278), 0.4209 (0.0233), 0.0964 (0.0140))$	

**Table A.4.10:** 3-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  and volatility  $\sigma_i$ . SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.2346 (0.1990, 0.2703)	0.01820
$\hat{\mu}_{cap,2}$	0.1062 (0.0621, 0.1502)	0.02246
$\hat{\mu}_{cap,3}$	-0.1031 (-0.2193, 0.0132)	0.05930
$\hat{\mu}_{cap,4}$	-0.3818 (-0.9260, 0.1623)	0.27763
$\hat{q}_1$	0.0306 (0.0286, 0.0327)	0.00104
$\hat{q}_2$	0.0336 (0.0316, 0.0355)	0.00099
$\hat{q}_3$	0.0306 (0.0279, 0.0333)	0.00139
$\hat{q}_4$	0.0498 (0.0413, 0.0583)	0.00432
$\hat{\mu}_{tot,1}$	0.2653 (—, —)	—
$\hat{\mu}_{tot,2}$	0.1398 (—, —)	—
$\hat{\mu}_{tot,3}$	-0.0725 (—, —)	—
$\hat{\mu}_{tot,4}$	-0.3320 (—, —)	—
$\hat{\sigma}_1$	0.0707 (0.0680, 0.0734)	0.00139
$\hat{\sigma}_2$	0.1270 (0.1229, 0.1311)	0.00210
$\hat{\sigma}_3$	0.2301 (0.2205, 0.2396)	0.00487
$\hat{\sigma}_4$	0.5651 (0.5329, 0.5974)	0.01646
$\hat{\Gamma}$	$\begin{pmatrix} 0.9649 (0.0044) & 0.0341 (0.0045) & 0.0002 (0.0006) & 0.0007 (0.0005) \\ 0.0207 (0.0030) & 0.9683 (0.0033) & 0.0110 (0.0015) & — (—) \\ 0.0000 (0.0000) & 0.0259 (0.0033) & 0.9632 (0.0039) & 0.0109 (0.0021) \\ 0.0000 (0.0000) & 0.0000 (0.0000) & 0.0517 (0.0094) & 0.9483 (0.0094) \end{pmatrix}$	
$\hat{\delta}$	$(0.2766 (0.0254), 0.4675 (0.0233), 0.2080 (0.0220), 0.0479 (0.0099))$	

**Table A.4.11:** 4-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  and volatility  $\sigma_i$ . SEs are marked in parenthesis for the HMM parameters for readability and space.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\mu}_{cap,1}$	0.2181 (0.1833, 0.2529)	0.01777
$\hat{\mu}_{cap,2}$	0.6049 (0.4268, 0.7831)	0.09089
$\hat{\mu}_{cap,3}$	-1.3869 (-1.7642, -1.0097)	0.19247
$\hat{\mu}_{cap,4}$	-0.0672 (-0.1861, 0.0518)	0.06068
$\hat{\mu}_{cap,5}$	-0.3986 (-0.9461, 0.1488)	0.27930
$\hat{q}_1$	0.0311 (0.0291, 0.0331)	0.00101
$\hat{q}_2$	0.0340 (0.0317, 0.0364)	0.00121
$\hat{q}_3$	0.0350 (0.0311, 0.0388)	0.00196
$\hat{q}_4$	0.0318 (0.0294, 0.0341)	0.00118
$\hat{q}_5$	0.0498 (0.0413, 0.0583)	0.00434
$\hat{\mu}_{tot,1}$	0.2492 (—, —)	—
$\hat{\mu}_{tot,2}$	0.6390 (—, —)	—
$\hat{\mu}_{tot,3}$	-1.3520 (—, —)	—
$\hat{\mu}_{tot,4}$	-0.0354 (—, —)	—
$\hat{\mu}_{tot,5}$	-0.3489 (—, —)	—
$\hat{\sigma}_1$	0.0702 (0.0676, 0.0729)	0.00136
$\hat{\sigma}_2$	0.1156 (0.1111, 0.1202)	0.00233
$\hat{\sigma}_3$	0.1186 (0.1103, 0.1270)	0.00425
$\hat{\sigma}_4$	0.2314 (0.2218, 0.2410)	0.00489
$\hat{\sigma}_5$	0.5665 (0.5340, 0.5990)	0.01659
$\hat{\Gamma}$	$\begin{pmatrix} 0.9626 (0.0047) & 0.0000 (0.0000) & 0.0368 (0.0047) & 0.0000 (0.0000) & 0.0006 (0.0005) \\ 0.0307 (0.0057) & 0.8341 (0.0345) & 0.1349 (0.0310) & — (—) & 0.0004 (0.0006) \\ — (—) & 0.4388 (0.0474) & 0.5190 (0.0526) & 0.0422 (0.0093) & — (—) \\ 0.0000 (0.0000) & 0.0261 (0.0034) & 0.0000 (0.0000) & 0.9631 (0.0040) & 0.0109 (0.0021) \\ 0.0000 (0.0000) & 0.0000 (0.0000) & 0.0000 (0.0000) & 0.0525 (0.0094) & 0.9475 (0.0094) \end{pmatrix}$	
$\hat{\delta}$	(0.2842 (0.0244), 0.3461 (0.0313), 0.1188 (0.0245), 0.2033 (0.0214), 0.0476 (0.0097))	

