

---

---

---

## *Half Title*

---

Introduction to Bayesian Econometrics:  
A GUIded toolkit using R



---

---

---

*Title Page*

---

Introduction to Bayesian Econometrics:  
A GUIded toolkit using R

by Andrés Ramírez-Hassan, PhD. Statistical Science.

---

LOC Page

*To my wife, Estephania, and to my parents, Nancy  
and Orlando.*



---

---

## *Contents*

---

<b>Foreword</b>	<b>xi</b>
<b>Preface</b>	<b>xiii</b>
<b>Introduction</b>	<b>xvii</b>
<b>Symbols</b>	<b>xxi</b>
<b>I Foundations: Theory, simulation methods and programming</b>	<b>1</b>
<b>1 Basic formal concepts</b>	<b>3</b>
1.1 The Bayes' rule . . . . .	3
1.2 Bayesian framework: A brief summary of theory . . . . .	9
1.2.1 Example: Health insurance . . . . .	15
1.3 Bayesian reports: Decision theory under uncertainty . . . . .	26
1.3.1 Example: Health insurance continues . . . . .	30
1.4 Summary . . . . .	31
1.5 Exercises . . . . .	32
<b>2 Conceptual differences between the Bayesian and Frequentist approaches</b>	<b>33</b>
2.1 The concept of probability . . . . .	33
2.2 Subjectivity is not the key . . . . .	34
2.3 Estimation, hypothesis testing and prediction . . . . .	35
2.4 The likelihood principle . . . . .	39
2.5 Why is not the Bayesian approach that popular? . . . . .	41
2.6 A simple working example . . . . .	43
2.6.1 Example: Math test . . . . .	44
2.7 Summary . . . . .	45
2.8 Exercises . . . . .	46
<b>3 Cornerstone models: Conjugate families</b>	<b>49</b>
3.1 Motivation of conjugate families . . . . .	49
3.1.1 Examples of exponential family distributions . . . . .	51
3.2 Conjugate prior to exponential family . . . . .	54
3.2.1 Examples: Theorem 4.2.1 . . . . .	55

3.3	Linear regression: The conjugate normal-normal/inverse gamma model . . . . .	70
3.4	Multivariate linear regression: The normal-normal/inverse Wishart model . . . . .	83
3.5	Summary . . . . .	86
3.6	Exercises . . . . .	87
<b>4</b>	<b>Simulation methods</b>	<b>91</b>
4.1	Markov chain Monte Carlo methods . . . . .	91
4.1.1	Gibbs sampler . . . . .	92
4.1.2	Metropolis-Hastings . . . . .	95
4.1.3	Hamiltonian Monte Carlo . . . . .	99
4.2	Importance sampling . . . . .	105
4.3	Particle filtering . . . . .	110
4.4	Convergence diagnostics . . . . .	122
4.4.1	Numerical standard error . . . . .	123
4.4.2	Effective number of simulation draws . . . . .	123
4.4.3	Tests of convergence . . . . .	124
4.4.4	Checking for errors in the posterior simulator . . . . .	125
4.5	Summary . . . . .	131
4.6	Exercises . . . . .	131
<b>II</b>	<b>Regression models: A GUIded toolkit</b>	<b>133</b>
<b>5</b>	<b>Graphical user interface</b>	<b>135</b>
5.1	Introduction . . . . .	135
5.2	Univariate models . . . . .	137
5.3	Multivariate models . . . . .	141
5.4	Time series models . . . . .	142
5.5	Longitudinal/panel models . . . . .	145
5.6	Bayesian model average . . . . .	146
5.7	Help . . . . .	149
5.8	Warning . . . . .	149
<b>6</b>	<b>Univariate models</b>	<b>151</b>
6.1	The Gaussian linear model . . . . .	151
6.2	The logit model . . . . .	155
6.3	The probit model . . . . .	160
6.4	The multinomial probit model . . . . .	162
6.5	The multinomial logit model . . . . .	166
6.6	Ordered probit model . . . . .	170
6.7	Negative binomial model . . . . .	174
6.8	Tobit model . . . . .	178
6.9	Quantile regression . . . . .	181
6.10	Bayesian bootstrap regression . . . . .	184
6.11	Summary . . . . .	186

<i>Contents</i>	ix
6.12 Exercises . . . . .	186
<b>7 Multivariate models</b>	<b>189</b>
7.1 Multivariate regression . . . . .	189
7.2 Seemingly unrelated regression . . . . .	195
7.3 Instrumental variable . . . . .	199
7.4 Multivariate probit model . . . . .	204
7.5 Summary . . . . .	209
7.6 Exercises . . . . .	209
<b>8 Time series models</b>	<b>213</b>
8.1 State-space representation . . . . .	214
8.2 ARMA processes . . . . .	225
8.3 Stochastic volatility models . . . . .	236
8.4 Vector Autoregressive models . . . . .	243
8.5 Summary . . . . .	253
8.6 Exercises . . . . .	254
<b>9 Longitudinal/Panel data models</b>	<b>257</b>
9.1 Normal model . . . . .	258
9.2 Logit model . . . . .	267
9.3 Poisson model . . . . .	270
9.4 Summary . . . . .	276
9.5 Exercises . . . . .	277
<b>10 Bayesian model averaging</b>	<b>281</b>
10.1 Foundation . . . . .	281
10.2 The Gaussian linear model . . . . .	286
10.3 Generalized linear models . . . . .	303
10.4 Calculating the marginal likelihood . . . . .	312
10.4.1 Savage-Dickey density ratio . . . . .	313
10.4.2 Chib's methods . . . . .	314
10.4.3 Gelfand-Dey method . . . . .	314
10.5 Summary . . . . .	321
10.6 Exercises . . . . .	322
<b>III Advanced methods: A brief introduction</b>	<b>327</b>
<b>11 Semi-parametric and non-parametric models</b>	<b>329</b>
11.1 Mixture models . . . . .	329
11.1.1 Finite Gaussian mixtures . . . . .	330
11.1.2 Dirichlet processes . . . . .	346
11.2 Splines . . . . .	360
11.3 Summary . . . . .	372
11.4 Exercises . . . . .	372

<b>12 Bayesian machine learning</b>	<b>375</b>
12.1 Cross-validation and Bayes factors . . . . .	375
12.2 Regularization . . . . .	378
12.2.1 Bayesian LASSO . . . . .	380
12.2.2 Stochastic search variable selection . . . . .	382
12.3 Bayesian additive regression trees . . . . .	386
12.4 Gaussian process . . . . .	394
12.5 Tall data problems . . . . .	403
12.5.1 Divide-and-conquer methods . . . . .	404
12.5.2 Subsampling-based algorithms . . . . .	405
12.6 Summary . . . . .	418
12.7 Exercises . . . . .	418
<b>13 Causal inference</b>	<b>421</b>
13.1 Identification setting . . . . .	421
13.2 Randomized controlled trial (RCT) . . . . .	423
13.3 Conditional independence assumption (CIA) . . . . .	431
13.4 Instrumental variables (IV) . . . . .	438
13.5 Difference-in-differences design (DiD) . . . . .	455
13.5.1 Classic DiD: two-group and two-period ( $2 \times 2$ ) setup . . . . .	456
13.5.2 Staggered (SDiD): G-group and T-period ( $G \times T$ ) setup . . . . .	461
13.6 Regression discontinuity design (RD) . . . . .	465
13.6.1 Sharp regression discontinuity design (SRD) . . . . .	465
13.6.2 Fuzzy regression discontinuity design (FRD) . . . . .	467
13.7 Sample selection . . . . .	477
13.8 Bayesian exponentially tilted empirical likelihood . . . . .	485
13.9 General Bayes posteriors . . . . .	497
13.10 Doubly robust Bayesian inferential framework (DRB) . . . . .	504
13.11 Summary . . . . .	511
13.12 Exercises . . . . .	512
<b>14 Approximate Bayesian methods</b>	<b>515</b>
14.1 Simulation-based approaches . . . . .	515
14.1.1 Approximate Bayesian computation . . . . .	516
14.1.2 Bayesian synthetic likelihood . . . . .	522
14.2 Optimization approaches . . . . .	529
14.2.1 Integrated nested Laplace approximations . . . . .	530
14.2.2 Variational Bayes . . . . .	542
14.3 Summary . . . . .	554
14.4 Exercises . . . . .	554
<b>Bibliography</b>	<b>557</b>
<b>A Appendix</b>	<b>593</b>

---

---

## *Foreword*

---



---

## Preface

---

The main goal of this book is to make the Bayesian inferential framework more approachable to students, researchers, and practitioners who wish to understand and apply this statistical/econometric approach but do not have the time to develop programming skills. I have aimed to strike a balance between applicability and theory. This book provides a very user-friendly graphical user interface (GUI) to implement the most common regression models, while also covering the basic mathematical developments and their code implementation for those interested in advancing to more complex models.

### To instructors and students

This book is divided into three parts: foundations (chapters 1 to 4), regression analysis (chapters 5 to 10), and *Advanced* methods (chapters 11 to 14). Our graphical user interface (GUI) is designed for the second part. The source code can be found at <https://github.com/besmarter/BSTApp>. Instructors and students can access all the code, along with simulated and real datasets. There are four ways to install our GUI:

1. Install *shiny* package, and then, type

```
shiny::runGitHub("besmarter/BSTApp", launch.browser=T)
```

in the **R** console or any **R** code editor and execute it. We recommend to type directly, no copy and paste.

2. Visit <https://andres-ramirez-hassan.shinyapps.io/BSTApp/>.  
Please note: the free Posit Cloud tier sometimes runs out of memory, which can cause the app to stop. Sorry for the inconvenience.
3. Visit <https://fly-besmarter.fly.dev/>. As with Posit Cloud, occasional memory limits on the free tier may affect performance.
4. Use a **Docker** image by typing in the **Command Prompt**

- (a) docker pull aramir21/besmartergui:latest
- (b) docker run -rm -p 3838:3838 aramir21/besmartergui

Then users can access our GUI going to <http://localhost:3838/>. See Chapter 5 for details.

Students should have a basic understanding of probability theory and statistics, as well as some background in econometrics and time series, particularly regression analysis. Familiarity with standard univariate and multivariate probability distributions is strongly recommended. See a nice summary

of useful probability distributions in [184, p. 182-191]. Additionally, students who wish to master the material in this book should have programming skills in **R** software.<sup>1</sup>

I have included both formal and computational exercises at the end of each chapter to help students gain a better understanding of the material presented. A solutions manual for these exercises accompanies this book.

Instructors can use this book as a textbook for a course on introductory Bayesian Econometrics/Statistics, with a strong emphasis on implementation and applications. This book is intended to be complementary, rather than a substitute, for excellent resources on the topic, such as [160], [76], [379], [184], [173], [255], and [250].

### Acknowledgments

I began developing our graphical user interface (GUI) in 2016, after being diagnosed with cervical dystonia. I worked on this side project during weekends, which I called “nerd weekends,” and it served as a form of release from my health condition. Once I began to recover, I invited Mateo Graciano, my former student, business partner, and friend, to join the project. He has been instrumental in developing our GUI, and I am enormously grateful to him.

I would also like to thank the members of the BEsmarter research group at Universidad EAFIT, as well as the NUMBATs members at Monash University, for their valuable feedback and recommendations to improve our GUI.

This book is an extension of the paper *A GUIided tour of Bayesian regression* [363], which serves as a brief user guide for our GUI. I decided to write this book to explain the underlying theory and code in our GUI, and to use it as a textbook in my course on Bayesian econometrics/statistics. I am grateful to my students in this course; their insights and thoughtful questions have deepened my understanding of the material.

I would like to thank Chris Parmeter for his valuable suggestions on how to present our user guide; Professors Raúl Pericchi and Juan Carlos Correa for introducing me to Bayesian statistics; and Liana Jacobi, Tomasz Wozniak, and Chun Fung Kwok (Jackson) from the University of Melbourne, as well as David Frazier from Monash University, for their engaging discussions and fruitful collaborations in Bayesian econometrics and statistics. I am also deeply grateful to Professor Peter Diggle for his unwavering support of my career, and especially to Professor Gael Martin, who gave me the opportunity to work with her and has been a constant source of intellectual inspiration.

I also wish to acknowledge my colleagues and staff at Universidad EAFIT for their continuous support.

Finally, I acknowledge the use of ChatGPT, which assisted me in improving the grammar, clarity, and flow of the text, as well as in tidying up the code presented in this book. Nevertheless, all concepts, mathematical developments, and underlying logic are entirely my own, based on my understanding

---

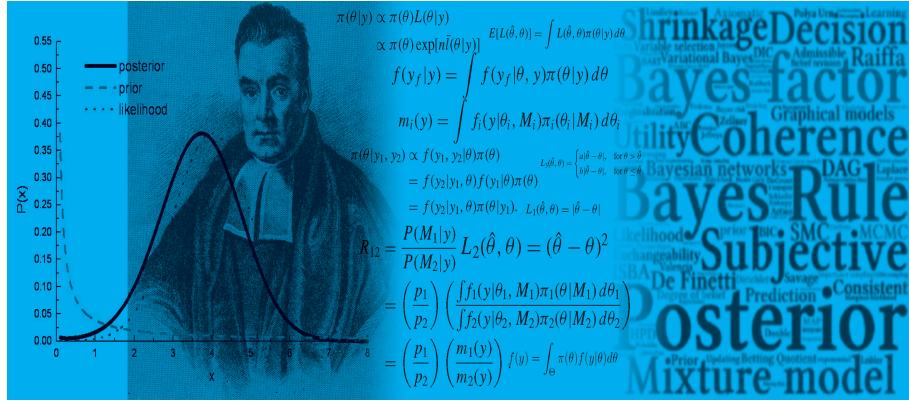
<sup>1</sup>An excellent starting point for **R** programming is the *R Introduction Manual*: <https://cran.r-project.org/doc/manuals/r-release/R-intro.pdf>.

and readings of the literature. Any remaining errors are solely my responsibility, for which I apologize in advance. I sincerely thank the reader and hope that this book proves useful.

To my parents, Orlando and Nancy, who have always been there for me with their unconditional support. They have taught me that the primary aspect of human spiritual evolution is humility, a lesson I am still learning every day. To my wife, Estephania, for her unwavering love and support.



# Introduction



**FIGURE 1**  
Supposedly portrait of Thomas Bayes.

Since the late 90s, Bayesian inference has gained significant popularity among researchers due to the computational revolution and the availability of algorithms to solve complex integrals. However, many researchers, students, and practitioners still lack a deep understanding and practical application of this inferential approach. The primary reason for this is the requirement for strong programming skills.

**Introduction to Bayesian Econometrics: A GUIided Toolkit using R** mainly targets those who want to apply Bayesian inference with a good conceptual and formal understanding, but who may not necessarily have the time to develop programming skills. Thus, this book provides a graphical user interface (GUI) to carry out Bayesian regression in a very user-friendly environment. The book also offers the basic theory and its code implementation using **R** software [345], along with some applications to highlight the potential of Bayesian inference. Additionally, it includes theory and computational exercises for those interested in developing more complex models. In particular, this book presents the mathematical proofs of the basic models step-by-step in the first part, which serve as the foundation for obtaining the most relevant mathematical results in the more complex models covered in the second and third parts.

Our GUI is based on an interactive web application using **shiny** [77],

along with several packages in **R**. Users can estimate univariate, multivariate, time series, longitudinal/panel data, and Bayesian model averaging models using our GUI. In addition, it provides basic summaries, as well as formal and graphical diagnostics of the posterior chains. Our GUI can be run on any operating system and is freely available at <https://github.com/besmarter/BSTApp>.

Users can access simulated and real datasets in the folders **DataSim** and **DataApp**, respectively. The former folder also includes the files used to simulate different processes, so the population parameters are available. As a result, these files can be used as a pedagogical tool to demonstrate various statistical properties. The latter folder contains the datasets used in our applications. Users are encouraged to use these datasets as templates for structuring their own datasets.

This book is divided into three parts. The first part covers the theory (both conceptual and mathematical), programming, and simulation foundations (chapters 1 to 4). The second part focuses on the applications of regression analysis, with particular emphasis on the computational aspect of obtaining draws from the posterior distributions at three levels of programming skills: no skills at all using our GUI, intermediate skills using specialized **R** packages for Bayesian inference, and relatively advanced programming skills to obtain posterior draws from scratch (chapters 5 to 10). The third part provides an introductory treatment of *advanced methods* in Bayesian inference (chapters 11 to 14).

I present some of the mathematical deductions in detail in the first part of the book, whereas I do not show most of the proofs in the second and third parts. However, the same mathematical steps in the first part can be used to derive the results in parts two and three.

In the first part, Chapter 1 begins with an introduction to formal concepts in Bayesian inference, starting with Bayes' rule, its components, their formal definitions, and basic examples. It then presents the basics of Bayesian inference based on decision theory under uncertainty. Chapter 2 discusses the conceptual differences between Bayesian and Frequentist statistical approaches, as well as provides a historical and philosophical perspective on Bayesian statistics and econometrics, highlighting the contrasts with the Frequentist approach. In Chapter 3, I introduce conjugate families in basic statistical models, solving them both analytically and computationally. Simulation-based methods are presented in Chapter 4; these algorithms are crucial in modern Bayesian inference since most realistic models do not have standard forms or analytical solutions.

In the second part, Chapter 5 presents our graphical user interface, and univariate and multivariate regression models are covered in Chapters 6 and 7. Chapter 8 focuses on univariate and multivariate time series models, while Chapter 9 covers Bayesian longitudinal/panel data models. Chapter 10 introduces Bayesian model averaging. In the third part, Chapter 11 explores semi-parametric and non-parametric models, Chapter 12 covers Bayesian methods

in machine learning algorithms, Chapter 13 discusses causal inference, and Chapter 14 describes approximation methods.

#### About me

My name is Andrés Ramírez-Hassan, and I am an applied and theoretical econometrician working as a Distinguished Professor in the School of Finance, Economics, and Government at Universidad EAFIT (Medellín, Colombia). I hold a PhD in Statistical Science, a Master's degree in Finance, another in Economics, and a Bachelor's degree in Economics. I was a research fellow at the Department of Econometrics and Business Statistics at Monash University, and a visiting Professor in the Department of Economics at the University of Melbourne and the University of Glasgow.

Since completing my PhD, much of my research has focused on Bayesian Econometrics, with applications in crime, finance, health, sports, and utilities. My work has been published (or is forthcoming) in highly regarded journals such as the *International Journal of Forecasting*, *Journal of Applied Econometrics*, *Econometric Reviews*, *Journal of Computational and Graphical Statistics*, *The R Journal*, *Economic Modelling*, *Spatial Economic Analysis*, *Economic Inquiry*, *World Development*, *Journal of Sport Economics*, *Empirical Economics*, *Australian and New Zealand Journal of Statistics*, *Brazilian Journal of Probability and Statistics*, and other prestigious international research outlets.

I founded **BEmarter** – Bayesian Econometrics: simulations, models, and applications to research, teaching, and encoding with responsibility. This is a research group whose **mission** is to *lead and excel in the generation and dissemination of Bayesian Econometric knowledge through research, teaching, and software*. We **envision** worldwide econometric research, teaching, and applications based on the Bayesian framework that:

- Inspires new econometric ideas,
- Creates a user-friendly environment for applications of Bayesian econometrics,
- Transforms classic econometric research, teaching, and applications,
- where one of the main concerns of science is to solve social problems.

**Contact:** aramir21@gmail.com / aramir21@eafit.edu.co

**Website:** <https://sites.google.com/view/arh-bayesian>



**FIGURE 2**

This book is licensed under the **Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License**.



---

## Symbols

---

### Symbol Description

$\neg$	Negation symbol	$\mathbf{0}_l$	$l$ -dimensional null vector
$\propto$	Proportional symbol	$\max_r$	Maximum over $r$
$\perp$	Independence symbol	$\min_r$	Minimum over $r$
$\mathcal{R}$	The Real set	$\nabla$	Gradiente operator
$\mathcal{R}^K$	Euclidean space of dimension $K$	$\stackrel{iid}{\sim}$	Independently and identically distributed
$\emptyset$	Empty set	$\stackrel{ind}{\sim}$	Independently and Not identically distributed
$\mathbb{1}$	Indicator function	$>$	Greater than
$P$	Probability measure	$<$	Less than
$:=$	Is defined as	$\approx$	Approximately equal to
$\text{argmax}$	Argument of the maximum	$\gtrsim$	Greater than or approximately equal to
$\text{argmin}$	Argument of the minimum	$\Delta$	Difference operator
$tr$	Trace operator	$\subseteq$	Subset
$vec$	Vectorization operator	$\subset$	Proper subset
$\lim$	Limit	$\xrightarrow{d}$	Convergence in distribution (law)
$\otimes$	Kronecker product	$\xrightarrow{p}$	Convergence in probability
$\text{diag}\{\cdot\}$	Diagonal matrix	$\xrightarrow{a.s.}$	Almost surely convergence
$\dim\{\cdot\}$	Dimension of an object		
$\mathbf{I}_l$	$l$ -dimensional identity matrix		
$\mathbf{i}_N$	$N$ -dimensional vector of ones		



# **Part I**

## **Foundations: Theory, simulation methods and programming**



# 1

---

## *Basic formal concepts*

---

We introduce formal concepts in Bayesian inference, beginning with Bayes' rule and its components, along with their formal definitions and basic examples. In addition, we present key features of Bayesian inference, such as Bayesian updating and asymptotic sampling properties. We also cover the basics of Bayesian inference from a decision-theoretic perspective under uncertainty, introducing important concepts like loss functions, risk functions, and optimal decision rules.