**Table A.4.12:** 5-state BS-HMM with state-dependent drift  $\mu_{cap,i}$  and volatility  $\sigma_i$ . SEs are marked in parenthesis for the HMM parameters for readability and space.

### Black-Scholes Hidden Markov Model

**Black-Scholes Continuous State Space Model** Note that standard errors of trailing 0's indicate ill Hessian behavior.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (1.0000, 1.0000)	0.000000
$\hat{\sigma}_\varepsilon$	0.0000 (0.0000, 0.0000)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	-6.0032 (-6.0032, -6.0032)	0.000000
$\hat{\mu}_{\text{tot}}$	-5.9703 (—, —)	—
$\hat{\sigma}$	0.2591 (0.2591, 0.2591)	0.000000

**Table A.4.13:** BS-SSM using an  $m = 20$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9990 (0.9984, 0.9996)	0.000304
$\hat{\sigma}_\varepsilon$	0.0144 (0.0131, 0.0156)	0.000652
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1112 (0.0756, 0.1468)	0.018159
$\hat{\mu}_{\text{tot}}$	0.1441 (—, —)	—
$\hat{\sigma}$	0.1735 (0.1711, 0.1759)	0.001205

**Table A.4.14:** BS-SSM using an  $m = 70$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9990 (0.9984, 0.9996)	0.000307
$\hat{\sigma}_\varepsilon$	0.0145 (0.0132, 0.0157)	0.000652
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1111 (0.0756, 0.1467)	0.018154
$\hat{\mu}_{\text{tot}}$	0.1441 (—, —)	—
$\hat{\sigma}$	0.1735 (0.1711, 0.1758)	0.001203

**Table A.4.15:** BS-SSM using an  $m = 100$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9990 (0.9984, 0.9996)	0.000309
$\hat{\sigma}_\varepsilon$	0.0145 (0.0132, 0.0158)	0.000653
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1111 (0.0755, 0.1467)	0.018150
$\hat{\mu}_{\text{tot}}$	0.1441 (—, —)	—
$\hat{\sigma}$	0.1734 (0.1711, 0.1758)	0.001202

**Table A.4.16:** BS-SSM using an  $m = 200$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9943 (0.9927, 0.9959)	0.000820
$\hat{\sigma}_\varepsilon$	0.0522 (0.0487, 0.0558)	0.001818
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1511 (0.1161, 0.1861)	0.017838
$\hat{\mu}_{\text{tot}}$	0.1841 (—, —)	—
$\hat{\sigma}$	0.1687 (0.1646, 0.1727)	0.002090

**Table A.4.17:** BS-SSM using an  $m = 40$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9943 (0.9927, 0.9959)	0.000823
$\hat{\sigma}_\varepsilon$	0.0524 (0.0488, 0.0559)	0.001820
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1509 (0.1160, 0.1858)	0.017822
$\hat{\mu}_{\text{tot}}$	0.1839 (—, —)	—
$\hat{\sigma}$	0.1685 (0.1644, 0.1726)	0.002079