---

### 1.1 The Bayes' rule

As expected, the starting point for performing Bayesian inference is Bayes' rule,<sup>1</sup> which provides the solution to the inverse problem of determining causes from observed effects. This rule combines prior beliefs with objective probabilities based on repeatable experiments, allowing us to move from observations to probable causes.

Formally, the conditional probability of  $A_i$  given  $B$  is equal to the conditional probability of  $B$  given  $A_i$ , multiplied by the marginal probability of  $A_i$ , divided by the marginal probability of  $B$ ,

$$\begin{aligned} P(A_i | B) &= \frac{P(A_i, B)}{P(B)} \\ &= \frac{P(B | A_i) \times P(A_i)}{P(B)}, \end{aligned} \tag{1.1}$$

where, by the law of total probability,  $P(B) = \sum_i P(B | A_i)P(A_i) \neq 0$ , and  $\{A_i, i = 1, 2, \dots\}$  is a finite or countably infinite partition of the sample space.

In the Bayesian framework,  $B$  represents sample information that updates a probabilistic statement about an unknown object  $A_i$  according to probability rules. This is done using Bayes' rule, which incorporates prior "beliefs"

---

<sup>1</sup>Note that I use the term "Bayes' rule" rather than "Bayes' theorem." It was Laplace [257] who generalized Bayes' theorem [33], and his generalization is referred to as Bayes' rule.

about  $A_i$ , i.e.,  $P(A_i)$ , sample information relating  $B$  to the particular state of nature  $A_i$  through a probabilistic statement,  $P(B | A_i)$ , and the probability of observing that specific sample information,  $P(B)$ .

Let's consider a simple example, *the base rate fallacy*:

Assume that the sample information comes from a positive result from a test whose true positive rate (sensitivity) is 98%, i.e.,  $P(+ | \text{disease}) = 0.98$ . In addition, the false positive rate is 2%, i.e.,  $P(+ | \neg\text{disease}) = 0.02$ . On the other hand, the prior probability of being infected with the disease is given by the base incidence rate,  $P(\text{disease}) = 0.002$ . The question is: *What is the probability of actually being infected, given a positive test result?*

This is an example of *the base rate fallacy*, where a positive test result for a disease with a very low base incidence rate still gives a low probability of actually having the disease.

The key to answering this question lies in understanding the difference between the probability of having the disease given a positive test result,  $P(\text{disease} | +)$ , and the probability of a positive result given the disease,  $P(+ | \text{disease})$ . The former is the crucial result, and Bayes' rule helps us to compute it. Using Bayes' rule (equation 1.1):

$$\begin{aligned} P(\text{disease} | +) &= \frac{P(+ | \text{disease}) \times P(\text{disease})}{P(+)} \\ &= \frac{0.98 \times 0.002}{0.98 \times 0.002 + 0.02 \times (1 - 0.002)} \\ &= 0.09, \end{aligned}$$

where  $P(+) = P(+ | \text{disease}) \times P(\text{disease}) + P(+ | \neg\text{disease}) \times P(\neg\text{disease})$ .<sup>2</sup>

The following code shows how to perform this exercise in **R**.

### **R code. The base rate fallacy**

```
1 PD <- 0.002 # Probability of disease
2 PPD <- 0.98 # True positive (Sensitivity)
3 PPF <- 0.02 # False positive
4 PDP <- PD * PPD / (PD * PPD + (1 - PD) * PPF)
5 paste("Probability of disease given a positive test is", sep
      = " ", round(PDP, 2))
6 "Probability of disease given a positive test is 0.09"
```

<sup>2</sup>  $\neg$  is the negation symbol.

We observe that despite having a positive result, the probability of actually having the disease remains low. This is due to the base rate being so small.

Another interesting example—one that lies at the heart of the origin of Bayes' theorem [33]—concerns the existence of God [424]. In Section X of David Hume's *An Inquiry concerning Human Understanding* (1748), titled *Of Miracles*, Hume argues that when someone claims to have witnessed a miracle, the claim provides weak evidence that the event actually occurred, since it contradicts our everyday experience. In response, Richard Price, who edited and published *An Essay Towards Solving a Problem in the Doctrine of Chances* in 1763 (following Bayes' death in 1761), criticizes Hume's argument by distinguishing between *logical* and *physical* impossibility. He illustrates this with the example of a die with a million sides: while rolling a specific number might seem impossible, it is merely *improbable*, whereas rolling a number not present on the die is *physically impossible*. The former may occur with enough trials; the latter never will.

*Note on the following example:*

*The next example is adapted from a modern blog post that illustrates the base rate fallacy using a resurrection scenario. While it is not taken from the original writings of Hume or Price, it reflects themes central to their philosophical debate on probability and miracles. It is included here to demonstrate how Bayes' rule behaves in cases involving extremely low prior probabilities. References to figures such as Jesus and Elvis are used purely for illustrative and pedagogical purposes. The goal is to explore the statistical reasoning, not to make any theological claims. Readers are encouraged to focus on the mathematical and conceptual content of the example.*

To illustrate the statistical implications of this discussion, let us consider the following scenario.<sup>3</sup>

The scenario involves two reported cases of resurrection: Jesus Christ and Elvis Presley. The base rate is therefore two out of the total number of people who have ever lived, estimated at approximately 117 billion.<sup>4</sup> That is,

$$P(\text{Res}) = \frac{2}{117 \times 10^9}.$$

Suppose now that the sample information comes from a highly reliable witness, with a true positive rate of 0.9999999. Unlike the original blog post, we assume a more conservative false positive rate of 50%.

We ask: *What is the probability that a resurrection actually occurred, given that a witness claims it did?*

Using Bayes' rule, and letting *Res* denote the event of resurrection, and *Witness* the event of a witness declaring a resurrection:

---

<sup>3</sup><https://blog.ephorie.de/base-rate-fallacy-or-why-no-one-is-justified-to-believe-that-jesus-rose>

<sup>4</sup>[https://www.prb.org/articles/how-many-people-have-ever-lived-on-earth/?utm\\_source=chatgpt.com](https://www.prb.org/articles/how-many-people-have-ever-lived-on-earth/?utm_source=chatgpt.com)

$$\begin{aligned}
 P(\text{Res} \mid \text{Witness}) &= \frac{P(\text{Witness} \mid \text{Res}) \times P(\text{Res})}{P(\text{Witness})} \\
 &= \frac{0.9999999 \times \frac{2}{117 \times 10^9}}{0.9999999 \times \frac{2}{117 \times 10^9} + 0.5 \times \left(1 - \frac{2}{117 \times 10^9}\right)} \\
 &\approx 3.42 \times 10^{-11}
 \end{aligned}$$

Here, the denominator is the marginal probability of a witness reporting a resurrection:

$$P(\text{Witness}) = P(\text{Witness} \mid \text{Res}) \times P(\text{Res}) + P(\text{Witness} \mid \neg\text{Res}) \times (1 - P(\text{Res}))$$

Thus, the probability that a resurrection actually occurred, even when reported by an extremely reliable witness, is approximately

$$3.42 \times 10^{-11}.$$

This value is exceedingly small, but importantly, it is not zero. A value of zero would imply *impossibility*, whereas this result reflects *extremely low probability*.

The following code shows how to perform this exercise in **R**.

### **R code. Of Miracles**

```

1 PR <- 2/(117 * 10^9) # Probability of resurrection
2 PWR <- 0.9999999 # True positive rate
3 PWNR <- 0.5 # False positive
4 PRW <- PR * PWR / (PR * PWR + (1 - PR) * PWNR) # Probability
   of resurrection given witness
5 paste("Probability of resurrection given witness is", sep =
   " ", PRW)
6 "Probability of resurrection given witness is
   3.41880307686464e-11"

```

Observe that we can condition on multiple events in Bayes' rule. Let's consider two conditioning events,  $B$  and  $C$ . Then, equation 1.1 becomes

$$\begin{aligned}
 P(A_i \mid B, C) &= \frac{P(A_i, B, C)}{P(B, C)} \\
 &= \frac{P(B \mid A_i, C) \times P(A_i \mid C) \times P(C)}{P(B \mid C)P(C)}. \tag{1.2}
 \end{aligned}$$

Let's use this rule in one of the most intriguing statistical puzzles, *the*

**FIGURE 1.1**

The Monty Hall show.

*Monty Hall problem*, to illustrate how to use equation 1.2 [404, 111]. This was the situation faced by a contestant in the American television game show *Let's Make a Deal*. In this game, the contestant was asked to choose a door, behind one of which there is a car, and behind the others, goats.

Let's say the contestant picks door No. 1, and the host (Monty Hall), who knows what is behind each door, opens door No. 3, revealing a goat (see Figure 1.1). Then, the host asks the tricky question: *Do you want to pick door No. 2?*

Let's define the following events:

- $P_i$ : the event **contestant picks door No.  $i$** , which stays closed,
- $H_i$ : the event **host picks door No.  $i$** , which is open and contains a goat,
- $C_i$ : the event **car is behind door No.  $i$** .

In this particular setting, the contestant is interested in the probability of the event  $P(C_2 | H_3, P_1)$ . A naive answer would be that it is irrelevant, as initially,  $P(C_i) = \frac{1}{3}$ ,  $i = 1, 2, 3$ , and now  $P(C_i | H_3) = \frac{1}{2}$ ,  $i = 1, 2$ , since the host opened door No. 3. So, why bother changing the initial guess if the odds are the same (1:1)?

The important point here is that the host knows what is behind each door and always picks a door where there is a goat, given the contestant's choice. In this particular setting:

$$P(H_3 | C_3, P_1) = 0, \quad P(H_3 | C_2, P_1) = 1, \quad P(H_3 | C_1, P_1) = \frac{1}{2}.$$

Then, using equation 1.2, we can calculate the posterior probability.

$$\begin{aligned} P(C_2 | H_3, P_1) &= \frac{P(C_2, H_3, P_1)}{P(H_3, P_1)} \\ &= \frac{P(H_3 | C_2, P_1)P(C_2 | P_1)P(P_1)}{P(H_3 | P_1) \times P(P_1)} \\ &= \frac{P(H_3 | C_2, P_1)P(C_2)}{P(H_3 | P_1)} \\ &= \frac{1 \times 1/3}{1/2}, \end{aligned}$$

where the third equation uses the fact that  $C_i$  and  $P_i$  are independent events, and  $P(H_3 | P_1) = \frac{1}{2}$  because this depends only on  $P_1$  (not on  $C_2$ ).

Therefore, changing the initial decision increases the probability of getting the car from  $\frac{1}{3}$  to  $\frac{2}{3}$ ! Thus, it is always a good idea to change the door.

Let's see a simulation exercise in R to check this answer:

### R code. The Monty Hall problem

```

1 set.seed(0101) # Set simulation seed
2 S <- 100000 # Simulations
3 Game <- function(switch = 0){
4   # switch = 0 is not change / switch = 1 is to change
5   opts <- 1:3
6   car <- sample(opts, 1) # car location
7   guess1 <- sample(opts, 1) # Initial guess
8   if(car != guess1) {
9     host <- opts[-c(car, guess1)]
10 } else {
11   host <- sample(opts[-c(car, guess1)], 1)
12 }
13 win1 <- guess1 == car # Win no change
14 guess2 <- opts[-c(host, guess1)]
15 win2 <- guess2 == car # Win change
16 if(switch == 0){
17   win <- win1
18 } else {win <- win2}
19 return(win)
20 }
21 Prob <- mean(replicate(S, Game(switch = 0))) # Not changing
22 Prob <- mean(replicate(S, Game(switch = 1))) # Changing

```

---

## 1.2 Bayesian framework: A brief summary of theory

Given an unknown parameter set  $\boldsymbol{\theta}$ , and a particular realization of the data  $\mathbf{y}$ , Bayes' rule may be applied analogously,<sup>5</sup>

$$\pi(\boldsymbol{\theta} | \mathbf{y}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})}{p(\mathbf{y})}, \quad (1.3)$$

where  $\pi(\boldsymbol{\theta} | \mathbf{y})$  is the posterior density function,  $\pi(\boldsymbol{\theta})$  is the prior density,  $p(\mathbf{y} | \boldsymbol{\theta})$  is the likelihood (econometric/statistical model), and

$$p(\mathbf{y}) = \int_{\Theta} p(\mathbf{y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} = \mathbb{E}[p(\mathbf{y} | \boldsymbol{\theta})] \quad (1.4)$$

is the marginal likelihood, prior predictive or evidence. Observe that for this expected value to be meaningful, the prior should be a proper density, that is, it must integrate to one; otherwise, it does not make sense.

Observe that  $p(\mathbf{y} | \boldsymbol{\theta})$  is not a density in  $\boldsymbol{\theta}$ . In addition,  $\pi(\boldsymbol{\theta})$  does not have to integrate to 1, that is,  $\pi(\boldsymbol{\theta})$  can be an improper density function,  $\int_{\Theta} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} = \infty$ . However,  $\pi(\boldsymbol{\theta} | \mathbf{y})$  is a proper density function, that is,  $\int_{\Theta} \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} = 1$ .

For instance, set  $\pi(\boldsymbol{\theta}) = c$ , where  $c$  is a constant, then  $\int_{\Theta} cd\boldsymbol{\theta} = \infty$ . However,

$$\int_{\Theta} \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} = \int_{\Theta} \frac{p(\mathbf{y} | \boldsymbol{\theta}) \times c}{\int_{\Theta} p(\mathbf{y} | \boldsymbol{\theta}) \times c d\boldsymbol{\theta}} d\boldsymbol{\theta} = 1$$

where  $c$  cancels out.

$\pi(\boldsymbol{\theta} | \mathbf{y})$  is a sample updated “probabilistic belief” version of  $\pi(\boldsymbol{\theta})$ , where  $\pi(\boldsymbol{\theta})$  is a prior probabilistic belief which can be constructed from previous empirical work, theoretical foundations, expert knowledge, and/or mathematical convenience. This prior usually depends on parameters, which are named *hyperparameters*. In addition, the Bayesian approach implies using a probabilistic model about  $\mathbf{Y}$  given  $\boldsymbol{\theta}$ , that is,  $p(\mathbf{y} | \boldsymbol{\theta})$ , where its integral over  $\Theta$ ,  $p(\mathbf{y})$ , is named *the model evidence* due to being a measure of model fit to the data.

Observe that the Bayesian inferential approach is conditional, that is, what can we learn about an unknown object  $\boldsymbol{\theta}$  given that we already observed  $\mathbf{Y} = \mathbf{y}$ ? The answer is also conditional on the probabilistic model, that is,  $p(\mathbf{y} | \boldsymbol{\theta})$ . So, what if we want to compare different models, say  $\mathcal{M}_m$ ,  $m = \{1, 2, \dots, M\}$ ? Then, we should make explicit this in the Bayes' rule formulation:

$$\pi(\boldsymbol{\theta} | \mathbf{y}, \mathcal{M}_m) = \frac{p(\mathbf{y} | \boldsymbol{\theta}, \mathcal{M}_m) \times \pi(\boldsymbol{\theta} | \mathcal{M}_m)}{p(\mathbf{y} | \mathcal{M}_m)}. \quad (1.5)$$

---

<sup>5</sup>From a Bayesian perspective,  $\boldsymbol{\theta}$  is fixed but unknown. Then, it is treated as a random object despite the lack of variability (see Chapter 2).

The posterior model probability is

$$\pi(\mathcal{M}_m | \mathbf{y}) = \frac{p(\mathbf{y} | \mathcal{M}_m) \times \pi(\mathcal{M}_m)}{p(\mathbf{y})}, \quad (1.6)$$

where  $p(\mathbf{y} | \mathcal{M}_m) = \int_{\boldsymbol{\Theta}} p(\mathbf{y} | \boldsymbol{\theta}, \mathcal{M}_m) \times \pi(\boldsymbol{\theta} | \mathcal{M}_m) d\boldsymbol{\theta}$  due to equation 1.5, and  $\pi(\mathcal{M}_m)$  is the prior model probability.

Calculating  $p(\mathbf{y})$  in equations 1.3 and 1.6 is very demanding in most of the realistic cases. Fortunately, it is not required when performing inference about  $\boldsymbol{\theta}$  as this is integrated out from it. Then, all you need to know about the shape of  $\boldsymbol{\theta}$  is in  $p(\mathbf{y} | \boldsymbol{\theta}, \mathcal{M}_m) \times \pi(\boldsymbol{\theta} | \mathcal{M}_m)$ , or without explicitly conditioning on  $\mathcal{M}_m$ ,

$$\pi(\boldsymbol{\theta} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta}). \quad (1.7)$$

Equation 1.7 is a very good shortcut to perform Bayesian inference about  $\boldsymbol{\theta}$ .

We can also avoid calculating  $p(\mathbf{y})$  when performing model selection (hypothesis testing) using the posterior odds ratio, that is, comparing models  $\mathcal{M}_1$  and  $\mathcal{M}_2$ ,

$$\begin{aligned} PO_{12} &= \frac{\pi(\mathcal{M}_1 | \mathbf{y})}{\pi(\mathcal{M}_2 | \mathbf{y})} \\ &= \frac{p(\mathbf{y} | \mathcal{M}_1)}{p(\mathbf{y} | \mathcal{M}_2)} \times \frac{\pi(\mathcal{M}_1)}{\pi(\mathcal{M}_2)}, \end{aligned} \quad (1.8)$$

where the first term in equation 1.8 is named the *Bayes factor*, and the second term is the *prior odds*. Observe that the Bayes factor is a ratio of ordinates for  $\mathbf{y}$  under different models. Then, the Bayes factor is a measure of relative sample evidence in favor of model 1 compared to model 2.

However, we still need to calculate  $p(\mathbf{y} | \mathcal{M}_m) = \int_{\boldsymbol{\Theta}} p(\mathbf{y} | \boldsymbol{\theta}, \mathcal{M}_m) \pi(\boldsymbol{\theta} | \mathcal{M}_m) d\boldsymbol{\theta} = \mathbb{E}[p(\mathbf{y} | \boldsymbol{\theta}, \mathcal{M}_m)]$ . For this integral to be meaningful, the prior must be proper. Using an improper prior has unintended consequences when comparing models; for instance, parsimonious models are favored by posterior odds or Bayes factors, and these values may depend on units of measure (see Chapter 3).

A nice feature of comparing models using posterior odds is that if we have an exhaustive set of competing models such that  $\sum_{m=1}^M \pi(\mathcal{M}_m | \mathbf{y}) = 1$ , then we can recover  $\pi(\mathcal{M}_m | \mathbf{y})$  without calculating  $p(\mathbf{y})$ . In particular, given two models  $\mathcal{M}_1$  and  $\mathcal{M}_2$  such that  $\pi(\mathcal{M}_1 | \mathbf{y}) + \pi(\mathcal{M}_2 | \mathbf{y}) = 1$ , we have:

$$\pi(\mathcal{M}_1 | \mathbf{y}) = \frac{PO_{12}}{1 + PO_{12}} \quad \text{and} \quad \pi(\mathcal{M}_2 | \mathbf{y}) = 1 - \pi(\mathcal{M}_1 | \mathbf{y}).$$

In general,

$$\pi(\mathcal{M}_m | \mathbf{y}) = \frac{p(\mathbf{y} | \mathcal{M}_m) \times \pi(\mathcal{M}_m)}{\sum_{l=1}^M p(\mathbf{y} | \mathcal{M}_l) \times \pi(\mathcal{M}_l)}.$$

**TABLE 1.1**  
Kass and Raftery guidelines.

$2 \times \log(PO_{12})$	$PO_{12}$	Evidence against $\mathcal{M}_2$
0 to 2	1 to 3	Not worth more than a bare mention
2 to 6	3 to 20	Positive
6 to 10	20 to 150	Strong
$> 10$	$> 150$	Very strong

Notes: [241] proposed these guidelines for model selection using posterior odds in a Bayesian framework.

These posterior model probabilities can be used to perform Bayesian model averaging.

Table 1.1 shows guidelines for the interpretation of  $2 \log(PO_{12})$  [241]. This transformation is done to replicate the structure of the likelihood ratio test statistic. However, posterior odds do not require nested models as the likelihood ratio test does.

Observe that the posterior odds ratio is a relative criterion, that is, we specify an exhaustive set of competing models and compare them. However, we may want to check the performance of a model on its own or use a non-informative prior. In this case, we can use *the posterior predictive p-value* [158, 159].<sup>6</sup>

The intuition behind the predictive p-value is simple: analyze the discrepancy between the model's assumptions and the data by checking a potential extreme tail-area probability. Observe that this approach does not check if a model is true; its focus is on potential discrepancies between the model and the data at hand.

This is done by simulating pseudo-data from our sampling model  $(\mathbf{y}^{(s)}, s = 1, 2, \dots, S)$  using draws from the posterior distribution, and then calculating a discrepancy measure,  $D(\mathbf{y}^{(s)}, \boldsymbol{\theta})$ , to estimate the posterior predictive p-value,

$$p_D(\mathbf{y}) = P[D(\mathbf{y}^{(s)}, \boldsymbol{\theta}) \geq D(\mathbf{y}, \boldsymbol{\theta})],$$

using the proportion of the  $S$  draws for which  $D(\mathbf{y}^{(s)}, \boldsymbol{\theta}^{(s)}) \geq D(\mathbf{y}, \boldsymbol{\theta}^{(s)})$ . Extreme tail probabilities ( $p_D(\mathbf{y}) \leq 0.05$  or  $p_D(\mathbf{y}) \geq 0.95$ ) suggest potential discrepancies between the data and the model. [159] also suggests the posterior predictive p-value based on the *minimum discrepancy*,

$$D_{\min}(\mathbf{y}) = \min_{\boldsymbol{\theta}} D(\mathbf{y}, \boldsymbol{\theta}),$$

and the *average discrepancy* statistic

$$D(\mathbf{y}) = \mathbb{E}[D(\mathbf{y}, \boldsymbol{\theta})] = \int_{\Theta} D(\mathbf{y}, \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}.$$

---

<sup>6</sup>[30] show potential issues due to using data twice in the construction of the predictive p-values. They also present alternative proposals, for instance, *the partial posterior predictive p-value*.

These alternatives can be more computationally demanding.

The Bayesian approach is also suitable to get probabilistic predictions, that is, we can obtain a posterior predictive density

$$\begin{aligned}\pi(\mathbf{y}_0 \mid \mathbf{y}, \mathcal{M}_m) &= \int_{\Theta} \pi(\mathbf{y}_0, \boldsymbol{\theta} \mid \mathbf{y}, \mathcal{M}_m) d\boldsymbol{\theta} \\ &= \int_{\Theta} \pi(\mathbf{y}_0 \mid \boldsymbol{\theta}, \mathbf{y}, \mathcal{M}_m) \pi(\boldsymbol{\theta} \mid \mathbf{y}, \mathcal{M}_m) d\boldsymbol{\theta}.\end{aligned}\quad (1.9)$$

Observe that equation 1.9 is again an expectation  $\mathbb{E}[\pi(\mathbf{y}_0 \mid \boldsymbol{\theta}, \mathbf{y}, \mathcal{M}_m)]$ , this time using the posterior distribution.<sup>7</sup> Therefore, the Bayesian approach takes estimation error into account when performing prediction.

As we have shown many times, expectation (integration) is a common feature in Bayesian inference. That is why the remarkable relevance of computation based on *Monte Carlo integration* in the Bayesian framework (see Chapter 4).

*Bayesian model averaging* (BMA) allows for considering model uncertainty in prediction or any unknown probabilistic object. In the case of the predictive density,

$$\pi(\mathbf{y}_0 \mid \mathbf{y}) = \sum_{m=1}^M \pi(\mathcal{M}_m \mid \mathbf{y}) \pi(\mathbf{y}_0 \mid \mathbf{y}, \mathcal{M}_m). \quad (1.10)$$

In the case of the posterior density of the parameters,

$$\pi(\boldsymbol{\theta} \mid \mathbf{y}) = \sum_{m=1}^M \pi(\mathcal{M}_m \mid \mathbf{y}) \pi(\boldsymbol{\theta} \mid \mathbf{y}, \mathcal{M}_m), \quad (1.11)$$

where

$$\mathbb{E}(\boldsymbol{\theta} \mid \mathbf{y}) = \sum_{m=1}^M \hat{\boldsymbol{\theta}}_m \pi(\mathcal{M}_m \mid \mathbf{y}), \quad (1.12)$$

and

$$\begin{aligned}Var(\theta_k \mid \mathbf{y}) &= \sum_{m=1}^M \pi(\mathcal{M}_m \mid \mathbf{y}) \widehat{Var}(\theta_{km} \mid \mathbf{y}, \mathcal{M}_m) \\ &\quad + \sum_{m=1}^M \pi(\mathcal{M}_m \mid \mathbf{y}) (\hat{\theta}_{km} - \mathbb{E}[\theta_{km} \mid \mathbf{y}])^2,\end{aligned}\quad (1.13)$$

---

<sup>7</sup>Computing expectations is a fundamental aspect of Bayesian inference. See [297] for a comprehensive review of the evolution of computational methods for this task.

$\hat{\theta}_m$  is the posterior mean and  $\widehat{Var}(\theta_{km} \mid \mathbf{y}, \mathcal{M}_m)$  is the posterior variance of the  $k$ -th element of  $\boldsymbol{\theta}$  under model  $\mathcal{M}_m$ .

Observe how the variance in equation 1.13 captures the extra variability due to potential differences between the mean posterior estimates associated with each model, and the posterior mean that incorporates model uncertainty in equation 1.12.

A significant advantage of the Bayesian approach, which is particularly useful in *state space representations* (see Chapter 8), is the way the posterior distribution updates with new sample information. Given  $\mathbf{y} = \mathbf{y}_{1:t+1}$  as a sequence of observations from 1 to  $t+1$ , then

$$\begin{aligned}\pi(\boldsymbol{\theta} \mid \mathbf{y}_{1:t+1}) &\propto p(\mathbf{y}_{1:t+1} \mid \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta}) \\ &= p(y_{t+1} \mid \mathbf{y}_{1:t}, \boldsymbol{\theta}) \times p(\mathbf{y}_{1:t} \mid \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta}) \\ &\propto p(y_{t+1} \mid \mathbf{y}_{1:t}, \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta} \mid \mathbf{y}_{1:t}).\end{aligned}\quad (1.14)$$

We observe in Equation 1.14 that the new prior is simply the posterior distribution based on the previous observations. This is particularly useful under the assumption of *conditional independence*, that is,  $Y_{t+1} \perp \mathbf{Y}_{1:t} \mid \boldsymbol{\theta}$ , so that  $p(y_{t+1} \mid \mathbf{y}_{1:t}, \boldsymbol{\theta}) = p(y_{t+1} \mid \boldsymbol{\theta})$ , allowing the posterior to be recovered recursively [336]. This facilitates online updating because all information up to time  $t$  is captured in  $\boldsymbol{\theta}$ . Therefore,  $\pi(\boldsymbol{\theta} \mid \mathbf{y}_{1:t+1}) \propto p(y_{t+1} \mid \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta} \mid \mathbf{y}_{1:t}) \propto \prod_{h=1}^{t+1} p(y_h \mid \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})$ . This recursive expression can be computed more efficiently at any specific point in time  $t$ , compared to a batch-mode algorithm, which requires processing all information up to time  $t$  simultaneously.

It is also important to consider the sampling properties of “Bayesian estimators”. This topic has attracted the attention of statisticians and econometricians for a long time. For instance, asymptotic posterior concentration on the population parameter vector is discussed by [44]. The convergence of posterior distributions is stated by the Bernstein-von Mises theorem [260, 439], which establishes a link between *credible intervals (sets)* and confidence intervals (sets), where a credible interval is an interval in the domain of the posterior distribution within which an unknown parameter falls with a particular probability. Credible intervals treat bounds as fixed and parameters as random, whereas confidence intervals reverse this. There are many settings in parametric models where Bayesian credible intervals with an  $\alpha$  level converge asymptotically to confidence intervals at the  $\alpha$  level. This suggests that Bayesian inference is asymptotically correct from a sampling perspective in these settings.

A heuristic approach to demonstrate this in the simplest case, where we assume random sampling and  $\theta \in \mathcal{R}$ , is the following:  $p(\mathbf{y} \mid \boldsymbol{\theta}) = \prod_{i=1}^N p(y_i \mid \boldsymbol{\theta})$ , so the log likelihood is  $l(\mathbf{y} \mid \boldsymbol{\theta}) \equiv \log p(\mathbf{y} \mid \boldsymbol{\theta}) = \sum_{i=1}^N \log p(y_i \mid \boldsymbol{\theta}) = N \times \bar{l}(\mathbf{y} \mid \boldsymbol{\theta})$ , where  $\bar{l} \equiv \frac{1}{N} \sum_{i=1}^N \log p(y_i \mid \boldsymbol{\theta})$  is the mean likelihood.<sup>8</sup> Then, the posterior distribution is proportional to

---

<sup>8</sup>Note that in the likelihood function the argument is  $\boldsymbol{\theta}$ , but we keep the notation for convenience in exposition.

$$\begin{aligned}\pi(\theta \mid \mathbf{y}) &\propto p(\mathbf{y} \mid \theta) \times \pi(\theta) \\ &= \exp\{\bar{l}(\mathbf{y} \mid \theta)\} \times \pi(\theta).\end{aligned}\quad (1.15)$$

Observe that as the sample size increases, that is, as  $N \rightarrow \infty$ , the exponential term should dominate the prior distribution as long as the prior does not depend on  $N$ , such that the likelihood determines the posterior distribution asymptotically.

*Maximum likelihood* theory shows that  $\lim_{N \rightarrow \infty} \bar{l}(\mathbf{y} \mid \theta) \rightarrow \bar{l}(\mathbf{y} \mid \theta_0)$ , where  $\theta_0$  is the population parameter of the data-generating process. In addition, performing a second-order Taylor expansion of the log likelihood at the maximum likelihood estimator,

$$\begin{aligned}l(\mathbf{y} \mid \theta) &\approx l(\mathbf{y} \mid \hat{\theta}) + \frac{dl(\mathbf{y} \mid \theta)}{d\theta}\Big|_{\hat{\theta}} (\theta - \hat{\theta}) + \frac{1}{2} \frac{d^2l(\mathbf{y} \mid \theta)}{d\theta^2}\Big|_{\hat{\theta}} (\theta - \hat{\theta})^2 \\ &= l(\mathbf{y} \mid \hat{\theta}) + \frac{1}{2} \sum_{i=1}^N \frac{d^2l(y_i \mid \theta)}{d\theta^2}\Big|_{\hat{\theta}} (\theta - \hat{\theta})^2 \\ &= l(\mathbf{y} \mid \hat{\theta}) - \frac{1}{2} N [-\bar{l}''\Big|_{\hat{\theta}}] (\theta - \hat{\theta})^2 \\ &= l(\mathbf{y} \mid \hat{\theta}) - \frac{N}{2\sigma^2} (\theta - \hat{\theta})^2\end{aligned}$$

where  $\frac{dl(\mathbf{y} \mid \theta)}{d\theta}\Big|_{\hat{\theta}} = 0$ ,  $\bar{l}'' \equiv \frac{1}{N} \sum_{i=1}^N \frac{d^2l(y_i \mid \theta)}{d\theta^2}\Big|_{\hat{\theta}}$  and  $\sigma^2 := [-\bar{l}''\Big|_{\hat{\theta}}]^{-1}$ .<sup>9</sup> Then,

$$\begin{aligned}\pi(\theta \mid \mathbf{y}) &\propto \exp\{l(\mathbf{y} \mid \theta)\} \times \pi(\theta) \\ &\approx \exp\left\{l(\mathbf{y} \mid \hat{\theta}) - \frac{N}{2\sigma^2} (\theta - \hat{\theta})^2\right\} \times \pi(\theta) \\ &\propto \exp\left\{-\frac{N}{2\sigma^2} (\theta - \hat{\theta})^2\right\} \times \pi(\theta)\end{aligned}$$

Observe that the posterior density is proportional to the kernel of a normal density with mean  $\hat{\theta}$  and variance  $\sigma^2/N$ , as long as  $\pi(\hat{\theta}) \neq 0$ . This kernel dominates as the sample size increases due to the  $N$  in the exponential term. It is also important to note that the prior should not exclude values of  $\theta$  that are logically possible, such as  $\hat{\theta}$ .

---

<sup>9</sup>The last definition follows from standard theory in maximum likelihood estimation (see [73, Chap. 10] and [455, Chap. 13]).

### 1.2.1 Example: Health insurance

Suppose that you are analyzing whether to buy health insurance next year. To make a better decision, you want to know *what is the probability that you will visit your doctor at least once next year?* To answer this question, you have records of the number of times you have visited your doctor over the last 5 years,  $\mathbf{y} = \{0, 3, 2, 1, 0\}$ . How should you proceed?

Assuming that this is a random sample<sup>10</sup> from a data-generating process (statistical model) that is Poisson, i.e.,  $Y_i \sim P(\lambda)$ , and your probabilistic prior beliefs about  $\lambda$  are well described by a Gamma distribution with shape and scale parameters  $\alpha_0$  and  $\beta_0$ , i.e.,  $\lambda \sim G(\alpha_0, \beta_0)$ , then you are interested in calculating the probability  $P(Y_0 > 0 | \mathbf{y})$ . To answer this, you need to calculate the posterior predictive density  $\pi(y_0 | \mathbf{y})$  in a Bayesian way.

In this example,  $p(\mathbf{y} | \lambda)$  is Poisson, and  $\pi(\lambda)$  is Gamma. Therefore, using Equation 1.9.

$$\pi(y_0 | \mathbf{y}) = \int_0^\infty \frac{\lambda^{y_0} \exp\{-\lambda\}}{y_0!} \times \pi(\lambda | \mathbf{y}) d\lambda,$$

where the posterior distribution is

$$\pi(\lambda | \mathbf{y}) \propto \lambda^{\sum_{i=1}^N y_i + \alpha_0 - 1} \exp\left\{-\lambda \left(\frac{\beta_0 N + 1}{\beta_0}\right)\right\}$$

by equation 1.3.

Observe that the last expression is the kernel of a Gamma distribution with parameters  $\alpha_n = \sum_{i=1}^N y_i + \alpha_0$  and  $\beta_n = \frac{\beta_0}{\beta_0 N + 1}$ . Given that  $\int_0^\infty \pi(\lambda | \mathbf{y}) d\lambda = 1$ , the constant of proportionality in the last expression is  $\Gamma(\alpha_n) \beta_n^{\alpha_n}$ , where  $\Gamma(\cdot)$  is the Gamma function. Thus, the posterior density function  $\pi(\lambda | \mathbf{y})$  is  $G(\alpha_n, \beta_n)$ .

Observe that

$$\begin{aligned} \mathbb{E}[\lambda | \mathbf{y}] &= \alpha_n \beta_n \\ &= \left( \sum_{i=1}^N y_i + \alpha_0 \right) \left( \frac{\beta_0}{\beta_0 N + 1} \right) \\ &= \bar{y} \left( \frac{N \beta_0}{N \beta_0 + 1} \right) + \alpha_0 \beta_0 \left( \frac{1}{N \beta_0 + 1} \right) \\ &= w \bar{y} + (1 - w) \mathbb{E}[\lambda], \end{aligned}$$

where  $\bar{y}$  is the sample mean estimate, which is the maximum likelihood estimate of  $\lambda$  in this example,  $w = \left(\frac{N \beta_0}{N \beta_0 + 1}\right)$ , and  $\mathbb{E}[\lambda] = \alpha_0 \beta_0$  is the prior mean.

---

<sup>10</sup>Independent and identically distributed draws.

The posterior mean is a weighted average of the maximum likelihood estimator (sample information) and the prior mean. Observe that  $\lim_{N \rightarrow \infty} w = 1$ , that is, the sample information asymptotically dominates.

The predictive distribution is

$$\begin{aligned}\pi(y_0 | \mathbf{y}) &= \int_0^\infty \frac{\lambda^{y_0} \exp\{-\lambda\}}{y_0!} \times \frac{1}{\Gamma(\alpha_n)\beta_n^{\alpha_n}} \lambda^{\alpha_n-1} \exp\{-\lambda/\beta_n\} d\lambda \\ &= \frac{1}{y_0! \Gamma(\alpha_n) \beta_n^{\alpha_n}} \int_0^\infty \lambda^{y_0+\alpha_n-1} \exp\left\{-\lambda \left(\frac{1+\beta_n}{\beta_n}\right)\right\} d\lambda \\ &= \frac{\Gamma(y_0 + \alpha_n) \left(\frac{\beta_n}{\beta_n+1}\right)^{y_0+\alpha_n}}{y_0! \Gamma(\alpha_n) \beta_n^{\alpha_n}} \\ &= \binom{y_0 + \alpha_n - 1}{y_0} \left(\frac{\beta_n}{\beta_n+1}\right)^{y_0} \left(\frac{1}{\beta_n+1}\right)^{\alpha_n}.\end{aligned}$$

The third equality follows from the kernel of a Gamma density, and the fourth from

$$\binom{y_0 + \alpha_n - 1}{y_0} = \frac{(y_0 + \alpha_n - 1)(y_0 + \alpha_n - 2) \dots \alpha_n}{y_0!} = \frac{\Gamma(y_0 + \alpha_n)}{\Gamma(\alpha_n) y_0!}$$

using a property of the Gamma function.

Observe that this is a Negative Binomial density, that is,  $Y_0 | \mathbf{y} \sim \text{NB}(\alpha_n, p_n)$  where  $p_n = \frac{\beta_n}{\beta_n+1}$ .

Up to this point, we have said nothing about the hyperparameters, which are required to give a concrete response to this exercise. Thus, we show two approaches to set them. First, we set  $\alpha_0 = 0.001$  and  $\beta_0 = \frac{1}{0.001}$ , which imply vague prior information about  $\lambda$  due to having a large degree of variability compared to the mean information.<sup>11</sup> In particular,  $\mathbb{E}[\lambda] = 1$  and  $\text{Var}[\lambda] = 1000$ .

In this setting,  $P(Y_0 > 0 | \mathbf{y}) = 1 - P(Y_0 = 0 | \mathbf{y}) \approx 0.67$ . That is, the probability of visiting the doctor at least once next year is approximately 0.67.

Another approach is using *Empirical Bayes*, where we set the hyperparameters maximizing the logarithm of the marginal likelihood,<sup>12</sup> that is,

$$\left[ \hat{\alpha}_0 \quad \hat{\beta}_0 \right]^\top = \underset{\alpha_0, \beta_0}{\operatorname{argmax}} \ln p(\mathbf{y})$$

where

$$p(\mathbf{y}) = \int_0^\infty \left\{ \frac{1}{\Gamma(\alpha_0)\beta_0^{\alpha_0}} \lambda^{\alpha_0-1} \exp\{-\lambda/\beta_0\} \prod_{i=1}^N \frac{\lambda^{y_i} \exp\{-\lambda\}}{y_i!} \right\} d\lambda$$

---

<sup>11</sup>We should be aware that there may be technical problems using this kind of hyperparameters in this setting [161].

<sup>12</sup>Empirical Bayes methods are criticized due to double-using the data. First to set the hyperparameters, and second, to perform Bayesian inference.

$$\begin{aligned}
&= \frac{\int_0^\infty \lambda^{\sum_{i=1}^N y_i + \alpha_0 - 1} \exp\left\{-\lambda \left(\frac{\beta_0 N + 1}{\beta_0}\right)\right\} d\lambda}{\Gamma(\alpha_0) \beta_0^{\alpha_0} \prod_{i=1}^N y_i!} \\
&= \frac{\Gamma(\sum_{i=1}^N y_i + \alpha_0) \left(\frac{\beta_0}{N\beta_0 + 1}\right)^{\sum_{i=1}^N y_i} \left(\frac{1}{N\beta_0 + 1}\right)^{\alpha_0}}{\Gamma(\alpha_0) \prod_{i=1}^N y_i}
\end{aligned}$$

Using the empirical Bayes approach, we get  $\hat{\alpha}_0 = 51.8$  and  $\hat{\beta}_0 = 0.023$ , then  $P(Y_0 > 0 | \mathbf{y}) = 1 - P(Y_0 = 0 | \mathbf{y}) \approx 0.70$ .

Observe that we can calculate the posterior odds comparing the model using an Empirical Bayes prior (model 1) versus the vague prior (model 2). We assume that  $\pi(\mathcal{M}_1) = \pi(\mathcal{M}_2) = 0.5$ , then

$$\begin{aligned}
PO_{12} &= \frac{p(\mathbf{y} | \text{Empirical Bayes})}{p(\mathbf{y} | \text{Vague prior})} \\
&= \frac{\frac{\Gamma(\sum_{i=1}^N y_i + 51.807) \left(\frac{0.023}{N \times 0.023 + 1}\right)^{\sum_{i=1}^N y_i} \left(\frac{1}{N \times 0.023 + 1}\right)^{51.807}}{\Gamma(51.807)}}{\frac{\Gamma(\sum_{i=1}^N y_i + 0.001) \left(\frac{1/0.001}{N/0.001 + 1}\right)^{\sum_{i=1}^N y_i} \left(\frac{1}{N/0.001 + 1}\right)^{0.001}}{\Gamma(0.001)}} \\
&\approx 919.
\end{aligned}$$

Then,  $2 \times \log(PO_{12}) = 13.64$ , which provides very strong evidence against the vague prior model (see Table 1.1). In particular,  $\pi(\text{Empirical Bayes} | \mathbf{y}) = \frac{919}{1+919} = 0.999$  and  $\pi(\text{Vague prior} | \mathbf{y}) = 1 - 0.999 = 0.001$ . These probabilities can be used to perform Bayesian model averaging (BMA). In particular,

$$\begin{aligned}
\mathbb{E}(\lambda | \mathbf{y}) &= 1.2 \times 0.999 + 1.2 \times 0.001 = 1.2 \\
Var(\lambda | \mathbf{y}) &= 0.025 \times 0.999 + 0.24 \times 0.001 \\
&\quad + (1.2 - 1.2)^2 \times 0.999 + (1.2 - 1.2)^2 \times 0.001 = 0.025.
\end{aligned}$$

The BMA predictive distribution is a mix of negative binomial distributions, that is,  $Y_0 | \mathbf{y} \sim 0.999 \times NB(57.8, 0.02) + 0.001 \times NB(6.001, 0.17)$ .

The following code shows how to perform this exercise in **R**.