**Table A.4.18:** BS-SSM using an  $m = 70$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9943 (0.9927, 0.9959)	0.000824
$\hat{\sigma}_\varepsilon$	0.0524 (0.0488, 0.0559)	0.001821
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1509 (0.1159, 0.1858)	0.017818
$\hat{\mu}_{\text{tot}}$	0.1838 (—, —)	—
$\hat{\sigma}$	0.1685 (0.1644, 0.1725)	0.002076

**Table A.4.19:** BS-SSM using an  $m = 100$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9943 (0.9927, 0.9959)	0.000825
$\hat{\sigma}_\varepsilon$	0.0524 (0.0488, 0.0560)	0.001821
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1508 (0.1159, 0.1858)	0.017816
$\hat{\mu}_{\text{tot}}$	0.1838 (—, —)	—
$\hat{\sigma}$	0.1685 (0.1644, 0.1725)	0.002074

**Table A.4.20:** BS-SSM using an  $m = 200$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (—, —)	—
$\hat{\sigma}_\varepsilon$	0.0000 (—, —)	—
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0409 (—, —)	—
$\hat{\mu}_{\text{tot}}$	0.0738 (—, —)	—
$\hat{\sigma}$	0.1418 (—, —)	—

**Table A.4.21:** BS-SSM using an  $m = 20$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9845 (0.9814, 0.9876)	0.001586
$\hat{\sigma}_\varepsilon$	0.0921 (0.0852, 0.0990)	0.003499
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1425 (0.1126, 0.1724)	0.015252
$\hat{\mu}_{\text{tot}}$	0.1755 (—, —)	—
$\hat{\sigma}$	0.1351 (0.1257, 0.1445)	0.004802

**Table A.4.22:** BS-SSM using an  $m = 40$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9845 (0.9814, 0.9876)	0.001585
$\hat{\sigma}_\varepsilon$	0.0921 (0.0852, 0.0989)	0.003493
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1427 (0.1128, 0.1726)	0.015266
$\hat{\mu}_{\text{tot}}$	0.1756 (—, —)	—
$\hat{\sigma}$	0.1353 (0.1259, 0.1447)	0.004797

**Table A.4.23:** BS-SSM using an  $m = 70$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9845 (0.9814, 0.9876)	0.001585
$\hat{\sigma}_\varepsilon$	0.0921 (0.0852, 0.0989)	0.003492
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1427 (0.1128, 0.1726)	0.015269
$\hat{\mu}_{\text{tot}}$	0.1757 (—, —)	—
$\hat{\sigma}$	0.1353 (0.1259, 0.1447)	0.004796

**Table A.4.24:** BS-SSM using an  $m = 100$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9845 (0.9814, 0.9876)	0.001585
$\hat{\sigma}_\varepsilon$	0.0921 (0.0852, 0.0989)	0.003493
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1427 (0.1128, 0.1727)	0.015271
$\hat{\mu}_{\text{tot}}$	0.1757 (—, —)	—
$\hat{\sigma}$	0.1353 (0.1259, 0.1447)	0.004796

**Table A.4.25:** BS-SSM using an  $m = 200$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (1.0000, 1.0000)	0.000000
$\hat{\sigma}_\varepsilon$	0.0001 (0.0001, 0.0001)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0800 (0.0435, 0.1165)	0.018602
$\hat{\mu}_{\text{tot}}$	0.1130 (—, —)	—
$\hat{\sigma}$	0.1637 (0.1637, 0.1637)	0.000000

**Table A.4.26:** BS-SSM using an  $m = 20$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (1.0000, 1.0000)	0.000000
$\hat{\sigma}_\varepsilon$	0.0000 (0.0000, 0.0000)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	7.1856 (7.1856, 7.1856)	0.000000
$\hat{\mu}_{\text{tot}}$	7.2186 (—, —)	—
$\hat{\sigma}$	0.4511 (0.4511, 0.4511)	0.000000