**R code. Health insurance, predictive distribution using vague hyperparameters**

---

```

1 set.seed(010101)
2 y <- c(0, 3, 2, 1, 0) # Data
3 N <- length(y)
4 ProbBo <- function(y, a0, b0){
5   N <- length(y) #sample size
6   an <- a0 + sum(y) # Posterior shape parameter
7   bn <- b0 / ((b0 * N) + 1) # Posterior scale parameter
8   p <- bn / (bn + 1) # Probability negative binomial density
9   Pr <- 1 - pnbinom(0, size=an,prob=(1 - p))
10  # Probability of visiting the Doctor at least once next
11  # year
12  # Observe that in R there is a slightly different
13  # parametrization.
14  return(Pr)
15 }
16 # Using a vague prior:
17 a0 <- 0.001 # Prior shape parameter
18 b0 <- 1 / 0.001 # Prior scale parameter
19 PriMeanV <- a0 * b0 # Prior mean
20 PriVarV <- a0 * b0^2 # Prior variance
21 Pp <- ProbBo(y, a0 = 0.001, b0 = 1 / 0.001)
22 # This setting is vague prior information.
23 # Using Empirical Bayes
24 LogMgLik <- function(theta, y){
25   N <- length(y)
26   #sample size
27   a0 <- theta[1]
28   # prior shape hyperparameter
29   b0 <- theta[2]
30   # prior scale hyperparameter
31   an <- sum(y) + a0
32   # posterior shape parameter
33   if(a0 <= 0 || b0 <= 0){
34     #Avoiding negative values
35     lnp <- -Inf
36   }else{
37     lnp <- lgamma(an) + sum(y)*log(b0/(N*b0+1)) - a0*log(N*
38     b0+1) - lgamma(a0)
39   }
40   # log marginal likelihood
41   return(-lnp)
42 }
43 theta0 <- c(0.01, 1/0.1)
44 # Initial values
45 control <- list(maxit = 1000)
46 # Number of iterations in optimization
47 EmpBay <- optim(theta0, LogMgLik, method = "BFGS", control =
48   control, hessian = TRUE, y = y)
49 a0EB <- EmpBay$par[1]; b0EB <- EmpBay$par[2]
50 PriMeanEB <- a0EB * b0EB # Prior mean
51 PriVarEB <- a0EB * b0EB^2 # Prior variance
52 PpEB <- ProbBo(y, a0 = a0EB, b0 = b0EB)

```

---

### R code. Health insurance, density plots

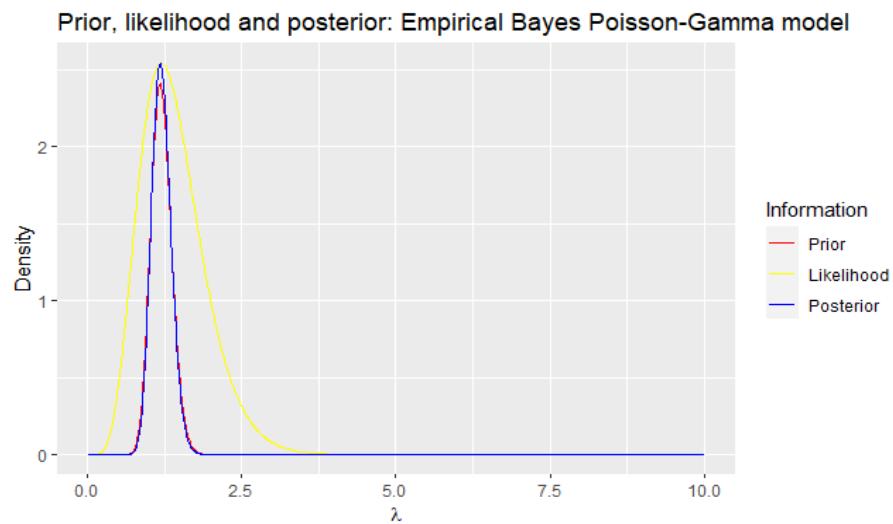
```

1 # Density figures:
2 lambda <- seq(0.01, 10, 0.01) # Values of lambda
3 VaguePrior <- dgamma(lambda, shape=a0, scale = b0)
4 EBPrior <- dgamma(lambda, shape=a0EB, scale = b0EB)
5 PosteriorV <- dgamma(lambda, shape = a0 + sum(y), scale = b0
   / ((b0 * N) + 1))
6 PosteriorEB <- dgamma(lambda, shape = a0EB+sum(y), scale =
   b0EB / ((b0EB * N) + 1))
7 # Likelihood function
8 Likelihood <- function(theta, y){
9   LogL <- dpois(y, theta, log = TRUE); Lik <- prod(exp(LogL))
10  return(Lik)
11 }
12 Liks <- sapply(lambda, function(par) {Likelihood(par, y = y)
   })
13 Sc <- max(PosteriorEB)/max(Liks)
14 LiksScale <- Liks * Sc #Scale for displaying in figure
15 data <- data.frame(cbind(lambda, VaguePrior, EBPrior,
   PosteriorV, PosteriorEB, LiksScale)) #Data frame
16 require(ggplot2); require(latex2exp); require(ggpubr)
17 fig1 <- ggplot(data = data, aes(lambda, VaguePrior)) + geom_
   line() + xlab(TeX("\lambda")) + ylab("Density") +
   ggtitle("Prior: Vague Gamma")
18 fig2 <- ggplot(data = data, aes(lambda, EBPrior)) + geom_
   line() + xlab(TeX("\lambda")) + ylab("Density") +
   ggtitle("Prior: Empirical Bayes Gamma")
19 fig3 <- ggplot(data = data, aes(lambda, PosteriorV)) + geom_
   line() + xlab(TeX("\lambda")) + ylab("Density") +
   ggtitle("Posterior: Vague Gamma")
20 fig4 <- ggplot(data = data, aes(lambda, PosteriorEB)) + geom_
   line() + xlab(TeX("\lambda")) + ylab("Density") +
   ggtitle("Posterior: Empirical Bayes Gamma")
21 FIG <- ggarrange(fig1, fig2, fig3, fig4, ncol = 2, nrow = 2)
22 annotate_figure(FIG, top = text_grob("Vague versus Empirical
   Bayes: Poisson-Gamma model", color = "black", face =
   "bold", size = 14))
23 dataNew <- data.frame(cbind(rep(lambda, 3), c(EBPrior,
   PosteriorEB, LiksScale), rep(1:3, each = 1000))) #Data
   frame
24 colnames(dataNew) <- c("Lambda", "Density", "Factor")
25 dataNew$Factor <- factor(dataNew$Factor, levels=c("1", "3",
   "2"),
26 labels=c("Prior", "Likelihood", "Posterior"))
27 ggplot(data = dataNew, aes_string(x = "Lambda", y = "Density
   ", group = "Factor")) + geom_line(aes(color = Factor)) +
   xlab(TeX("\lambda")) + ylab("Density") + ggtitle("Prior,
   likelihood and posterior: Empirical Bayes Poisson
   -Gamma model") + guides(color=guide_legend(title="Information")) + scale_color_manual(values = c("red", "yellow", "blue"))

```

**FIGURE 1.2**

Vague versus Empirical Bayes: Poisson-Gamma model.

**FIGURE 1.3**

Prior, likelihood and posterior: Empirical Bayes Poisson-Gamma model.

Figure 1.2 displays the prior and posterior densities based on vague and Empirical Bayes hyperparameters. We observe that the prior and posterior densities using the latter are more informative, as expected.

Figure 1.3 shows the prior, scaled likelihood, and posterior densities of  $\lambda$  based on the hyperparameters from the Empirical Bayes approach. The posterior density is a compromise between prior and sample information.

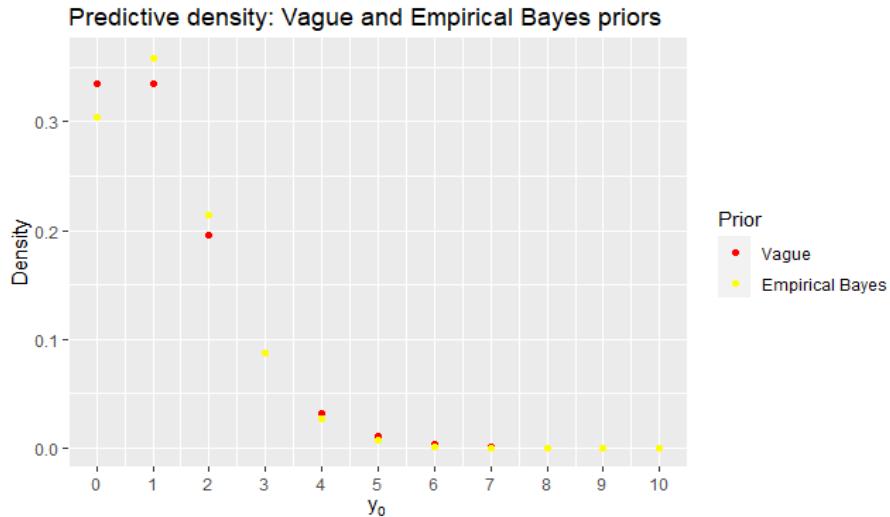
### R code. Health insurance, Predictive density

```

1 # Predictive distributions
2 PredDen <- function(y, y0, a0, b0){
3   N <- length(y)
4   #sample size
5   an <- a0 + sum(y)
6   # Posterior shape parameter
7   bn <- b0 / ((b0 * N) + 1)
8   # Posterior scale parameter
9   p <- bn / (bn + 1)
10  # Probability negative binomial density
11  Pr <- dnb(y0, size=an, prob=(1 - p))
12  # Predictive density
13  # Observe that in R there is a slightly different
     parametrization.
14  return(Pr)
15 }
16 y0 <- 0:10
17 PredVague <- PredDen(y=y, y0=y0, a0=a0, b0=b0)
18 PredEB <- PredDen(y=y, y0=y0, a0=aOEB, b0=bOEB)
19 dataPred <- as.data.frame(cbind(y0, PredVague, PredEB))
20 colnames(dataPred) <- c("y0", "PredictiveVague", "
     PredictiveEB")
21 ggplot(data = dataPred) + geom_point(aes(y0, PredictiveVague
     , color = "red")) +
22 xlab(TeX("$y_0$")) + ylab("Density") + ggtitle("Predictive
     density: Vague and Empirical Bayes priors") + geom_point
     (aes(y0, PredictiveEB, color = "yellow")) +
23 guides(color = guide_legend(title="Prior")) + scale_color_
     manual(labels = c("Vague", "Empirical Bayes"), values =
     c("red", "yellow")) + scale_x_continuous(breaks=seq
     (0,10,by=1))

```

Figure 1.4 displays the predictive probability mass of not having any visits to a physician next year, as well as having one, two, and so on, using Empirical Bayes and vague hyperparameters. The predictive probabilities of not having any visits are approximately 30% and 33% based on the Empirical Bayes and vague hyperparameters, respectively.



**FIGURE 1.4**  
Predictive density: Vague and Empirical Bayes.

### R code. Health insurance, Bayesian model average

```

1 # Posterior odds: Vague vs Empirical Bayes
2 P012 <- exp(-LogMgLik(c(a0EB, b0EB), y = y))/exp(-LogMgLik(c
  (a0, b0), y = y))
3 PostProMEM <- P012/(1 + P012)
4 # Posterior model probability Empirical Bayes
5 PostProbMV <- 1 - PostProMEM
6 # Posterior model probability vague prior
7 # Bayesian model average (BMA)
8 PostMeanEB <- (a0EB + sum(y)) * (b0EB / (b0EB * N + 1))
9 # Posterior mean Empirical Bayes
10 PostMeanV <- (a0 + sum(y)) * (b0 / (b0 * N + 1))
11 # Posterior mean vague priors
12 BMAMean <- PostProMEM * PostMeanEB + PostProbMV * PostMeanV
13 # BMA posterior mean
14 PostVarEB <- (a0EB + sum(y)) * (b0EB/(b0EB * N + 1))^2
15 # Posterior variance Empirical Bayes
16 PostVarV <- (a0 + sum(y)) * (b0 / (b0 * N + 1))^2
17 # Posterior variance vague prior
18 BMAMVar <- PostProMEM * PostVarEB + PostProbMV*PostVarV +
  PostProMEM * (PostMeanEB - BMAMean)^2 + PostProbMV * (
    PostMeanV - BMAMean)^2
19 # BMA posterior variance

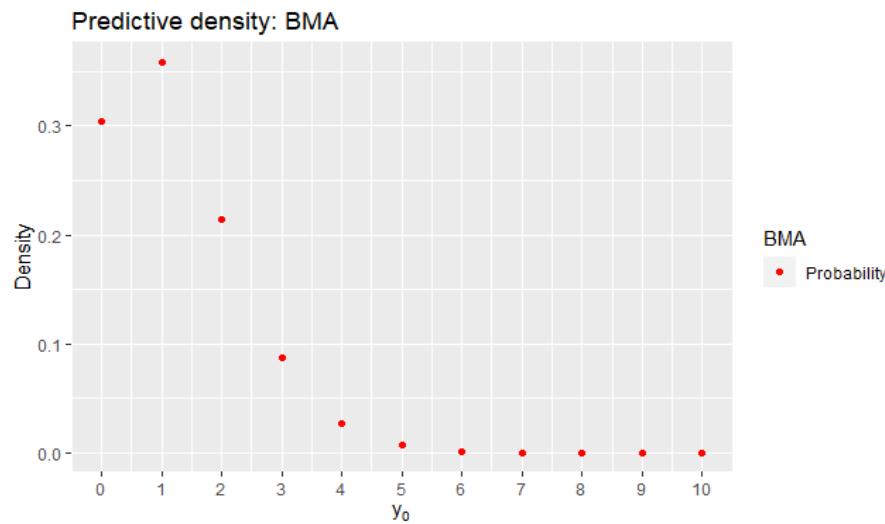
```

**R code. Health insurance, Bayesian model average**

```

1 # BMA: Predictive
2 BMAPred <- PostProMEM * PredEB+PostProbMV * PredVague
3 dataPredBMA <- as.data.frame(cbind(y0, BMAPred))
4 colnames(dataPredBMA) <- c("y0", "PredictiveBMA")
5 ggplot(data = dataPredBMA) + geom_point(aes(y0,
  PredictiveBMA, color = "red")) + xlab(TeX("\$y_0\$")) +
  ylab("Density") + ggtitle("Predictive density: BMA") +
  guides(color = guide_legend(title="BMA")) + scale_color_
  manual(labels = c("Probability"), values = c("red")) +
  scale_x_continuous(breaks=seq(0,10,by=1))
6

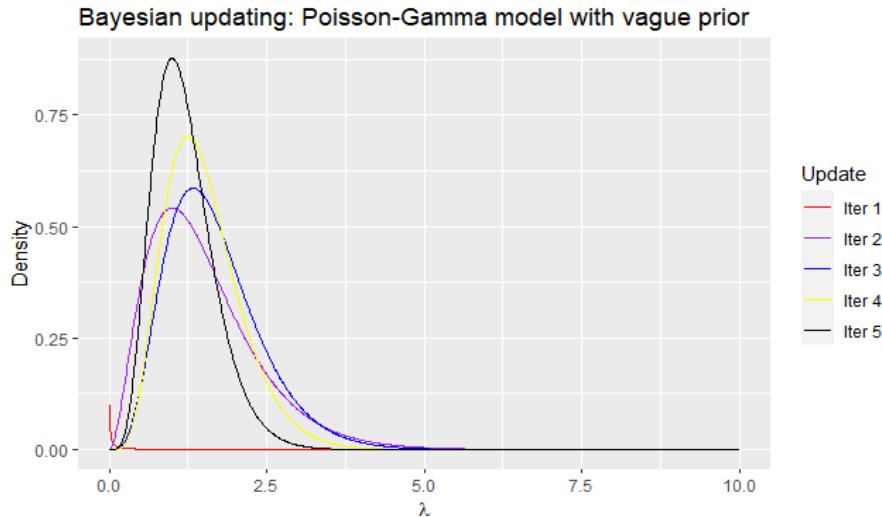
```



**FIGURE 1.5**  
Bayesian model average: Predictive density.

Figure 1.5 displays the predictive density using Bayesian model averaging based on the vague and Empirical Bayes hyperparameters. This figure closely resembles the predictive probability mass function based on the Empirical Bayes framework, as the posterior model probability for that setting is nearly one.

Figure 1.6 shows how the posterior distribution updates with new sample information, starting from an initial non-informative prior (iteration 1). We



**FIGURE 1.6**  
Bayesian updating: Posterior densities.

observe that iteration 5 incorporates all the sample information in our example. As a result, the posterior density in iteration 5 is identical to the posterior density shown in Figure 1.3.

***R code. Health insurance, Bayes updating***

```

1 # Bayesian updating
2 BayUp <- function(y, lambda, a0, b0){
3   N <- length(y)
4   #sample size
5   an <- a0 + sum(y)
6   # Posterior shape parameter
7   bn <- b0 / ((b0 * N) + 1)
8   # Posterior scale parameter
9   p <- dgamma(lambda, shape = an, scale = bn)
10  # Posterior density
11  return(list(Post = p, a0New = an, b0New = bn))
12 }
```

### R code. Health insurance, Bayes updating

```

1 PostUp <- NULL
2 for(i in 1:N){
3   if(i == 1){
4     PostUpi <- BayUp(y[i], lambda, a0 = 0.001, b0 = 1/0.001)
5   }
6   else{
7     PostUpi <- BayUp(y[i], lambda, a0 = PostUpi$a0New, b0 =
8       PostUpi$b0New)
9   }
10 PostUp <- cbind(PostUp, PostUpi$Post)
11 DataUp <- data.frame(cbind(rep(lambda, 5), c(PostUp), rep
12   (1:5, each = 1000))) #Data frame
13 colnames(DataUp) <- c("Lambda", "Density", "Factor")
14 DataUp$Factor <- factor(DataUp$Factor, levels=c("1", "2", "3",
15   "4", "5"),
16   labels=c("Iter 1", "Iter 2", "Iter 3", "Iter 4", "Iter 5"))
17 ggplot(data = DataUp, aes_string(x = "Lambda", y = "Density"
18   , group = "Factor")) + geom_line(aes(color = Factor)) +
19   xlab(TeX("\lambda")) + ylab("Density") + ggtitle(""
20   Bayesian updating: Poisson-Gamma model with vague prior"
21   ) + guides(color=guide_legend(title="Update")) + scale_
22   color_manual(values = c("red", "purple", "blue", "yellow
23   ", "black"))
24 S <- 100000 # Posterior draws
25 PostMeanLambdaUps <- sapply(1:N, function(i) {mean(sample(
26   lambda, S, replace = TRUE, prob = PostUp[, i]))}) #
27   Posterior mean update i
28 paste("Posterior means using all information and sequential
29   updating are:", round(PostMeanV, 2), "and", round(
30   PostMeanLambdaUps[5], 2), sep = " ")
31 Posterior means using all information and sequential
32   updating are: 1.2 and 1.2
33 PostVarLambdaUps <- sapply(1:N, function(i) {var(sample(
34   lambda, S, replace = TRUE, prob = PostUp[, i]))}) #
35   Posterior variance update i
36 paste("Posterior variances using all information and
37   sequential updating are:", round(PostVarV, 2), "and",
38   round(PostVarLambdaUps[5], 2), sep = " ")
39 Posterior variances using all information and sequential
40   updating are: 0.24 and 0.24

```

### 1.3 Bayesian reports: Decision theory under uncertainty

The Bayesian framework allows reporting the full posterior distributions. However, some situations require reporting a specific value of the posterior distribution (point estimate), an informative interval (set), point or interval predictions, and/or selecting a specific model. Decision theory offers an elegant framework to make decisions regarding the optimal posterior values to report [39].

The starting point is a *loss function*, which is a non-negative real-valued function whose arguments are the unknown *state of nature* ( $\Theta$ ), and a set of *actions* to be taken ( $\mathcal{A}$ ), that is,

$$L(\boldsymbol{\theta}, a) : \Theta \times \mathcal{A} \rightarrow \mathcal{R}^+.$$

This function is a mathematical representation of the loss incurred from making mistakes. In particular, selecting action  $a \in \mathcal{A}$  when  $\boldsymbol{\theta} \in \Theta$  is the true state. In our case, the unknown state of nature can refer to parameters, functions of them, future or unknown realizations, models, etc.

From a Bayesian perspective, we should choose the action that minimizes the posterior expected loss ( $a^*(\mathbf{y})$ ), that is, the *posterior risk function* ( $\mathbb{E}[L(\boldsymbol{\theta}, a) | \mathbf{y}]$ ),

$$a^*(\mathbf{y}) = \operatorname{argmin}_{a \in \mathcal{A}} \mathbb{E}[L(\boldsymbol{\theta}, a) | \mathbf{y}],$$

where  $\mathbb{E}[L(\boldsymbol{\theta}, a) | \mathbf{y}] = \int_{\Theta} L(\boldsymbol{\theta}, a) \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}$ .<sup>13</sup>

Different loss functions imply different optimal decisions. We illustrate this assuming  $\boldsymbol{\theta} \in \mathcal{R}$ .

- The quadratic loss function,  $L(\boldsymbol{\theta}, a) = [\boldsymbol{\theta} - a]^2$ , gives as the optimal decision the posterior mean,  $a^*(\mathbf{y}) = \mathbb{E}[\boldsymbol{\theta} | \mathbf{y}]$ , that is:

$$\mathbb{E}[\boldsymbol{\theta} | \mathbf{y}] = \operatorname{argmin}_{a \in \mathcal{A}} \int_{\Theta} [\boldsymbol{\theta} - a]^2 \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}.$$

To obtain this result, let's use the first-order condition, differentiate the risk function with respect to  $a$ , interchange the differential and integral order, and set the result equal to zero:

$$-2 \int_{\Theta} [\boldsymbol{\theta} - a^*] \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} = 0.$$

---

<sup>13</sup>[80] propose Laplace-type estimators (LTE) based on the *quasi-posterior*,  $p(\boldsymbol{\theta}) = \frac{\exp\{L_n(\boldsymbol{\theta})\}\pi(\boldsymbol{\theta})}{\int_{\Theta} \exp\{L_n(\boldsymbol{\theta})\}\pi(\boldsymbol{\theta})d\boldsymbol{\theta}}$ , where  $L_n(\boldsymbol{\theta})$  is not necessarily a log-likelihood function. The LTE minimizes the *quasi-posterior risk*.

This implies that

$$a^* \int_{\Theta} \pi(\theta | \mathbf{y}) d\theta = a^*(\mathbf{y}) = \int_{\Theta} \theta \pi(\theta | \mathbf{y}) d\theta = \mathbb{E}[\theta | \mathbf{y}],$$

that is, the posterior mean is the Bayesian optimal action. This means that we should report the posterior mean as a point estimate of  $\theta$  when facing the quadratic loss function.

- The generalized quadratic loss function,  $L(\theta, a) = w(\theta)[\theta - a]^2$ , where  $w(\theta) > 0$  is a weighting function, gives as the optimal decision rule the weighted mean. We should follow the same steps as the previous result to obtain

$$a^*(\mathbf{y}) = \frac{\mathbb{E}[w(\theta) \times \theta | \mathbf{y}]}{\mathbb{E}[w(\theta) | \mathbf{y}]}.$$

Observe that the weighted average is driven by the weighting function  $w(\theta)$ .

- The absolute error loss function,  $L(\theta, a) = |\theta - a|$ , gives as the optimal action the posterior median (Exercise 5).
- The generalized absolute error function,

$$L(\theta, a) = \begin{cases} K_0(\theta - a), & \theta - a \geq 0, \\ K_1(a - \theta), & \theta - a < 0, \end{cases} \quad K_0, K_1 > 0,$$

implies the following risk function:

$$\mathbb{E}[L(\theta, a) | \mathbf{y}] = \int_{-\infty}^a K_1(a - \theta) \pi(\theta | \mathbf{y}) d\theta + \int_a^{\infty} K_0(\theta - a) \pi(\theta | \mathbf{y}) d\theta.$$

Differentiating with respect to  $a$ , interchanging differentials and integrals, and equating to zero, we get:

$$K_1 \int_{-\infty}^{a^*} \pi(\theta | \mathbf{y}) d\theta - K_0 \int_{a^*}^{\infty} \pi(\theta | \mathbf{y}) d\theta = 0.$$

Thus, we have

$$\int_{-\infty}^{a^*} \pi(\theta | \mathbf{y}) d\theta = \frac{K_0}{K_0 + K_1},$$

that is, any  $\frac{K_0}{K_0 + K_1}$ -percentile of  $\pi(\theta | \mathbf{y})$  is an optimal Bayesian estimate of  $\theta$ .

We can also use decision theory under uncertainty in hypothesis testing. In particular, testing  $H_0 : \theta \in \Theta_0$  versus  $H_1 : \theta \in \Theta_1$ , where  $\Theta = \Theta_0 \cup \Theta_1$  and  $\emptyset = \Theta_0 \cap \Theta_1$ , there are two actions of interest,  $a_0$  and  $a_1$ , where  $a_j$  denotes not rejecting  $H_j$ , for  $j = \{0, 1\}$ .

Given the  $0 - K_j$  loss function:

$$L(\theta, a_j) = \begin{cases} 0, & \text{if } \theta \in \Theta_j, \\ K_j, & \text{if } \theta \in \Theta_i, j \neq i, \end{cases}$$

where there is no loss if the right decision is made, for instance, not rejecting  $H_0$  when  $\theta \in \Theta_0$ , and the loss is  $K_j$  when an error is made. For example, a type I error occurs when rejecting the null hypothesis ( $H_0$ ) when it is true ( $\theta \in \Theta_0$ ), which results in a loss of  $K_1$  due to choosing action  $a_1$ , not rejecting  $H_1$ .

The posterior expected loss associated with decision  $a_j$ , i.e., not rejecting  $H_j$ , is:

$$\mathbb{E}[L(\theta, a_j) | \mathbf{y}] = 0 \times P(\Theta_j | \mathbf{y}) + K_j P(\Theta_i | \mathbf{y}) = K_j P(\Theta_i | \mathbf{y}), \quad j \neq i.$$

Therefore, the Bayes optimal decision is the one that minimizes the posterior expected loss. That is, the null hypothesis is rejected ( $a_1$  is not rejected) when

$$K_0 P(\Theta_1 | \mathbf{y}) > K_1 P(\Theta_0 | \mathbf{y}).$$

Given our framework,  $\Theta = \Theta_0 \cup \Theta_1$  and  $\emptyset = \Theta_0 \cap \Theta_1$ , we have  $P(\Theta_0 | \mathbf{y}) = 1 - P(\Theta_1 | \mathbf{y})$ . As a result, the rejection region of the Bayesian test is:

$$R = \left\{ \mathbf{y} : P(\Theta_1 | \mathbf{y}) > \frac{K_1}{K_1 + K_0} \right\}.$$

Decision theory also helps to construct interval (region) estimates. Let  $\Theta_{C(\mathbf{y})} \subset \Theta$  be a *credible set* for  $\theta$ , and let the loss function be defined as:

$$L(\theta, \Theta_{C(\mathbf{y})}) = 1 - \mathbb{1}\{\theta \in \Theta_{C(\mathbf{y})}\},$$

where

$$\mathbb{1}\{\theta \in \Theta_{C(\mathbf{y})}\} = \begin{cases} 1, & \text{if } \theta \in \Theta_{C(\mathbf{y})}, \\ 0, & \text{if } \theta \notin \Theta_{C(\mathbf{y})}. \end{cases}$$

Thus, the loss function becomes:

$$L(\theta, \Theta_{C(\mathbf{y})}) = \begin{cases} 0, & \text{if } \theta \in \Theta_{C(\mathbf{y})}, \\ 1, & \text{if } \theta \notin \Theta_{C(\mathbf{y})}. \end{cases}$$

This is a 0–1 loss function, which equals zero when  $\theta \in \Theta_{C(\mathbf{y})}$  and equals one when  $\theta \notin \Theta_{C(\mathbf{y})}$ . Consequently, the risk function is:

$$1 - P(\theta \in \Theta_{C(\mathbf{y})}).$$

Given a *measure of credibility*  $\alpha(\mathbf{y})$  that defines the level of trust that  $\theta \in \Theta_{C(\mathbf{y})}$ , we can measure the accuracy of the report by the loss function:

$$L(\theta, \alpha(\mathbf{y})) = [\mathbb{1}\{\theta \in \Theta_{C(\mathbf{y})}\} - \alpha(\mathbf{y})]^2.$$

This loss function could be used to suggest a choice of the report  $\alpha(\mathbf{y})$ . Given that this is a quadratic loss function, the optimal action is the posterior mean, that is,

$$\mathbb{E}[\mathbb{1}\{\theta \in \Theta_{C(\mathbf{y})}\} | \mathbf{y}] = P(\theta \in \Theta_{C(\mathbf{y})} | \mathbf{y}).$$

This probability can be calculated given the posterior distribution as

$$P(\theta \in \Theta_{C(\mathbf{y})} | \mathbf{y}) = \int_{\Theta_{C(\mathbf{y})}} \pi(\theta | \mathbf{y}) d\theta.$$

This represents a measure of the belief that  $\theta \in \Theta_{C(\mathbf{y})}$  given the prior beliefs and sample information.

The set  $\Theta_{C(\mathbf{y})} \subset \Theta$  is a  $100(1 - \alpha)\%$  credible set with respect to  $\pi(\theta | \mathbf{y})$  if

$$P(\theta \in \Theta_{C(\mathbf{y})} | \mathbf{y}) = \int_{\Theta_{C(\mathbf{y})}} \pi(\theta | \mathbf{y}) d\theta = 1 - \alpha.$$

Two alternatives for reporting credible sets are the *symmetric credible set* and the *highest posterior density set* (HPD). The former is based on the  $\frac{\alpha}{2}\%$  and  $(1 - \frac{\alpha}{2})\%$  percentiles of the posterior distribution, and the latter is a  $100(1 - \alpha)\%$  credible interval for  $\theta$  with the property that it has the smallest distance compared to any other  $100(1 - \alpha)\%$  credible interval for  $\theta$  based on the posterior distribution. Specifically,

$$C(\mathbf{y}) = \{\theta : \pi(\theta | \mathbf{y}) \geq k(\alpha)\},$$

where  $k(\alpha)$  is the largest number such that

$$\int_{\theta: \pi(\theta | \mathbf{y}) \geq k(\alpha)} \pi(\theta | \mathbf{y}) d\theta = 1 - \alpha.$$

The HPD set can be a collection of disjoint intervals when working with multimodal posterior densities. Additionally, HPD sets have the limitation of not necessarily being invariant under transformations.

Decision theory can also be used to perform prediction (point, sets, or probabilistic). Suppose that there is a loss function  $L(Y_0, a)$  involving the prediction of  $Y_0$ . Then, the expected loss is

$$\mathbb{E}_{Y_0}[L(Y_0, a)] = \int_{\mathcal{Y}_0} L(y_0, a) \pi(y_0 | \mathbf{y}) dy_0,$$

where  $\pi(y_0 | \mathbf{y})$  is the predictive density function. Thus, we make an

optimal choice for prediction that minimizes the risk function given a specific loss function.

Although Bayesian Model Averaging (BMA) allows for incorporating model uncertainty in a regression framework, sometimes it is desirable to select just one model. A compelling alternative is to choose the model with the highest posterior model probability. This model is the best alternative for prediction in the case of a 0-1 loss function [103].

### 1.3.1 Example: Health insurance continues

We show some optimal rules in the health insurance example, specifically the best point estimates of  $\lambda$  under the quadratic, absolute, and generalized absolute loss functions. For the generalized absolute loss function, we assume that underestimating  $\lambda$  is twice as costly as overestimating it, i.e.,  $K_0 = 2$  and  $K_1 = 1$ .

Given that the posterior distribution of  $\lambda$  is  $G(\alpha_0 + \sum_{i=1}^N y_i, \frac{\beta_0}{\beta_0 N + 1})$ , and using the hyperparameters from empirical Bayes, we obtain the following optimal point estimates:

- The posterior mean:  $\mathbb{E}[\lambda | \mathbf{y}] = \alpha_n \beta_n = 1.2$ ,
- The posterior median: 1.19,
- The 2/3-th quantile: 1.26.

These are the optimal point estimates for the quadratic, absolute, and generalized absolute loss functions, respectively.

In addition, we test the null hypothesis  $H_0 : \lambda \in [0, 1)$  versus the alternative hypothesis  $H_1 : \lambda \in [1, \infty)$ , setting  $K_0 = K_1 = 1$ . We should reject the null hypothesis since  $P(\lambda \in [1, \infty)) = 0.9 > \frac{K_1}{K_0 + K_1} = 0.5$ .

The 95% symmetric credible interval is (0.91, 1.53), and the highest posterior density (HPD) interval is (0.90, 1.51). Finally, the optimal point prediction under the quadratic loss function is 1.2, which is the mean value of the posterior predictive distribution. The optimal model, assuming a 0-1 loss function, is the model using the hyperparameters from the empirical Bayes procedure, since the posterior model probability of this model is approximately 1, whereas the posterior model probability of the model using vague hyperparameters is approximately 0.

**R code. Health insurance, Bayesian reports**

```

1 an <- sum(y) + a0EB # Posterior shape parameter
2 bn <- b0EB / (N*b0EB + 1) # Posterior scale parameter
3 S <- 1000000
4 # Number of posterior draws
5 Draws <- rgamma(1000000, shape = an, scale = bn)
6 ##### Point estimation #####
7 OptQua <- an*bn # Mean: Optimal choice quadratic loss
    function
8 OptAbs <- qgamma(0.5, shape = an, scale = bn) # Median:
    Optimal choice absolute loss function
9 # Setting K0 = 2 and K1 = 1, that is, to underestimate
    lambda is twice as costly as to overestimate it.
10 K0 <- 2; K1 <- 1
11 OptGenAbs <- quantile(Draws, K0/(K0 + K1)) # Median: Optimal
    choice generalized absolute loss function
12 ##### Hypothesis test #####
13 # H0: lambda in [0,1) vs H1: lambda in [1, Inf]
14 K0 <- 1; K1 <- 1
15 ProbH0 <- pgamma(1, shape = an, scale = bn) # Posterior
    probability H0
16 ProbH1 <- 1 -ProbH0 # Posterior probability H1
17 # We should reject H0 given ProbH1 > K1 / (K0 + K1)
18 ##### Credible intervals #####
19 LimInf <- qgamma(0.025, shape = an, scale = bn) # Lower
    bound
20 LimSup <- qgamma(0.975, shape = an, scale = bn) # Upper
    bound
21 HDI <- HDInterval::hdi(Draws, credMass = 0.95) # Highest
    posterior density credible interval
22 ##### Predictive optimal choices #####
23 p <- bn / (bn + 1) # Probability negative binomial density
24 OptPred <- p/(1-p)*an # Optimal point prediction given a
    quadratic loss function in prediction

```

## 1.4 Summary

We introduce Bayes' rule to update probabilistic statements using humorous examples. We then study the three key probabilistic objects in Bayesian inference: the posterior distribution, the marginal likelihood, and the predictive density. The posterior distribution allows for inference regarding parameters, the marginal likelihood is required for hypothesis testing and model selection

using the Bayes factor, and the predictive density enables probabilistic predictions. We also review some sampling properties of Bayesian estimators and the process of Bayes updating. All of these concepts were illustrated using a simple example in **R** software. Finally, we introduce decision theory concepts that can be applied to report summary statistics while minimizing posterior expected losses.

## 1.5 Exercises

### 1. *The Court Case: The Blue or Green Cab*

A cab was involved in a hit-and-run accident at night. There are two cab companies in the town: Blue and Green. The Blue company has 150 cabs, while the Green company has 850 cabs. A witness stated that a blue cab was involved in the accident. The court tested the reliability of the witness under similar circumstances and found that the witness correctly identified the color of the cab 80% of the time, but made an incorrect identification 25% of the time. *What is the probability that the cab involved in the accident was actually blue, given that the witness said it was blue?*

### 2. *The Monty Hall Problem*

What is the probability of winning a car in the *Monty Hall problem* if you switch your decision, when there are four doors, three goats, and one car? Solve this problem both analytically and computationally. What if there are  $n$  doors,  $n - 1$  goats, and one car?

3. Solve the health insurance example using a Gamma prior in the rate parametrization, that is,  $\pi(\lambda) = \frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} \lambda^{\alpha_0-1} \exp\{-\lambda\beta_0\}$ .
4. Suppose you are analyzing the decision to buy car insurance for the next year. To make a better decision, you want to know: *What is the probability that you will have a car claim next year?* You have the records of your car claims over the last 15 years,  $\mathbf{y} = \{0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0\}$ .

Assume that this is a random sample from a data-generating process (statistical model) that is Bernoulli,  $Y_i \sim \text{Ber}(p)$ . Your prior beliefs about  $p$  are well described by a Beta distribution with parameters  $\alpha_0$  and  $\beta_0$ , i.e.,  $p \sim B(\alpha_0, \beta_0)$ . You are interested in calculating the probability of a claim the next year,  $P(Y_0 = 1 | \mathbf{y})$ .

Solve this using both an empirical Bayes approach and a non-informative approach where  $\alpha_0 = \beta_0 = 1$  (uniform distribution).

5. Show that, given the loss function  $L(\theta, a) = |\theta - a|$ , the optimal decision rule minimizing the risk function,  $a^*(\mathbf{y})$ , is the median.

## 2

---

### *Conceptual differences between the Bayesian and Frequentist approaches*

---

We outline some of the conceptual differences between the Bayesian and Frequentist inferential approaches. We emphasize Bayesian concepts, as most readers may already be familiar with the Frequentist statistical framework. We illustrate the differences between these two inferential approaches using a simple example. In addition, we provide some potential explanations for why the Bayesian inferential framework is not well known at the introductory level among practitioners and applied researchers.

---

#### **2.1 The concept of probability**

Let's begin with the following thought experiment: Assume that you are watching the international game show "Who Wants to Be a Millionaire?". The contestant is asked to answer a very simple question: **What is the last name of the brothers who are credited with inventing the world's first successful motor-operated airplane?**

- What is the probability that the contestant answers this question correctly?

Unless you have:

1. watched this particular contestant participate in this show many times,
2. seen her asked this same question each time,
3. and computed the relative frequency with which she gives the correct answer,

you need to answer this question as a Bayesian!

Uncertainty about the event *answering this question* needs to be expressed as a "degree of belief," informed by both data on the skill of the particular participant and how much she knows about inventors, as well as possibly prior knowledge of her performance in other game shows. Of course, your prior knowledge of the contestant may be minimal, or it may be well-informed.

Either way, your final answer remains a degree of belief about an uncertain, and inherently unrepeatable, state of nature.

The point of this hypothetical, light-hearted scenario is simply to highlight that a key distinction between the Frequentist and Bayesian approaches to inference is not the use (or nature) of prior information, but the manner in which probability is used. To the Bayesian, probability is the mathematical construct used to quantify uncertainty about an unknown state of nature, conditional on observed data and prior knowledge about the context in which that state occurs. To the Frequentist, probability is intrinsically linked to the concept of a repeated experiment, and the relative frequency with which a particular outcome occurs, conditional on that unknown state. This distinction remains key whether the Bayesian chooses to be *informative or subjective* in the specification of prior information, or chooses to be *non-informative or objective*.

Frequentists consider probability to be a physical phenomenon, like mass or wavelength, whereas Bayesians stipulate that probability exists in the mind of scientists, as any scientific construct [331].

It seems that the understanding of the concept of probability for the common human being is more associated with “degrees of belief” rather than relative frequency. Peter Diggle, President of The Royal Statistical Society between 2014 and 2016, was asked in an interview, “A different trend which has surged upwards in statistics during Peter’s career is the popularity of Bayesian statistics. Does Peter consider himself a Bayesian?” He replied, “... you can’t not believe in Bayes’ theorem because it’s true. But that doesn’t make you a Bayesian in the philosophical sense. When people are making personal decisions – even if they don’t formally process Bayes’ theorem in their mind – they are adapting what they think they should believe in response to new evidence as it comes in. Bayes’ theorem is just the formal mathematical machinery for doing that.”

However, we should mention that psychological experiments suggest that human beings suffer from *anchoring*, a cognitive bias that causes us to rely too heavily on previous information (the prior), so that the updating process (posterior) due to new information (likelihood) is not as strong as Bayes’ rule would suggest [235].

---

## 2.2 Subjectivity is not the key

The concepts of *subjectivity* and *objectivity* indeed characterize both statistical paradigms in differing ways. Among Bayesians, there are those who are immersed in *subjective* rationality [362, 112, 396, 269], but others who adopt *objective* prior distributions such as Jeffreys’, reference, empirical, or robust priors [32, 256, 223, 38] to operationalize Bayes’ rule and thereby weight quan-

titative (data-based) evidence. Among Frequentists, there are choices made about significance levels which, if not explicitly subjective, are typically not grounded in any objective and documented assessment of the relative losses of Type I and Type II errors.<sup>1</sup> In addition, both Frequentist and Bayesian Econometrician/Statisticians make decisions about the form of the data generating process, or “model”, which – if not subject to rigorous diagnostic assessment – retains a subjective element that potentially influences the final inferential outcome. Although we all know that by definition, a model is a schematic and simplified approximation to reality,

“Since all models are wrong, the scientist cannot obtain a *correct* one by excessive elaboration. On the contrary, following William of Occam, he should seek an economical description of natural phenomena.” [54].

We also know that “All models are wrong, but some are useful” [57], which is why model diagnostics are important. This task can be performed in both approaches. Particularly, the Bayesian framework can use predictive *p*-values for absolute testing [158, 30] or posterior odds ratios for relative statements [222, 241]. This is because the marginal likelihood, conditional on data, is interpreted as the evidence for the prior distribution [37].

In addition, what does objectivity mean in a Frequentist approach? For example, why should we use a 5% or 1% significance level rather than any other value? As someone said, the apparent objectivity is really a consensus [269]. In fact, “Student” (William Gosset) saw statistical significance at any level as being “nearly valueless” in itself [470]. But, this is not just a situation in the Frequentist approach. The cut-offs used to “establish” scientific evidence against a null hypothesis, in terms of  $\log_{10}$  scale [223] or  $\log_e$  scale [241] as shown in Table 1.1, are also *ad hoc*.

Although the true state of nature in Bayesian inference is expressed in “degrees of belief”, the distinction between the two paradigms does not reside in one being more, or less, *subjective* than the other. Rather, the differences are philosophical, pedagogical, and methodological.

---

### 2.3 Estimation, hypothesis testing and prediction

All that is required to perform estimation, hypothesis testing (model selection), and prediction in the Bayesian approach is to apply Bayes’ rule. This ensures coherence under a probabilistic view. However, there is no free lunch: coherence reduces flexibility. On the other hand, the Frequentist approach may not be coherent from a probabilistic point of view, but it is highly flexible. This approach can be seen as a toolkit that offers inferential solutions under

---

<sup>1</sup>Type I error is rejecting the null hypothesis when it is true, and Type II error is not rejecting the null hypothesis when it is false.

the umbrella of understanding probability as relative frequency. For instance, a point estimator in a Frequentist approach is found such that it satisfies good sampling properties like unbiasedness, efficiency, or a large sample property such as consistency.

A notable difference is that optimal Bayesian decisions are calculated by minimizing the expected value of the loss function with respect to the posterior distribution, i.e., conditional on observed data. In contrast, Frequentist “optimal” actions are based on the expected values over the distribution of the estimator (a function of data), conditional on the unknown parameters. This involves considering sampling variability.

The Bayesian approach allows for the derivation of the posterior distribution of any unknown object, such as parameters, latent variables, future or unobserved variables, or models. A major advantage is that predictions can account for estimation error, and predictive distributions (probabilistic forecasts) can be easily derived.