**Table A.4.27:** BS-SSM using an  $m = 40$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9842 (0.9810, 0.9873)	0.001614
$\hat{\sigma}_\varepsilon$	0.0933 (0.0862, 0.1003)	0.003590
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1390 (0.1097, 0.1684)	0.014977
$\hat{\mu}_{\text{tot}}$	0.1720 (—, —)	—
$\hat{\sigma}$	0.1314 (0.1214, 0.1414)	0.005124

**Table A.4.28:** BS-SSM using an  $m = 70$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9842 (0.9810, 0.9873)	0.001614
$\hat{\sigma}_\varepsilon$	0.0933 (0.0862, 0.1003)	0.003590
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1390 (0.1097, 0.1684)	0.014976
$\hat{\mu}_{\text{tot}}$	0.1720 (—, —)	—
$\hat{\sigma}$	0.1314 (0.1214, 0.1414)	0.005123

**Table A.4.29:** BS-SSM using an  $m = 100$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9842 (0.9810, 0.9873)	0.001614
$\hat{\sigma}_\varepsilon$	0.0933 (0.0862, 0.1003)	0.003590
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1390 (0.1097, 0.1684)	0.014976
$\hat{\mu}_{\text{tot}}$	0.1720 (—, —)	—
$\hat{\sigma}$	0.1314 (0.1214, 0.1414)	0.005123

**Table A.4.30:** BS-SSM using an  $m = 200$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9983 (0.9983, 0.9983)	0.000003
$\hat{\sigma}_\varepsilon$	0.0002 (0.0002, 0.0002)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	-2.7355 (-3.5777, -1.8933)	0.429717
$\hat{\mu}_{\text{tot}}$	-2.7025 (—, —)	—
$\hat{\sigma}$	3.0458 (2.7486, 3.3430)	0.151629

**Table A.4.31:** BS-SSM using an  $m = 40$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9994 (0.9994, 0.9994)	0.000000
$\hat{\sigma}_\varepsilon$	0.0023 (0.0023, 0.0024)	0.000018
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0550 (0.0492, 0.0609)	0.002972
$\hat{\mu}_{\text{tot}}$	0.0880 (—, —)	—
$\hat{\sigma}$	0.0276 (0.0276, 0.0276)	0.000000

**Table A.4.32:** BS-SSM using an  $m = 70$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9842 (0.9810, 0.9873)	0.001614
$\hat{\sigma}_\varepsilon$	0.0933 (0.0862, 0.1003)	0.003590
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1391 (0.1097, 0.1684)	0.014977
$\hat{\mu}_{\text{tot}}$	0.1720 (—, —)	—
$\hat{\sigma}$	0.1314 (0.1214, 0.1414)	0.005123

**Table A.4.33:** BS-SSM using an  $m = 100$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9842 (0.9810, 0.9873)	0.001614
$\hat{\sigma}_\varepsilon$	0.0933 (0.0862, 0.1003)	0.003590
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1391 (0.1097, 0.1684)	0.014976
$\hat{\mu}_{\text{tot}}$	0.1720 (—, —)	—
$\hat{\sigma}$	0.1314 (0.1214, 0.1414)	0.005123

**Table A.4.34:** BS-SSM using an  $m = 200$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

**Black-Scholes Continuous State Space Beta Model** Note that standard errors of trailing 0's indicate ill Hessian behavior.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9999 (0.9999, 0.9999)	0.000003
$\hat{\sigma}_\varepsilon$	0.0001 (0.0001, 0.0001)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0815 (0.0815, 0.0815)	0.000001
$\hat{\mu}_{\text{tot}}$	0.1145 (—, —)	—
$\hat{\sigma}$	0.0204 (0.0204, 0.0204)	0.000000
$\hat{\beta}_\mu$	-1.5964 (-1.6129, -1.5799)	0.008424
$\hat{\beta}_\sigma$	29.7184 (29.7184, 29.7184)	0.000000

**Table A.4.35:** BS-SSM $_\beta$  using an  $m = 20$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9999 (0.9999, 0.9999)	0.000000
$\hat{\sigma}_\varepsilon$	0.0000 (0.0000, 0.0000)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	-0.2213 (-0.3395, -0.1032)	0.060291
$\hat{\mu}_{\text{tot}}$	-0.1884 (—, —)	—
$\hat{\sigma}$	0.0494 (0.0494, 0.0494)	0.000000
$\hat{\beta}_\mu$	-2.2842 (-5.2731, 0.7047)	1.524932
$\hat{\beta}_\sigma$	34.0251 (34.0251, 34.0251)	0.000000

**Table A.4.36:** BS-SSM $_\beta$  using an  $m = 40$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (1.0000, 1.0000)	0.000000
$\hat{\sigma}_\varepsilon$	0.0000 (0.0000, 0.0000)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0244 (0.0244, 0.0244)	0.000000
$\hat{\mu}_{\text{tot}}$	0.0574 (—, —)	—
$\hat{\sigma}$	0.2005 (0.2005, 0.2005)	0.000000
$\hat{\beta}_\mu$	6.2974 (6.2974, 6.2974)	0.000000
$\hat{\beta}_\sigma$	21.3210 (21.3210, 21.3210)	0.000000

**Table A.4.37:** BS-SSM $_\beta$  using an  $m = 70$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0130 (0.0130, 0.0130)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017520
$\hat{\mu}_{\text{tot}}$	0.1239 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004942
$\hat{\beta}_\mu$	-1.4539 (-1.4539, -1.4539)	0.000000
$\hat{\beta}_\sigma$	7.5824 (7.5824, 7.5824)	0.000000

**Table A.4.38:** BS-SSM $_\beta$  using an  $m = 100$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0130 (-0.3410, 0.3670)	0.180599
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0546, 0.1272)	0.018539
$\hat{\mu}_{\text{tot}}$	0.1239 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004949
$\hat{\beta}_\mu$	-1.4547 (-40.8457, 37.9362)	20.097417
$\hat{\beta}_\sigma$	7.5868 (-199.0721, 214.2457)	105.438209

**Table A.4.39:** BS-SSM $_\beta$  using an  $m = 200$  point grid and truncation  $b_{\max} = 0.5$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9990 (0.9990, 0.9990)	0.000000
$\hat{\sigma}_\varepsilon$	0.0001 (0.0001, 0.0001)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	-0.5994 (-0.8299, -0.3690)	0.117563
$\hat{\mu}_{\text{tot}}$	-0.5665 (—, —)	—
$\hat{\sigma}$	0.1062 (0.0249, 0.1875)	0.041498
$\hat{\beta}_\mu$	0.0166 (-4.5620, 4.5952)	2.336028
$\hat{\beta}_\sigma$	3.3083 (-12.0061, 18.6227)	7.813446

**Table A.4.40:** BS-SSM $_\beta$  using an  $m = 20$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (1.0000, 1.0000)	0.000000
$\hat{\sigma}_\varepsilon$	0.0000 (0.0000, 0.0000)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	-5.2982 (-5.9325, -4.6640)	0.323609
$\hat{\mu}_{\text{tot}}$	-5.2653 (—, —)	—
$\hat{\sigma}$	28.8391 (28.0949, 29.5833)	0.379718
$\hat{\beta}_\mu$	-37.7873 (-39.7743, -35.8002)	1.013803
$\hat{\beta}_\sigma$	12.5466 (12.4915, 12.6017)	0.028107

**Table A.4.41:** BS-SSM $_\beta$  using an  $m = 40$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0365 (0.0365, 0.0365)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017531
$\hat{\mu}_{\text{tot}}$	0.1239 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_\mu$	-0.5183 (-0.5183, -0.5183)	0.000000
$\hat{\beta}_\sigma$	2.7028 (2.7028, 2.7028)	0.000000

**Table A.4.42:** BS-SSM $_\beta$  using an  $m = 70$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0365 (0.0365, 0.0365)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017528
$\hat{\mu}_{\text{tot}}$	0.1239 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_\mu$	-0.5186 (-0.5186, -0.5186)	0.000000
$\hat{\beta}_\sigma$	2.7043 (2.7043, 2.7043)	0.000000

**Table A.4.43:** BS-SSM $_\beta$  using an  $m = 100$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0364 (0.0364, 0.0364)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1253)	0.017535
$\hat{\mu}_{\text{tot}}$	0.1239 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_\mu$	-0.5190 (-0.5190, -0.5190)	0.000000
$\hat{\beta}_\sigma$	2.7063 (2.7063, 2.7063)	0.000000