Hypothesis testing (model selection) in the Bayesian framework is based on *inductive logic* reasoning (*inverse probability*). Based on observed data, we evaluate which hypothesis is most tenable, performing this evaluation using posterior odds. These odds are in turn based on Bayes factors, which assess the evidence in favor of a null hypothesis while explicitly considering the alternative [241], following the rules of probability [269]. This approach compares how well hypotheses predict data [180], minimizes the weighted sum of type I and type II error probabilities [115, 334], and takes into account the implicit balance of losses [223, 41]. Posterior odds allow for the use of the same framework to analyze nested and non-nested models and perform model averaging.

However, Bayes factors cannot be based on improper or vague priors [250], the practical interplay between model selection and posterior distributions is not as straightforward as it may be in the Frequentist approach, and the computational burden can be more demanding due to the need to solve potentially difficult integrals.

On the other hand, the Frequentist approach establishes most of its estimators as the solution to a system of equations. Observe that optimization problems often reduce to solving systems. We can potentially obtain the distribution of these estimators, but most of the time, asymptotic arguments or resampling techniques are required. Hypothesis testing relies on pivotal quantities and/or resampling, and prediction is typically based on a *plug-in approach*, which means that estimation error is not taken into account.<sup>2</sup>

Comparing models depends on their structure. For instance, there are different Frequentist statistical approaches to compare nested and non-nested models. A nice feature in some situations is that there is a practical interplay between hypothesis testing and confidence intervals. For example, in the normal population mean hypothesis framework, you cannot reject a null hy-

---

<sup>2</sup>A pivot quantity is a function of unobserved parameters and observations whose probability distribution does not depend on the unknown parameters.

pothesis  $H_0 : \mu = \mu^0$  at the  $\alpha$  significance level (Type I error) if  $\mu^0$  is in the  $1 - \alpha$  confidence interval. Specifically,

$$P\left(\mu \in \left[\hat{\mu} - |t_{N-1}^{\alpha/2}| \times \hat{\sigma}_{\hat{\mu}}, \hat{\mu} + |t_{N-1}^{\alpha/2}| \times \hat{\sigma}_{\hat{\mu}}\right]\right) = 1 - \alpha,$$

where  $\hat{\mu}$  and  $\hat{\sigma}_{\hat{\mu}}$  are the maximum likelihood estimators of the mean and standard error,  $t_{N-1}^{\alpha/2}$  is the quantile value of the Student's  $t$ -distribution at the  $\alpha/2$  probability level with  $N - 1$  degrees of freedom, and  $N$  is the sample size.

A remarkable difference between the Bayesian and Frequentist inferential frameworks is the interpretation of credible/confidence intervals. Observe that once we have estimates, such that for example the previous interval is  $[0.2, 0.4]$  given a 95% confidence level, we cannot say that  $P(\mu \in [0.2, 0.4]) = 0.95$  in the Frequentist framework. In fact, this probability is either 0 or 1 in this approach, as  $\mu$  is either in the interval or it is not. The problem is that we will never know for certain in applied settings. This is because  $P(\mu \in [\hat{\mu} - |t_{N-1}^{0.025}| \times \hat{\sigma}_{\hat{\mu}}, \hat{\mu} + |t_{N-1}^{0.025}| \times \hat{\sigma}_{\hat{\mu}}]) = 0.95$  is interpreted in the context of repeated sampling. On the other hand, once we have the posterior distribution in the Bayesian framework, we can say that  $P(\mu \in [0.2, 0.4]) = 0.95$ .

Following common practice, most researchers and practitioners conduct hypothesis testing based on the  $p$ -value in the Frequentist framework. But **what is a  $p$ -value?** Most users do not know the answer, as statistical inference is often not performed by statisticians [38].<sup>3</sup> A  $p$ -value is the probability of obtaining a statistical summary of the data equal to or *more extreme* than what was actually observed, assuming that the null hypothesis is true.

Therefore,  $p$ -value calculations involve not just the observed data, but also more *extreme* hypothetical observations. Thus,

“What the use of  $p$  implies, therefore, is that a hypothesis that may be true may be rejected because it has not predicted observable results that have not occurred.” [223]

Some researchers and practitioners using Frequentist inference often intertwines two distinct logical frameworks: Fisher's  $p$ -value approach [140] and the Neyman–Pearson significance testing framework [322]. The  $p$ -value serves as an informal, data-dependent measure of evidence against the null hypothesis. It is rooted in *reduction to absurdity* reasoning, where the extremeness of the observed data is assessed under the assumption that the null hypothesis is true. However, the  $p$ -value is frequently misinterpreted as the probability that the null hypothesis is false—a misconception known as the  $p$ -value fallacy [180].

In contrast, the Neyman–Pearson framework adopts a deductive, long-run perspective: it defines decision rules that control the frequency of Type I errors over repeated sampling, irrespective of the evidence in any particular case. Conflating these two frameworks leads to interpretational inconsistencies, especially when the  $p$ -value is used both as a measure of evidence and

---

<sup>3</sup><https://fivethirtyeight.com/features/not-even-scientists-can-easily-explain-p-values/>

as a decision-making threshold. A clearer separation of these paradigms is essential for coherent statistical reasoning.

The American Statistical Association has several concerns regarding the use of the  $p$ -value as a cornerstone for hypothesis testing in science. This concern motivates the ASA's statement on  $p$ -values [447], which can be summarized in the following principles:

- “P-values can indicate how incompatible the data are with a specified statistical model.”
- “P-values do not measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone.”
- “Scientific conclusions and business or policy decisions should not be based solely on whether a  $p$ -value passes a specific threshold.”
- “Proper inference requires full reporting and transparency.”
- “A  $p$ -value, or statistical significance, does not measure the size of an effect or the importance of a result.”
- “By itself, a  $p$ -value does not provide a good measure of evidence regarding a model or hypothesis.”

To sum up, Fisher proposed the  $p$ -value as a witness rather than a judge. So, a  $p$ -value lower than the significance level means more inspection of the null hypothesis, but it is not a final conclusion about it.

Another key distinction between frequentist and Bayesian approaches lies in how scientific hypotheses are evaluated. Users of the Frequentist approach rely on the  $p$ -value, which quantifies the probability of observing data as extreme as—or more extreme than—the sample under the assumption that the null hypothesis is true. Bayesians, in contrast, use the Bayes factor, which compares the predictive performance of two competing hypotheses by evaluating the ratio of their marginal likelihoods. While the  $p$ -value reflects  $P(\text{data} \mid \text{hypothesis})$ , the Bayes factor is more aligned with  $P(\text{hypothesis} \mid \text{data})$ , though not equivalent. Notably, there exists an approximate relationship between the  $t$ -statistic and the Bayes factor in the context of regression coefficients [347], which offers some practical interpretability across paradigms. In particular,  $|t| > (\log(N) + 6)^{1/2}$  corresponds to strong evidence in favor of rejecting the null hypothesis of no relevance of a control in a regression. Observe that, in this setting, the threshold of the  $t$  statistic, and as a consequence the significance level, depends on the sample size. This setting agrees with the idea in experimental designs of selecting the sample size such that we control Type I and Type II errors. In observational studies, we cannot control the sample size, but we can select the significance level.

See also [403] and [36] for nice exercises that reveal potential flaws of the

$p$ -value ( $p$ ) due to  $p \sim U[0, 1]$  under the null hypothesis,<sup>4</sup> and calibrations of the  $p$ -value to interpret it as the odds ratio and the error probability. In particular,  $B(p) = -e \times p \times \log(p)$  when  $p < e^{-1}$ , and interpret this as the Bayes factor of  $H_0$  to  $H_1$ , where  $H_1$  denotes the unspecified alternative to  $H_0$ , and  $\alpha(p) = \left(1 + [-e \times p \times \log(p)]^{-1}\right)^{-1}$  as the error probability  $\alpha$  in rejecting  $H_0$ . Take into account that  $B(p)$  and  $\alpha(p)$  are lower bounds.

The logic of argumentation in the Frequentist approach is based on *deductive logic*, which means that it starts from a statement about the true state of nature (null hypothesis) and predicts what should be observed if this statement were true. On the other hand, the Bayesian approach is based on *inductive logic*, which means that it defines which hypothesis is more consistent with what is observed. The former inferential approach establishes that the truth of the premises implies the truth of the conclusion, which is why we reject or fail to reject hypotheses. The latter establishes that the premises supply some evidence, but not full assurance, of the truth of the conclusion, which is why we get probabilistic statements.

Here, there is a distinction between the effects of causes (forward causal inference) and the causes of effects (reverse causal inference) [163, 111]. To illustrate this point, imagine that a firm increases the price of a specific good. Economic theory would suggest that, as a result, demand for the good decreases. In this case, the premise (null hypothesis) is the price increase, and the consequence is the decrease in the firm's demand.

Alternatively, one could observe a reduction in a firm's demand and attempt to identify the cause behind it. For example, a reduction in quantity could be due to a negative supply shock. The Frequentist approach typically follows the first view (effects of causes), while Bayesian reasoning focuses on determining the probability of potential causes (causes of effects).

## 2.4 The likelihood principle

The **likelihood principle** states that in making inferences or decisions about the state of nature, all the relevant *experimental* information is given by the *likelihood function*. The Bayesian framework follows this statement, i.e., it is conditional on observed data.

We follow [37], who in turn followed [270], to illustrate the likelihood principle. We are given a coin and are interested in the probability,  $\theta$ , of it landing heads when flipped. We wish to test  $H_0 : \theta = 1/2$  versus  $H_1 : \theta > 1/2$ . An experiment is conducted by flipping the coin (independently) in a series of trials, with the result being the observation of 9 heads and 3 tails.

---

<sup>4</sup><https://joyeuserrancce.wordpress.com/2011/04/22/proof-that-p-values-under-the-null-are-uniformly-distributed/> for a simple proof.

This is not yet enough information to specify  $p(y|\theta)$ , since the series of trials has not been explained. Two possibilities arise:

1. The experiment consisted of a predetermined 12 flips, so that  $Y = [\text{Heads}]$  follows a  $B(12, \theta)$  distribution. In this case,  $p_1(y|\theta) = \binom{12}{y} \theta^y (1-\theta)^{12-y} = 220 \times \theta^9 (1-\theta)^3$ .
2. The experiment consisted of flipping the coin until 3 tails were observed ( $r = 3$ ). In this case,  $Y$ , the number of heads (failures) before obtaining 3 tails, follows a  $NB(3, 1-\theta)$  distribution. Here,  $p_2(y|\theta) = \binom{y+r-1}{r-1} (1-(1-\theta)^y) (1-\theta)^r = 55 \times \theta^9 (1-\theta)^3$ .

Using a Frequentist approach, the significance level of  $y = 9$  using the Binomial model against  $\theta = 1/2$  would be:

$$\alpha_1 = P_{1/2}(Y \geq 9) = p_1(9|1/2) + p_1(10|1/2) + p_1(11|1/2) + p_1(12|1/2) = 0.073.$$

#### **R code. The likelihood principle: Binomial model**

```
1 success <- 9 # Number of observed success in n trials
2 n <- 12 # Number of trials
3 siglevel <- sum(sapply(9:n, function(y) dbinom(y,n,0.5)))
```

For the Negative Binomial model, the significance level would be:

$$\alpha_2 = P_{1/2}(Y \geq 9) = p_2(9|1/2) + p_2(10|1/2) + \dots = 0.0327.$$

#### **R code. The likelihood principle: Negative Binomial model**

```
1 success <- 3 # Number of target success (tails)
2 failures <- 9 # Number of failures
3 siglevel <- 1 - pnbinom((failures - 1),success,0.5)
```

We arrive at different conclusions using a significance level of 5%, whereas we obtain the same outcomes using a Bayesian approach because the kernels of both distributions are identical ( $\theta^9 \times (1-\theta)^3$ ).

## 2.5 Why is not the Bayesian approach that popular?

At this stage, one might wonder why the Bayesian statistical framework is not the dominant inferential approach, despite its historical origin in 1763 [33], whereas the Frequentist statistical framework was largely developed in the early 20th century. The scientific debate over the Bayesian inferential approach lasted for 150 years, and this may be explained by some of the following factors.

One issue is the *apparent subjectivity* of the Bayesian approach, which runs counter to the strong conviction that science demands objectivity. Bayesian probability is considered a measure of degrees of belief, where the initial prior may be just a guess. This was not accepted as objective and rigorous science. Initial critics argued that Bayes was quantifying ignorance by assigning equal probabilities to all potential outcomes. As a consequence, prior distributions were dismissed [306].

*Bayes himself seemed not to have believed in his idea.* Although it seems that Bayes made his breakthrough in the late 1740s, he did not submit it for publication to the Royal Society. It was his friend, Richard Price, another Presbyterian minister, who rediscovered Bayes' idea, polished it, and published it.

However, it was Laplace who independently generalized Bayes' theorem in 1781. Initially, he applied it to gambling problems and soon thereafter to astronomy, combining various sources of information to advance research in situations where data was scarce. He later sought to apply his discovery to finding the probability of causes, which he thought required large datasets, thus turning to demography. In this field, he had to perform large-scale calculations, leading to the development of Laplace's approximation and the central limit theorem [256]. Unfortunately, this came at the cost of abandoning his research on Bayesian inference.

Once *Laplace passed away in 1827*, Bayes' rule disappeared from the scientific discourse for almost a century. In part, personal attacks against Laplace led to the rule being forgotten. Moreover, there was a prevailing belief that statistics should not address causation, and that the prior was too subjective to be compatible with science. Nonetheless, practitioners continued to use Bayes' rule to solve problems in astronomy, communication, medicine, military affairs, and social issues with remarkable results.

Thus, the concept of degrees of belief to operationalize probability was abandoned in favor of scientific objectivity. Probability was then defined as the frequency with which an event occurs in many repeatable trials, which became the accepted norm. Critics of Laplace argued that these two concepts were diametrically opposed, although Laplace considered them to be basically equivalent when large sample sizes are involved [306].

*The era of Frequentists, or sampling theorists, began*, led by Karl Pearson

and his nemesis, Ronald Fisher. Both were brilliant and persuasive characters, opposing the inverse probability (Bayesian) approach and making it nearly impossible to argue against their ideas. Pearson's legacy was carried on by his son, Egon, and Egon's friend, Jerzy Neyman. Both inherited the anti-Bayesian and anti-Fisher sentiments.

Despite the anti-Bayesian campaign among statisticians, some independent thinkers continued to develop Bayesian ideas, including Borel, Ramsey, and de Finetti, who were isolated in different countries: France, England, and Italy. However, the anti-Bayesian trio of Fisher, Neyman, and Egon Pearson dominated the spotlight in the 1920s and 1930s. Only a geophysicist, Harold Jeffreys, kept Bayesian inference alive during the 1930s and 1940s. Jeffreys was a quiet, reserved gentleman working in the astronomy department at Cambridge. He was Fisher's friend due to their shared character, although they were intellectual opposites when it came to Bayesian inference, leading to intense intellectual battles. Unfortunately for the Bayesian approach, *Jeffreys lost*. His work was highly technical, using confusing high-level mathematics. He focused on inference from scientific evidence, rather than guiding future actions based on decision theory, which was crucial in that era for mathematical statistics, especially during the Second World War. In contrast, Fisher was a dominant figure, persuasive in public and a master of practical applications, with his techniques written in a popular style with minimal mathematics.

Nevertheless, Bayes' rule achieved remarkable results in applied settings such as at AT&T and the U.S. Social Security system. Bayesian inference also played a significant role during the Second World War and the Cold War. Alan Turing used inverse probability at Bletchley Park to crack German messages encoded using the Enigma machine, which was employed by U-boats. Andrei Kolmogorov used Bayesian methods to improve firing tables for Russian artillery. Bernard Koopman applied it for searching targets at sea, and the RAND Corporation used it during the Cold War. Unfortunately, *these Bayesian developments remained top secret for almost 40 years*, keeping the contribution of inverse probability hidden from modern history.

In the 1950s and 1960s, three mathematicians led the resurgence of the Bayesian approach: Good, Savage, and Lindley. However, it seems that they were reluctant to apply their theories to real-world problems. Despite the fact that the Bayesian approach proved its worth in various areas such as business decisions, naval searches, and lung cancer detection, it was largely applied to simple models due to its *mathematical complexity and requirement for large computations*. However, some breakthroughs changed this. First, hierarchical models were introduced by Lindley and Smith, where a complex model is decomposed into many smaller, easier-to-solve models. Second, Markov chain Monte Carlo (MCMC) methods were developed by Hastings in the 1970s [196] and the Geman brothers in the 1980s [165]. These methods were incorporated into the Bayesian inferential framework in the 1990s by Gelfand and Smith [156], and Tierney [433], when desktop computers gained sufficient computational power to solve complex models. Since then, the Bayesian inferential

framework has gained increasing popularity among both practitioners and scientists.

## 2.6 A simple working example

We will illustrate some conceptual differences between the Bayesian and Frequentist statistical approaches by performing inference on a random sample  $\mathbf{Y} = [Y_1, Y_2, \dots, Y_N]$ , where  $Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$  for  $i = 1, 2, \dots, N$ .

In particular, we set  $\pi(\mu, \sigma) = \pi(\mu)\pi(\sigma) \propto \frac{1}{\sigma}$ . This is a standard *non-informative improper* prior (Jeffreys prior, see Chapter 3). That is, this prior is perfectly compatible with the sample information. Thus, the posterior distribution is

$$\begin{aligned}\pi(\mu, \sigma | \mathbf{y}) &\propto \frac{1}{\sigma} \times (\sigma^2)^{-N/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right\} \\ &= \frac{1}{\sigma} \times (\sigma^2)^{-N/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N ((y_i - \bar{y}) - (\mu - \bar{y}))^2 \right\} \\ &= \frac{1}{\sigma} \exp \left\{ -\frac{N}{2\sigma^2} (\mu - \bar{y})^2 \right\} \times (\sigma)^{-N} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{y})^2 \right\} \\ &= \frac{1}{\sigma} \exp \left\{ -\frac{N}{2\sigma^2} (\mu - \bar{y})^2 \right\} \times (\sigma)^{-(\alpha_n+1)} \exp \left\{ -\frac{\alpha_n \hat{\sigma}^2}{2\sigma^2} \right\},\end{aligned}$$

where  $\bar{y} = \frac{\sum_{i=1}^N y_i}{N}$ ,  $\alpha_n = N - 1$  and  $\hat{\sigma}^2 = \frac{\sum_{i=1}^N (y_i - \bar{y})^2}{N-1}$ .

The first term in the last expression is the kernel of a normal density,  $\mu | \sigma, \mathbf{y} \sim N(\bar{y}, \sigma^2/N)$ . The second term is the kernel of an inverted gamma density [465, p. 371],  $\sigma | \mathbf{y} \sim IG(\alpha_n, \hat{\sigma}^2)$ . Therefore,  $\pi(\mu | \sigma, \mathbf{y}) = (2\pi\sigma^2/N)^{-1/2} \exp \left\{ -\frac{N}{2\sigma^2} (\mu - \bar{y})^2 \right\}$  and  $\pi(\sigma | \mathbf{y}) = \frac{2}{\Gamma(\alpha_n/2)} \left( \frac{\alpha_n \hat{\sigma}^2}{2} \right)^{\alpha_n/2} \frac{1}{\sigma^{\alpha_n+1}} \times \exp \left\{ -\frac{\alpha_n \hat{\sigma}^2}{2\sigma^2} \right\}$ .

Observe that  $\mathbb{E}[\mu | \sigma, \mathbf{y}] = \bar{y}$ ; this is also the maximum likelihood (Frequentist) point estimate of  $\mu$  in this setting. In addition, the Frequentist  $(1 - \alpha)\%$  confidence interval and the Bayesian  $(1 - \alpha)\%$  credible interval have exactly the same form,  $\bar{y} \pm |z_{\alpha/2}| \frac{\sigma}{\sqrt{N}}$ , where  $z_{\alpha/2}$  is the  $\alpha/2$  percentile of a standard normal distribution. However, the interpretations are entirely different. The confidence interval has a probabilistic interpretation under sampling variability of  $\bar{Y}$ : in repeated sampling,  $(1 - \alpha)\%$  of the intervals  $\bar{Y} \pm |z_{\alpha/2}| \frac{\sigma}{\sqrt{N}}$  would include  $\mu$ . However, given an observed realization of  $\bar{Y}$ , say  $\bar{y}$ , the probability of  $\bar{y} \pm |z_{\alpha/2}| \frac{\sigma}{\sqrt{N}}$  including  $\mu$  is either 1 or 0. This is why we refer to it

## 44 Conceptual differences between the Bayesian and Frequentist approaches

as a  $(1 - \alpha)\%$  confidence interval. On the other hand,  $\bar{y} \pm |z_{\alpha/2}| \frac{\sigma}{\sqrt{N}}$  has a straightforward probabilistic interpretation in the Bayesian framework: there is a  $(1 - \alpha)\%$  probability that  $\mu$  lies within this interval.