**Table A.4.44:** BS-SSM $_\beta$  using an  $m = 200$  point grid and truncation  $b_{\max} = 1.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9993 (0.9993, 0.9993)	0.000000
$\hat{\sigma}_\varepsilon$	0.0001 (0.0001, 0.0001)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.4009 (0.4009, 0.4009)	0.000000
$\hat{\mu}_{\text{tot}}$	0.4338 (—, —)	—
$\hat{\sigma}$	0.1353 (0.0179, 0.2527)	0.059892
$\hat{\beta}_\mu$	-0.0659 (-0.0659, -0.0659)	0.000000
$\hat{\beta}_\sigma$	3.6550 (-5.0218, 12.3317)	4.426908

**Table A.4.45:** BS-SSM $_\beta$  using an  $m = 20$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0705 (0.0705, 0.0705)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017529
$\hat{\mu}_{\text{tot}}$	0.1238 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_\mu$	-0.2683 (-0.2683, -0.2683)	0.000000
$\hat{\beta}_\sigma$	1.3990 (1.3990, 1.3990)	0.000000

**Table A.4.46:** BS-SSM $_\beta$  using an  $m = 70$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0705 (0.0705, 0.0705)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017531
$\hat{\mu}_{\text{tot}}$	0.1238 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004942
$\hat{\beta}_\mu$	-0.2683 (-0.2683, -0.2683)	0.000000
$\hat{\beta}_\sigma$	1.3993 (1.3993, 1.3993)	0.000000

**Table A.4.47:** BS-SSM $_\beta$  using an  $m = 100$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0704 (0.0704, 0.0704)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017530
$\hat{\mu}_{\text{tot}}$	0.1238 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004942
$\hat{\beta}_\mu$	-0.2684 (-0.2684, -0.2684)	0.000000
$\hat{\beta}_\sigma$	1.3999 (1.3999, 1.3999)	0.000000

**Table A.4.48:** BS-SSM $_\beta$  using an  $m = 200$  point grid and truncation  $b_{\max} = 2.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9981 (0.9981, 0.9981)	0.000000
$\hat{\sigma}_\varepsilon$	0.0002 (0.0002, 0.0002)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.1012 (0.1012, 0.1012)	0.000000
$\hat{\mu}_{\text{tot}}$	0.1342 (—, —)	—
$\hat{\sigma}$	0.1140 (-0.0001, 0.2282)	0.058230
$\hat{\beta}_\mu$	-0.0297 (-0.0297, -0.0297)	0.000000
$\hat{\beta}_\sigma$	2.0544 (-4.6174, 8.7262)	3.403985

**Table A.4.49:** BS-SSM $_\beta$  using an  $m = 20$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9995 (0.9995, 0.9995)	0.000001
$\hat{\sigma}_\varepsilon$	0.0001 (0.0001, 0.0001)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	1.1720 (0.8863, 1.4577)	0.145788
$\hat{\mu}_{\text{tot}}$	1.2050 (—, —)	—
$\hat{\sigma}$	1.0067 (1.0067, 1.0067)	0.000000
$\hat{\beta}_\mu$	-0.8772 (-3.2273, 1.4729)	1.199027
$\hat{\beta}_\sigma$	-0.4718 (-0.4718, -0.4718)	0.000000

**Table A.4.50:** BS-SSM $_\beta$  using an  $m = 40$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9986 (0.9986, 0.9986)	0.000003
$\hat{\sigma}_\varepsilon$	0.0001 (0.0001, 0.0001)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.4133 (0.4133, 0.4133)	0.000000
$\hat{\mu}_{\text{tot}}$	0.4463 (—, —)	—
$\hat{\sigma}$	0.1283 (-0.1555, 0.4122)	0.144821
$\hat{\beta}_\mu$	-1.7006 (-1.7006, -1.7006)	0.000000
$\hat{\beta}_\sigma$	10.5171 (-41.0901, 62.1244)	26.330239

**Table A.4.51:** BS-SSM $_\beta$  using an  $m = 70$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0769 (0.0769, 0.0769)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017532
$\hat{\mu}_{\text{tot}}$	0.1238 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_\mu$	-0.2459 (-0.2459, -0.2459)	0.000000
$\hat{\beta}_\sigma$	1.2824 (1.2824, 1.2824)	0.000000

**Table A.4.52:** BS-SSM $_\beta$  using an  $m = 100$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0769 (0.0769, 0.0769)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017531
$\hat{\mu}_{\text{tot}}$	0.1238 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_\mu$	-0.2459 (-0.2459, -0.2459)	0.000000
$\hat{\beta}_\sigma$	1.2825 (1.2825, 1.2825)	0.000000

**Table A.4.53:** BS-SSM $_\beta$  using an  $m = 200$  point grid and truncation  $b_{\max} = 3.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (1.0000, 1.0000)	0.000001
$\hat{\sigma}_\varepsilon$	0.0000 (0.0000, 0.0000)	0.000002
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	-700.6359 (-702.2372, -699.0346)	0.816979
$\hat{\mu}_{\text{tot}}$	-700.6029 (—, —)	—
$\hat{\sigma}$	0.0000 (0.0000, 0.0000)	0.000000
$\hat{\beta}_\mu$	-3508.6879 (-3516.6923, -3500.6835)	4.083880
$\hat{\beta}_\sigma$	-3527.8868 (-3527.8868, -3527.8868)	0.000000

**Table A.4.54:** BS-SSM $_\beta$  using an  $m = 20$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (1.0000, 1.0000)	0.000000
$\hat{\sigma}_\varepsilon$	0.0000 (0.0000, 0.0000)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.2385 (-0.9011, 1.3781)	0.581422
$\hat{\mu}_{\text{tot}}$	0.2715 (—, —)	—
$\hat{\sigma}$	1.9752 (1.9752, 1.9752)	0.000000
$\hat{\beta}_\mu$	1.6992 (-9.3980, 12.7963)	5.661823
$\hat{\beta}_\sigma$	7.5654 (7.5654, 7.5654)	0.000000

**Table A.4.55:** BS-SSM $_\beta$  using an  $m = 40$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	1.0000 (1.0000, 1.0000)	0.000000
$\hat{\sigma}_\varepsilon$	0.0000 (0.0000, 0.0000)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0064 (-7.1635, 7.1764)	3.658142
$\hat{\mu}_{\text{tot}}$	0.0394 (—, —)	—
$\hat{\sigma}$	0.0035 (0.0035, 0.0035)	0.000000
$\hat{\beta}_\mu$	-0.8171 (-2.9652, 1.3311)	1.095999
$\hat{\beta}_\sigma$	1.0604 (1.0604, 1.0604)	0.000000

**Table A.4.56:** BS-SSM $_\beta$  using an  $m = 70$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0769 (0.0769, 0.0769)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017532
$\hat{\mu}_{\text{tot}}$	0.1238 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_\mu$	-0.2459 (-0.2459, -0.2459)	0.000000
$\hat{\beta}_\sigma$	1.2824 (1.2824, 1.2824)	0.000000

**Table A.4.57:** BS-SSM $_\beta$  using an  $m = 100$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.

Parameter	Estimate (95% CI)	Std. Error
$\hat{\rho}$	0.9827 (0.9794, 0.9861)	0.001701
$\hat{\sigma}_\varepsilon$	0.0769 (0.0769, 0.0769)	0.000000
$\hat{q}$	0.0330 (0.0321, 0.0339)	0.000917
$\hat{\mu}_{\text{cap}}$	0.0909 (0.0565, 0.1252)	0.017531
$\hat{\mu}_{\text{tot}}$	0.1238 (—, —)	—
$\hat{\sigma}$	0.1306 (0.1209, 0.1403)	0.004941
$\hat{\beta}_\mu$	-0.2459 (-0.2459, -0.2459)	0.000000
$\hat{\beta}_\sigma$	1.2824 (1.2824, 1.2824)	0.000000

**Table A.4.58:** BS-SSM $_\beta$  using an  $m = 200$  point grid and truncation  $b_{\max} = 4.0$ . Parameter estimates with 95% confidence intervals in parentheses.