If we want to get the marginal posterior density of  $\mu$ ,

$$\begin{aligned}\pi(\mu|\mathbf{y}) &= \int_0^\infty \pi(\mu, \sigma|\mathbf{y}) d\sigma \\ &\propto \int_0^\infty \frac{1}{\sigma} \times (\sigma^2)^{-N/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right\} d\sigma \\ &= \int_0^\infty \left( \frac{1}{\sigma} \right)^{N+1} \exp \left\{ -\frac{N}{2\sigma^2} \frac{\sum_{i=1}^N (y_i - \mu)^2}{N} \right\} d\sigma \\ &= \left[ \frac{2}{\Gamma(N/2)} \left( \frac{N \sum_{i=1}^N (y_i - \mu)^2}{2N} \right)^{N/2} \right]^{-1} \\ &\propto \left[ \sum_{i=1}^N (y_i - \mu)^2 \right]^{-N/2} \\ &= \left[ \sum_{i=1}^N ((y_i - \bar{y}) - (\mu - \bar{y}))^2 \right]^{-N/2} \\ &= [\alpha_n \hat{\sigma}^2 + N(\mu - \bar{y})^2]^{-N/2} \\ &\propto \left[ 1 + \frac{1}{\alpha_n} \left( \frac{\mu - \bar{y}}{\hat{\sigma}/\sqrt{N}} \right)^2 \right]^{-(\alpha_n+1)/2}.\end{aligned}$$

The fourth line arises from the kernel of an inverted gamma density with  $N$  degrees of freedom in the integral [465].

The last expression represents the kernel of a Student's  $t$  density function with  $\alpha_n = N - 1$  degrees of freedom, expected value equal to  $\bar{y}$ , and variance  $\frac{\hat{\sigma}^2}{N} \left( \frac{\alpha_n}{\alpha_n - 2} \right)$ . Therefore,  $\mu|\mathbf{y} \sim t\left(\bar{y}, \frac{\hat{\sigma}^2}{N} \left( \frac{\alpha_n}{\alpha_n - 2} \right), \alpha_n\right)$ .

Observe that a  $(1 - \alpha)\%$  confidence interval and a  $(1 - \alpha)\%$  credible interval have the same form,  $\bar{y} \pm |t_{\alpha/2}^{\alpha_n}| \frac{\hat{\sigma}}{\sqrt{N}}$ , where  $t_{\alpha/2}^{\alpha_n}$  is the  $\alpha/2$  percentile of a Student's  $t$  distribution. However, the interpretations are entirely different.

The mathematical similarity between the Frequentist and Bayesian expressions in this example arises from the use of an improper prior.

### 2.6.1 Example: Math test

You have a random sample of math scores of size  $N = 50$  from a normal distribution,  $Y_i \sim N(\mu, \sigma^2)$ . The sample mean and variance are equal to 102 and 10, respectively. Assuming an improper prior equal to  $1/\sigma$ , we proceed with the following tasks:

- Compute the 95% confidence and credible intervals for  $\mu$ .
- Determine the posterior probability that  $\mu > 103$ .

Using the fact that  $\mu|\mathbf{y} \sim t\left(\bar{y}, \frac{\hat{\sigma}^2}{N} \left(\frac{\alpha_n}{\alpha_n - 2}\right), \alpha_n\right)$ , which implies that the confidence and credible intervals for  $\mu$  are given by

$$\bar{y} \pm |t_{\alpha/2}^{\alpha_n}| \frac{\hat{\sigma}}{\sqrt{N}},$$

where  $\bar{y} = 102$ ,  $\hat{\sigma}^2 = 10$ , and  $\alpha_n = 49$ . Thus, the 95% confidence and credible intervals for  $\mu$  are the same, namely (101.1, 102.9), and the posterior probability that  $\mu > 103$  is 1.49% given the sample information.

### R code. Example: Math test

```

1 N <- 50 # Sample size
2 y_bar <- 102 # Sample mean
3 s2 <- 10 # Sample variance
4 alpha <- N - 1
5 serror <- (s2/N)^0.5
6 LimInf <- y_bar - abs(qt(0.025, alpha)) * serror # Lower
    bound
7 LimSup <- y_bar + abs(qt(0.025, alpha)) * serror # Upper
    bound
8 y.cut <- 103
9 P <- 1-metRology::pt.scaled(y.cut, df = alpha, mean = y_bar,
    sd = serror) # Probability of mu greater than y.cut

```

## 2.7 Summary

The differences between the Bayesian and Frequentist inferential approaches are philosophical, particularly with regard to the role of probability; pedagogical, especially in relation to the use of inference for decision-making; and methodological, due to differences in their mathematical and computational frameworks. Although, at the methodological level, the debate has become considerably muted—except for certain aspects of inference—there is widespread recognition that each approach has much to contribute to statistical practice [179, 31, 240]. As Bradley Efron stated, “Computer-age statistical inference at its most successful **combines** elements of the two philosophies” [133].

## 2.8 Exercises

### 1. Jeffreys-Lindley's Paradox

The **Jeffreys-Lindley's paradox** [223, 271] represents an apparent disagreement between the Bayesian and Frequentist frameworks in a hypothesis testing scenario.

In particular, assume that in a city, 49,581 boys and 48,870 girls have been born over 20 years. Assume that the male births follow a Binomial distribution with probability  $\theta$ . We wish to test the null hypothesis  $H_0 : \theta = 0.5$  versus the alternative hypothesis  $H_1 : \theta \neq 0.5$ .

- Show that the posterior model probability for the null model is approximately 0.95. Assume  $\pi(H_0) = \pi(H_1) = 0.5$ , and that  $\pi(\theta)$  follows a uniform distribution, i.e.,  $U(0, 1)$ , under  $H_1$ .
  - Show that the  $p$ -value for this hypothesis test is 0.0235 using the normal approximation,  $Y \sim N(N \times \theta, N \times \theta \times (1 - \theta))$ .
2. We want to test  $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$  given  $Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ .

Assume  $\pi(H_0) = \pi(H_1) = 0.5$ , and that  $\pi(\mu, \sigma) \propto 1/\sigma$  under the alternative hypothesis.

Show that

$$p(\mathbf{y}|\mathcal{M}_1) = \frac{\pi^{-N/2}}{2} \Gamma(N/2) \left( \frac{1}{\alpha_n \hat{\sigma}^2} \right)^{N/2} \left( \frac{N}{\alpha_n \hat{\sigma}^2} \right)^{-1/2} \frac{\Gamma(1/2)\Gamma(\alpha_n/2)}{\Gamma((\alpha_n + 1)/2)}$$

and

$$p(\mathbf{y}|\mathcal{M}_0) = (2\pi)^{-N/2} \left[ \frac{2}{\Gamma(N/2)} \left( \frac{N}{2} \frac{\sum_{i=1}^N (y_i - \mu_0)^2}{N} \right)^{N/2} \right]^{-1}.$$

Then, the posterior odds ratio is:

$$\begin{aligned} PO_{01} &= \frac{p(\mathbf{y}|\mathcal{M}_0)}{p(\mathbf{y}|\mathcal{M}_1)} \\ &= \frac{\Gamma((\alpha_n + 1)/2)}{\Gamma(1/2)\Gamma(\alpha_n/2)} (\alpha_n \hat{\sigma}^2 / N)^{-1/2} \left[ 1 + \frac{(\mu_0 - \bar{y})^2}{\alpha_n \hat{\sigma}^2 / N} \right]^{-\left(\frac{\alpha_n + 1}{2}\right)}, \end{aligned}$$

where  $\alpha_n = N - 1$  and  $\hat{\sigma}^2 = \frac{\sum_{i=1}^N (y_i - \bar{y})^2}{N-1}$ .

Find the relationship between the posterior odds ratio and the classical test statistic for the null hypothesis.

### 3. Math Test Continues

Using the setting of the **Example: Math Test** in subsection 2.6.1, test  $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$  where  $\mu_0 = \{100, 100.5, 101, 101.5, 102\}$ .

- What is the  $p$ -value for these hypothesis tests?
- Find the posterior model probability of the null model for each  $\mu_0$ .



# 3

---

## *Cornerstone models: Conjugate families*

---

We will introduce conjugate families, which are distributions for which the posterior distribution belongs to the same family as the prior distribution, given the likelihood. We provide some examples and solve them both analytically and computationally. We begin with simple examples of discrete and continuous distributions and then study the linear model in detail, both univariate and multivariate, deriving the posterior distributions, the marginal likelihood, and the predictive distribution analytically. Additionally, we will include mathematical and computational exercises in **R**.

---

### 3.1 Motivation of conjugate families

By observing the three fundamental pieces of Bayesian analysis –the posterior distribution (parameter inference), the marginal likelihood (hypothesis testing), and the predictive distribution (prediction)– as given in equations 3.1, 3.2, and 3.3, respectively, we can understand that some of the initial limitations of Bayesian analysis were due to the absence of algorithms for sampling from non-standard posterior distributions (equation 3.1), and the lack of analytical solutions for the marginal likelihood (equation 3.2) and the predictive distribution (equation 3.3), both of which require significant computational power.

$$\pi(\boldsymbol{\theta} | \mathbf{y}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})}{p(\mathbf{y})}, \quad (3.1)$$

$$p(\mathbf{y}) = \int_{\Theta} p(\mathbf{y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}, \quad (3.2)$$

and

$$p(\mathbf{y}_0 | \mathbf{y}) = \int_{\Theta} p(\mathbf{y}_0 | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}, \quad (3.3)$$

Although algorithms for sampling from non-standard posterior distributions have existed since the second half of the last century [308, 196, 165],

their application within the Bayesian framework emerged later [156, 433], likely coinciding with the rise in computational power of desktop computers. However, it is still common practice today to use models with standard conditional posterior distributions in order to reduce computational requirements. In addition, mathematical techniques coupled with computational algorithms [157, 89, 96] and approximations [434, 231] are employed to obtain the marginal likelihood (prior predictive).

Despite these advances, two potentially conflicting desirable model specification features are evident from equations 3.1, 3.2, and 3.3: (1) analytical solutions and (2) the posterior distribution belonging to the same family as the prior distribution for a given likelihood. The latter is known as *conjugate priors*, a family of priors that is closed under sampling [399].

These features are desirable because the former facilitates hypothesis testing and predictive analysis, while the latter ensures invariance in the prior-to-posterior updating process. Both features reduce computational burden.

Although each of these features can be achieved independently –such as using improper priors for analytical tractability and broadly defining the family of priors for conjugacy– these features are in conflict.

Fortunately, we can achieve both characteristics if we assume that the data-generating process follows a distribution function in the *exponential family*. That is, given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$ , a probability density function  $p(\mathbf{y} \mid \boldsymbol{\theta})$  belongs to the exponential family if it has the form:

$$\begin{aligned} p(\mathbf{y} \mid \boldsymbol{\theta}) &= \prod_{i=1}^N h(y_i) C(\boldsymbol{\theta}) \exp \{ \eta(\boldsymbol{\theta})^\top \mathbf{T}(y_i) \} \\ &= h(\mathbf{y}) C(\boldsymbol{\theta})^N \exp \{ \eta(\boldsymbol{\theta})^\top \mathbf{T}(\mathbf{y}) \} \\ &= h(\mathbf{y}) \exp \{ \eta(\boldsymbol{\theta})^\top \mathbf{T}(\mathbf{y}) - A(\boldsymbol{\theta}) \}, \end{aligned} \quad (3.4)$$

where  $h(\mathbf{y}) = \prod_{i=1}^N h(y_i)$  is a non-negative function,  $\eta(\boldsymbol{\theta})$  is a known function of the parameters, and  $A(\boldsymbol{\theta}) = \log \{ \int_{\mathbf{Y}} h(\mathbf{y}) \exp \{ \eta(\boldsymbol{\theta})^\top \mathbf{T}(\mathbf{y}) \} d\mathbf{y} \} = -N \log(C(\boldsymbol{\theta}))$  is the normalization factor. Additionally,  $\mathbf{T}(\mathbf{y}) = \sum_{i=1}^N \mathbf{T}(y_i)$  is the vector of sufficient statistics for the distribution (by the factorization theorem).

If the support of  $\mathbf{Y}$  is independent of  $\boldsymbol{\theta}$ , the family is said to be *regular*; otherwise, it is *irregular*. Furthermore, if we set  $\eta = \eta(\boldsymbol{\theta})$ , the exponential family is said to be in the *canonical form*.

$$\begin{aligned} p(\mathbf{y} \mid \boldsymbol{\eta}) &= h(\mathbf{y}) D(\boldsymbol{\eta})^N \exp \{ \boldsymbol{\eta}^\top \mathbf{T}(\mathbf{y}) \} \\ &= h(\mathbf{y}) \exp \{ \boldsymbol{\eta}^\top \mathbf{T}(\mathbf{y}) - B(\boldsymbol{\eta}) \}. \end{aligned}$$

A nice feature of this representation is that  $\mathbb{E}[\mathbf{T}(\mathbf{y}) \mid \boldsymbol{\eta}] = \nabla B(\boldsymbol{\eta})$  and  $Var[\mathbf{T}(\mathbf{y}) \mid \boldsymbol{\eta}] = \nabla^2 B(\boldsymbol{\eta})$ .

### 3.1.1 Examples of exponential family distributions

#### 1. Discrete distributions

Let's show that some of the most common distributions for random variables, which can take values on a finite or countably infinite set, are part of the exponential family.

##### Poisson distribution

Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from a *Poisson distribution* let's show that  $p(\mathbf{y} | \lambda)$  is in the exponential family.

$$\begin{aligned} p(\mathbf{y} | \lambda) &= \prod_{i=1}^N \frac{\lambda^{y_i} \exp(-\lambda)}{y_i!} \\ &= \frac{\lambda^{\sum_{i=1}^N y_i} \exp(-N\lambda)}{\prod_{i=1}^N y_i!} \\ &= \frac{\exp(-N\lambda) \exp(\sum_{i=1}^N y_i \log(\lambda))}{\prod_{i=1}^N y_i!}, \end{aligned}$$

then  $h(\mathbf{y}) = \left[ \prod_{i=1}^N y_i! \right]^{-1}$ ,  $\eta(\lambda) = \log(\lambda)$ ,  $T(\mathbf{y}) = \sum_{i=1}^N y_i$  (sufficient statistic) and  $C(\lambda) = \exp(-\lambda)$ .

If we set  $\eta = \log(\lambda)$ , then

$$p(\mathbf{y} | \eta) = \frac{\exp(\eta \sum_{i=1}^N y_i - N \exp(\eta))}{\prod_{i=1}^N y_i!},$$

such that  $B(\eta) = N \exp(\eta)$ , then  $\nabla(B(\eta)) = N \exp(\eta) = N\lambda = \mathbb{E} \left[ \sum_{i=1}^N y_i \middle| \lambda \right]$ , that is,  $\mathbb{E} \left[ \frac{\sum_{i=1}^N y_i}{N} \middle| \lambda \right] = \mathbb{E}[\bar{y} \mid \lambda] = \lambda$ , and  $\nabla^2(B(\eta)) = N \exp(\eta) = N\lambda = \text{Var} \left[ \sum_{i=1}^N y_i \middle| \lambda \right] = N^2 \times \text{Var} [\bar{y} \mid \lambda]$ , then  $\text{Var} [\bar{y} \mid \lambda] = \frac{\lambda}{N}$ .

##### Bernoulli distribution

Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from a *Bernoulli distribution* let's show that  $p(\mathbf{y} | \theta)$  is in the exponential family.

$$\begin{aligned} p(\mathbf{y} | \theta) &= \prod_{i=1}^N \theta^{y_i} (1-\theta)^{1-y_i} \\ &= \theta^{\sum_{i=1}^N y_i} (1-\theta)^{N-\sum_{i=1}^N y_i} \\ &= (1-\theta)^N \exp \left\{ \sum_{i=1}^N y_i \log \left( \frac{\theta}{1-\theta} \right) \right\}, \end{aligned}$$

then  $h(\mathbf{y}) = \mathbb{1}[y_i \in \{0, 1\}]$  (indicator function),  $\eta(\theta) = \log\left(\frac{\theta}{1-\theta}\right)$ ,  
 $T(\mathbf{y}) = \sum_{i=1}^N y_i$  and  $C(\theta) = 1 - \theta$ .

Write this distribution in the canonical form, and find the mean and variance of the sufficient statistic (Exercise 1).

### Multinomial distribution

Given a random sample  $\mathbf{Y} = [\mathbf{Y}_1 \ \mathbf{Y}_2 \ \dots \ \mathbf{Y}_N]$  from a *m-dimensional multinomial distribution*, where  $\mathbf{Y}_i = [Y_{i1} \ Y_{i2} \ \dots \ Y_{im}]$ ,  $\sum_{l=1}^m Y_{il} = n$ ,  $n$  independent trials each of which leads to a success for exactly one of  $m$  categories with probabilities  $\boldsymbol{\theta} = [\theta_1 \ \theta_2 \ \dots \ \theta_m]$ ,  $\sum_{l=1}^m \theta_l = 1$ . Let's show that  $p(\mathbf{y} | \boldsymbol{\theta})$  is in the exponential family.

$$\begin{aligned} p(\mathbf{y} | \boldsymbol{\theta}) &= \prod_{i=1}^N \frac{n!}{\prod_{l=1}^m y_{il}!} \prod_{l=1}^m \theta_l^{y_{il}} \\ &= \frac{(n!)^N}{\prod_{i=1}^N \prod_{l=1}^m y_{il}!} \exp \left\{ \sum_{i=1}^N \sum_{l=1}^m y_{il} \log(\theta_l) \right\} \\ &= \frac{(n!)^N}{\prod_{i=1}^N \prod_{l=1}^m y_{il}!} \exp \left\{ \left( N \times n - \sum_{i=1}^N \sum_{l=1}^{m-1} y_{il} \right) \log(\theta_m) \right. \\ &\quad \left. + \sum_{i=1}^N \sum_{l=1}^{m-1} y_{il} \log(\theta_l) \right\} \\ &= \frac{(n!)^N}{\prod_{i=1}^N \prod_{l=1}^m y_{il}!} \theta_m^{N \times n} \exp \left\{ \sum_{i=1}^N \sum_{l=1}^{m-1} y_{il} \log(\theta_l / \theta_m) \right\}, \end{aligned}$$

then  $h(\mathbf{y}) = \frac{(n!)^N}{\prod_{i=1}^N \prod_{l=1}^m y_{il}!}$ ,  $\eta(\boldsymbol{\theta}) = \left[ \log\left(\frac{\theta_1}{\theta_m}\right) \dots \log\left(\frac{\theta_{m-1}}{\theta_m}\right) \right]$ ,  
 $T(\mathbf{y}) = \left[ \sum_{i=1}^N y_{i1} \dots \sum_{i=1}^N y_{im-1} \right]$  and  $C(\boldsymbol{\theta}) = \theta_m^n$ .

## 2. Continuous distributions

Let's show that some of the most common distributions for random variables, which can take any value within a certain range or interval –an infinite number of possible values– are part of the exponential family.

### Normal distribution

Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from a *normal distribution* let's show that  $p(\mathbf{y} | \mu, \sigma^2)$  is in the exponential family.

$$p(\mathbf{y} | \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{2\pi\sigma^2} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - \mu)^2 \right\}$$

$$\begin{aligned}
&= (2\pi)^{-N/2} (\sigma^2)^{-N/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right\} \\
&= (2\pi)^{-N/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N y_i^2 + \frac{\mu}{\sigma^2} \sum_{i=1}^N y_i \right. \\
&\quad \left. - N \frac{\mu^2}{2\sigma^2} - \frac{N}{2} \log(\sigma^2) \right\},
\end{aligned}$$

then  $h(\mathbf{y}) = (2\pi)^{-N/2}$ ,  $\eta(\mu, \sigma^2) = \begin{bmatrix} \frac{\mu}{\sigma^2} & \frac{-1}{2\sigma^2} \end{bmatrix}$ ,  $T(\mathbf{y}) = \begin{bmatrix} \sum_{i=1}^N y_i & \sum_{i=1}^N y_i^2 \end{bmatrix}$  and  $C(\mu, \sigma^2) = \exp \left\{ -\frac{\mu^2}{2\sigma^2} - \frac{\log(\sigma^2)}{2} \right\}$ .

Observe that

$$p(\mathbf{y} | \mu, \sigma^2) = (2\pi)^{-N/2} \exp \left\{ \eta_1 \sum_{i=1}^N y_i + \eta_2 \sum_{i=1}^N y_i^2 - \frac{N}{2} \log(-2\eta_2) + \frac{N}{4} \frac{\eta_1^2}{\eta_2} \right\},$$

where  $B(\boldsymbol{\eta}) = \frac{N}{2} \log(-2\eta_2) - \frac{N}{4} \frac{\eta_1^2}{\eta_2}$ . Then,

$$\nabla B(\boldsymbol{\eta}) = \begin{bmatrix} -\frac{N}{2} \frac{\eta_1}{\eta_2} \\ -\frac{N}{2} \frac{1}{\eta_2} + \frac{N}{4} \frac{\eta_1^2}{\eta_2^2} \end{bmatrix} = \begin{bmatrix} N \times \mu \\ N \times (\mu^2 + \sigma^2) \end{bmatrix} = \begin{bmatrix} \mathbb{E} \left[ \sum_{i=1}^N y_i | \mu, \sigma^2 \right] \\ \mathbb{E} \left[ \sum_{i=1}^N y_i^2 | \mu, \sigma^2 \right] \end{bmatrix}.$$

### Multivariate normal distribution

Given  $\mathbf{Y} = [\mathbf{Y}_1 \ \mathbf{Y}_2 \ \dots \ \mathbf{Y}_p]$  a  $N \times p$  matrix such that  $\mathbf{Y}_i \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $i = 1, 2, \dots, N$ , that is, each  $i$ -th row of  $\mathbf{Y}$  follows a *multivariate normal distribution*. Then, assuming independence between rows, let's show that  $p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is in the exponential family.

$$\begin{aligned}
p(\mathbf{y}_1, \dots, \mathbf{y}_N | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \prod_{i=1}^N (2\pi)^{-p/2} |\boldsymbol{\Sigma}|^{-1/2} \\
&\quad \times \exp \left\{ -\frac{1}{2} (\mathbf{y}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \right\} \\
&= (2\pi)^{-pN/2} |\boldsymbol{\Sigma}|^{-N/2} \\
&\quad \times \exp \left\{ -\frac{1}{2} \text{tr} \left[ \sum_{i=1}^N (\mathbf{y}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \right] \right\} \\
&= (2\pi)^{-pN/2} |\boldsymbol{\Sigma}|^{-N/2} \\
&\quad \times \exp \left\{ -\frac{1}{2} \text{tr} \left[ (\mathbf{S} + N(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})^\top) \boldsymbol{\Sigma}^{-1} \right] \right\}
\end{aligned}$$

$$\begin{aligned}
&= (2\pi)^{-pN/2} \\
&\times \exp \left\{ -\frac{1}{2} \left[ \left( \text{vec}(\mathbf{S})^\top + N \text{vec}(\hat{\boldsymbol{\mu}} \hat{\boldsymbol{\mu}}^\top)^\top \right) \text{vec}(\boldsymbol{\Sigma}^{-1}) \right. \right. \\
&- 2N \hat{\boldsymbol{\mu}}^\top \text{vec}(\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}) + N \text{tr}(\boldsymbol{\mu} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}) \\
&\left. \left. + N \log(|\boldsymbol{\Sigma}|) \right] \right\},
\end{aligned}$$

where the second line uses the trace operator ( $\text{tr}$ ), and its invariance under cyclic permutation is applied in the third line. Additionally, we add and subtract  $\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i$  inside each parenthesis, resulting in  $\mathbf{S} = \sum_{i=1}^N (\mathbf{y}_i - \hat{\boldsymbol{\mu}})(\mathbf{y}_i - \hat{\boldsymbol{\mu}})^\top$ . The fourth line is obtained after collecting terms and using properties of the trace operator to introduce the vectorization operator ( $\text{vec}$ ), specifically,  $\text{tr}(\mathbf{A}^\top \mathbf{B}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$ , and  $\text{vec}(\mathbf{A} + \mathbf{B}) = \text{vec}(\mathbf{A}) + \text{vec}(\mathbf{B})$ .

Then  $h(\mathbf{y}) = (2\pi)^{-pN/2}$ ,  $\eta(\boldsymbol{\mu}, \boldsymbol{\Sigma})^\top = \left[ (\text{vec}(\boldsymbol{\Sigma}^{-1}))^\top \quad (\text{vec}(\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}))^\top \right]$ ,  $T(\mathbf{y}) = \left[ -\frac{1}{2} \left( \text{vec}(\mathbf{S})^\top + N \text{vec}(\hat{\boldsymbol{\mu}} \hat{\boldsymbol{\mu}}^\top)^\top \right) \quad -N \hat{\boldsymbol{\mu}}^\top \right]^\top$  and  $C(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \exp \left\{ -\frac{1}{2} (\text{tr}(\boldsymbol{\mu} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}) + \log(|\boldsymbol{\Sigma}|)) \right\}$ .

### 3.2 Conjugate prior to exponential family

#### Theorem 4.2.1

The prior distribution  $\pi(\boldsymbol{\theta}) \propto C(\boldsymbol{\theta})^{b_0} \exp \{ \eta(\boldsymbol{\theta})^\top \mathbf{a}_0 \}$  is conjugate to the exponential family (equation 3.4).

#### Proof

$$\begin{aligned}
\pi(\boldsymbol{\theta} | \mathbf{y}) &\propto C(\boldsymbol{\theta})^{b_0} \exp \{ \eta(\boldsymbol{\theta})^\top \mathbf{a}_0 \} \times h(\mathbf{y}) C(\boldsymbol{\theta})^N \exp \{ \eta(\boldsymbol{\theta})^\top \mathbf{T}(\mathbf{y}) \} \\
&\propto C(\boldsymbol{\theta})^{N+b_0} \exp \{ \eta(\boldsymbol{\theta})^\top (\mathbf{T}(\mathbf{y}) + \mathbf{a}_0) \}.
\end{aligned}$$

Observe that the posterior is in the exponential family,  $\pi(\boldsymbol{\theta} | \mathbf{y}) \propto C(\boldsymbol{\theta})^{\beta_n} \exp \{ \eta(\boldsymbol{\theta})^\top \boldsymbol{\alpha}_n \}$ ,  $\beta_n = N + b_0$  and  $\boldsymbol{\alpha}_n = \mathbf{T}(\mathbf{y}) + \mathbf{a}_0$ .

#### Remarks

We observe, by comparing the prior and the likelihood, that  $b_0$  plays the role of a hypothetical sample size, and  $\mathbf{a}_0$  plays the role of hypothetical sufficient statistics. This perspective aids the elicitation process, that is, integrating non-sample information into the prior distribution.

We established this result in the *standard form* of the exponential family. We can also establish it in the *canonical form* of the exponential family. Observe that, given  $\boldsymbol{\eta} = \eta(\boldsymbol{\theta})$ , another way to derive a prior for  $\boldsymbol{\eta}$  is to use the change of variable theorem, given a bijective function.

In the case where there is a regular conjugate prior, [118] show that the posterior expectation of the sufficient statistics is a weighted average between the prior expectation and the likelihood estimate.

### 3.2.1 Examples: Theorem 4.2.1

#### 1. Likelihood functions from discrete distributions

##### The Poisson-gamma model

Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from a Poisson distribution then a conjugate prior density for  $\lambda$  has the form

$$\begin{aligned}\pi(\lambda) &\propto (\exp(-\lambda))^{b_0} \exp\{a_0 \log(\lambda)\} \\ &= \exp(-\lambda b_0) \lambda^{a_0} \\ &= \exp(-\lambda \beta_0) \lambda^{\alpha_0 - 1}.\end{aligned}$$

This is the kernel of a gamma density in the *rate parametrization*,  $G(\alpha_0, \beta_0)$ , where  $\alpha_0 = a_0 + 1$  and  $\beta_0 = b_0$ .<sup>1</sup> Thus, a prior conjugate distribution for the Poisson likelihood is a gamma distribution.

Since  $\sum_{i=1}^N Y_i$  is a sufficient statistic for the Poisson distribution, we can interpret  $a_0$  as the number of occurrences in  $b_0$  experiments.

Observe that

$$\begin{aligned}\pi(\lambda \mid \mathbf{y}) &\propto \exp(-\lambda \beta_0) \lambda^{\alpha_0 - 1} \times \exp(-N\lambda) \lambda^{\sum_{i=1}^N y_i} \\ &= \exp(-\lambda(N + \beta_0)) \lambda^{\sum_{i=1}^N y_i + \alpha_0 - 1}.\end{aligned}$$

As expected, this is the kernel of a gamma distribution, which means  $\lambda \mid \mathbf{y} \sim G(\alpha_n, \beta_n)$ ,  $\alpha_n = \sum_{i=1}^N y_i + \alpha_0$  and  $\beta_n = N + \beta_0$ .

Observe that  $\alpha_0/\beta_0$  is the prior mean, and  $\alpha_0/\beta_0^2$  is the prior variance. Then,  $\alpha_0 \rightarrow 0$  and  $\beta_0 \rightarrow 0$  imply a non-informative prior such that the posterior mean converges to the maximum likelihood estimate  $\bar{y} = \frac{\sum_{i=1}^N y_i}{N}$ ,

$$\begin{aligned}\mathbb{E}[\lambda \mid \mathbf{y}] &= \frac{\alpha_n}{\beta_n} \\ &= \frac{\sum_{i=1}^N y_i + \alpha_0}{N + \beta_0} \\ &= \frac{N\bar{y}}{N + \beta_0} + \frac{\alpha_0}{N + \beta_0}.\end{aligned}$$

The posterior mean is a weighted average of the sample and prior

---

<sup>1</sup>Another parametrization of the gamma density is the *scale parametrization*, where  $\kappa_0 = 1/\beta_0$ . See the health insurance example in Chapter 1.

information. This is a general result for regular conjugate priors [118]. Note that  $\lim_{N \rightarrow \infty} \mathbb{E}[\lambda | \mathbf{y}] = \bar{y}$ .

Additionally,  $\alpha_0 \rightarrow 0$  and  $\beta_0 \rightarrow 0$  corresponds to  $\pi(\lambda) \propto \frac{1}{\lambda}$ , which is an improper prior. Improper priors may have undesirable consequences for Bayes factors (hypothesis testing); see below for a discussion of this in the linear regression framework. In this example, we can obtain analytical solutions for the marginal likelihood and the predictive distribution (see the health insurance example and Exercise 3 in Chapter 1).

### The Bernoulli-beta model

Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from a Bernoulli distribution then a conjugate prior density for  $\theta$  has the form

$$\begin{aligned}\pi(\theta) &\propto (1 - \theta)^{b_0} \exp \left\{ a_0 \log \left( \frac{\theta}{1 - \theta} \right) \right\} \\ &= (1 - \theta)^{b_0 - a_0} \theta^{a_0} \\ &= \theta^{\alpha_0 - 1} (1 - \theta)^{\beta_0 - 1}.\end{aligned}$$

This is the kernel of a beta density,  $B(\alpha_0, \beta_0)$ , where  $\alpha_0 = a_0 + 1$  and  $\beta_0 = b_0 - a_0 + 1$ . A prior conjugate distribution for the Bernoulli likelihood is a beta distribution. Given that  $b_0$  is the hypothetical sample size and  $a_0$  is the hypothetical sufficient statistic (the number of successes),  $b_0 - a_0$  represents the number of failures. This implies that  $\alpha_0$  is the number of prior successes plus one, and  $\beta_0$  is the number of prior failures plus one.

Since the mode of a beta-distributed random variable is given by  $\frac{\alpha_0 - 1}{\alpha_0 + \beta_0 - 2} = \frac{a_0}{b_0}$ , we can interpret this as the prior probability of success. Setting  $\alpha_0 = 1$  and  $\beta_0 = 1$ , which corresponds to a uniform distribution on the interval  $[0, 1]$ , represents a setting with 0 successes (and 0 failures) in 0 experiments.

Observe that

$$\begin{aligned}\pi(\theta | \mathbf{y}) &\propto \theta^{\alpha_0 - 1} (1 - \theta)^{\beta_0 - 1} \times \theta^{\sum_{i=1}^N y_i} (1 - \theta)^{N - \sum_{i=1}^N y_i} \\ &= \theta^{\alpha_0 + \sum_{i=1}^N y_i - 1} (1 - \theta)^{\beta_0 + N - \sum_{i=1}^N y_i - 1}.\end{aligned}$$

The posterior distribution is beta,  $\theta | \mathbf{y} \sim B(\alpha_n, \beta_n)$ ,  $\alpha_n = \alpha_0 + \sum_{i=1}^N y_i$  and  $\beta_n = \beta_0 + N - \sum_{i=1}^N y_i$ , where the posterior mean  $\mathbb{E}[\theta | \mathbf{y}] = \frac{\alpha_n}{\alpha_n + \beta_n} = \frac{\alpha_0 + N \bar{y}}{\alpha_0 + \beta_0 + N} = \frac{\alpha_0 + \beta_0}{\alpha_0 + \beta_0 + N} \frac{\alpha_0}{\alpha_0 + \beta_0} + \frac{N}{\alpha_0 + \beta_0 + N} \bar{y}$ . The posterior mean is a weighted average between the prior mean and the maximum likelihood estimate.

The marginal likelihood in this setting is

$$p(\mathbf{y}) = \int_0^1 \frac{\theta^{\alpha_0 - 1} (1 - \theta)^{\beta_0 - 1}}{B(\alpha_0, \beta_0)} \times \theta^{\sum_{i=1}^N y_i} (1 - \theta)^{N - \sum_{i=1}^N y_i} d\theta$$

$$= \frac{B(\alpha_n, \beta_n)}{B(\alpha_0, \beta_0)},$$

where  $B(\cdot, \cdot)$  is the beta function.

In addition, the predictive density is

$$\begin{aligned} p(y_0 | \mathbf{y}) &= \int_0^1 \theta^{y_0} (1-\theta)^{1-y_0} \times \frac{\theta^{\alpha_n-1} (1-\theta)^{\beta_n-1}}{B(\alpha_n, \beta_n)} d\theta \\ &= \frac{B(\alpha_n + y_0, \beta_n + 1 - y_0)}{B(\alpha_n, \beta_n)} \\ &= \frac{\Gamma(\alpha_n + \beta_n) \Gamma(\alpha_n + y_0) \Gamma(\beta_n + 1 - y_0)}{\Gamma(\alpha_n + \beta_n + 1) \Gamma(\alpha_n) \Gamma(\beta_n)} \\ &= \begin{cases} \frac{\alpha_n}{\alpha_n + \beta_n}, & y_0 = 1 \\ \frac{\beta_n}{\alpha_n + \beta_n}, & y_0 = 0 \end{cases}. \end{aligned}$$

This is a Bernoulli distribution with probability of success equal to  $\frac{\alpha_n}{\alpha_n + \beta_n}$ .

### The multinomial-Dirichlet model

Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from a multinomial distribution then a conjugate prior density for  $\boldsymbol{\theta} = [\theta_1 \ \theta_2 \ \dots \ \theta_m]$  has the form

$$\begin{aligned} \pi(\boldsymbol{\theta}) &\propto \theta_m^{b_0} \exp \{ \boldsymbol{\eta}(\boldsymbol{\theta})^\top \mathbf{a}_0 \} \\ &= \prod_{l=1}^{m-1} \theta_l^{a_{0l}} \theta_m^{b_0 - \sum_{l=1}^{m-1} a_{0l}} \\ &= \prod_{l=1}^m \theta_l^{\alpha_{0l}-1}, \end{aligned}$$

where  $\boldsymbol{\eta}(\boldsymbol{\theta}) = \left[ \log \left( \frac{\theta_1}{\theta_m} \right) \ \dots \ \log \left( \frac{\theta_{m-1}}{\theta_m} \right) \right]$ ,  $\mathbf{a}_0 = [a_{01} \ \dots \ a_{0m-1}]^\top$ ,  $\boldsymbol{\alpha}_0 = [\alpha_{01} \ \alpha_{02} \ \dots \ \alpha_{0m}]$ ,  $\alpha_{0l} = a_{0l} + 1$ ,  $l = 1, 2, \dots, m-1$  and  $\alpha_{0m} = b_0 - \sum_{l=1}^{m-1} a_{0l} + 1$ .

This is the kernel of a Dirichlet distribution, that is, the prior distribution is  $D(\boldsymbol{\alpha}_0)$ .

Observe that  $a_{0l}$  is the hypothetical number of times outcome  $l$  is observed over the hypothetical  $b_0$  trials. Setting  $\alpha_{0l} = 1$ , which corresponds to a uniform distribution over the open standard simplex, implicitly sets  $a_{0l} = 0$ , meaning that there are 0 occurrences of category  $l$  in  $b_0 = 0$  experiments.

The posterior distribution of the multinomial-Dirichlet model is given by

$$\begin{aligned}\pi(\boldsymbol{\theta} \mid \mathbf{y}) &\propto \prod_{l=1}^m \theta_l^{\alpha_{0l}-1} \times \prod_{l=1}^m \theta_l^{\sum_{i=1}^N y_{il}} \\ &= \prod_{l=1}^m \theta_l^{\alpha_{0l} + \sum_{i=1}^N y_{il} - 1}.\end{aligned}$$

This is the kernel of a Dirichlet distribution  $D(\boldsymbol{\alpha}_n)$ ,  $\boldsymbol{\alpha}_n = [\alpha_{n1} \alpha_{n2} \dots \alpha_{nm}]$ ,  $\alpha_{nl} = \alpha_{0l} + \sum_{i=1}^N y_{il}$ ,  $l = 1, 2, \dots, m$ . Observe that

$$\begin{aligned}\mathbb{E}[\theta_j \mid \mathbf{y}] &= \frac{\alpha_{nj}}{\sum_{l=1}^m [\alpha_{0l} + \sum_{i=1}^N y_{il}]} \\ &= \frac{\sum_{l=1}^m \alpha_{0l}}{\sum_{l=1}^m [\alpha_{0l} + \sum_{i=1}^N y_{il}]} \frac{\alpha_{0j}}{\sum_{l=1}^m \alpha_{0l}} \\ &+ \frac{\sum_{l=1}^m \sum_{i=1}^N y_{il}}{\sum_{l=1}^m [\alpha_{0l} + \sum_{i=1}^N y_{il}]} \frac{\sum_{i=1}^N y_{ij}}{\sum_{l=1}^m \sum_{i=1}^N y_{il}}.\end{aligned}$$

We have again that the posterior mean is a weighted average between the prior mean and the maximum likelihood estimate.

The marginal likelihood is

$$\begin{aligned}p(\mathbf{y}) &= \int_{\Theta} \frac{\prod_{l=1}^m \theta_l^{\alpha_{0l}-1}}{B(\boldsymbol{\alpha}_0)} \times \prod_{i=1}^N \frac{n!}{\prod_{l=1}^m y_{il}!} \prod_{l=1}^m \theta_l^{y_{il}} d\boldsymbol{\theta} \\ &= \frac{N \times n!}{B(\boldsymbol{\alpha}_0) \prod_{i=1}^N \prod_{l=1}^m y_{il}!} \int_{\Theta} \prod_{l=1}^m \theta_l^{\alpha_{0l} + \sum_{i=1}^N y_{il} - 1} d\boldsymbol{\theta} \\ &= \frac{N \times n!}{B(\boldsymbol{\alpha}_0) \prod_{i=1}^N \prod_{l=1}^m y_{il}!} B(\boldsymbol{\alpha}_n) \\ &= \frac{N \times n! \Gamma(\sum_{l=1}^m \alpha_{0l})}{\Gamma(\sum_{l=1}^m \alpha_{0l} + N \times n) \prod_{l=1}^m \Gamma(\alpha_{0l}) \prod_{i=1}^N y_{il}!},\end{aligned}$$

where  $B(\boldsymbol{\alpha}) = \frac{\prod_{l=1}^m \Gamma(\alpha_l)}{\Gamma(\sum_{l=1}^m \alpha_l)}$ .

Following similar steps we get the predictive density

$$p(y_0 \mid \mathbf{y}) = \frac{n! \Gamma(\sum_{l=1}^m \alpha_{nl})}{\Gamma(\sum_{l=1}^m \alpha_{nl} + n)} \prod_{l=1}^m \frac{\Gamma(\alpha_{nl} + y_{0l})}{\Gamma(\alpha_{nl}) y_{0l}!}.$$

This is a Dirichlet-multinomial distribution with parameters  $\boldsymbol{\alpha}_n$ .

### Example: English premier league, Liverpool vs Manchester city

Let's consider an example using data from the English Premier League. In particular, we want to calculate the probability that, in the next five matches between Liverpool and Manchester City, Liverpool wins two games and Manchester City wins three. This calculation is based on historical data from the last five matches where Liverpool played at home between January 14th, 2018, and April 10th, 2022. In those matches, Liverpool secured two wins, there were two draws, and Manchester City won one match.<sup>2</sup>

We use two strategies to estimate the hyperparameters. First, we estimate the hyperparameters of the Dirichlet distribution using betting odds from bookmakers at 19:05 on October 6th, 2022 (Colombia time). We obtained data from 24 bookmakers (see file *DataOddsLIVvsMAN.csv*),<sup>3</sup> and we transform these odds into probabilities using a simple standardization approach. Then, we apply maximum likelihood estimation to estimate the hyperparameters.

Second, we use empirical Bayes, where we estimate the hyperparameters by optimizing the marginal likelihood.

#### **R code. Multinomial-Dirichlet model: Liverpool vs Manchester city**

```

1 # Multinomial-Dirichlet example: Liverpool vs Manchester
  city
2 Data <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTAApp/refs/heads/master/DataApp/
  DataOddsLIVvsMAN.csv", sep = ",", header = TRUE, quote =
  "")
3 attach(Data); library(dplyr)
4 Probs <- Data %>% mutate(pns1 = 1/home, pns2 = 1/draw, pns3
  = 1/away) %>% mutate(SumInvOdds = pns1 + pns2 + pns3)
  %>% mutate(p1 = pns1/SumInvOdds, p2 = pns2/SumInvOdds,
  p3 = pns3/SumInvOdds) %>% select(p1, p2, p3)
5 # We get probabilities using simple standardization. There
  are more technical approaches to do this. See for
  instance Shin (1993) and Strumbelj (2014).

```

<sup>2</sup><https://www.11v11.com/teams/manchester-city/tab/opposingTeams/opposition/Liverpool/>.

<sup>3</sup>[https://www.oddsportal.com/soccer/england/premier-league/liverpool-manchester-](https://www.oddsportal.com/soccer/england/premier-league/liverpool-manchester-city-WrqgEz5S/)

### R code. Multinomial-Dirichlet model: Liverpool vs Manchester city

```

1 DirMLE <- sirt::dirichlet.mle(Probs) # Use maximum
   likelihood to estimate parameters of the Dirichlet
   distribution
2 alpha0odds <- DirMLE$alpha
3 y <- c(2, 2, 1)
4 # Historical records last five matches: Liverpool wins (2),
   draws (2) and Manchester city wins (1)
5 # Marginal likelihood
6 MarLik <- function(a0){
7   n <- sum(y)
8   Res1 <- sum(sapply(1:length(y),
9     function(l){lgamma(a0[1]+y[l])-lgamma(a0[1]))})
10  Res <- lgamma(sum(a0))-lgamma(sum(a0)+n)+Res1
11  return(-Res)
12 }
13 EmpBay <- optim(alpha0odds, MarLik, method = "BFGS")
14 alpha0EB <- EmpBay$par
15 # Bayes factor empirical Bayes vs betting odds.
16 # This is greater than 1 by construction
17 BF <- exp(-MarLik(alpha0EB))/exp(-MarLik(alpha0odds))
18 # Posterior distribution based on empirical Bayes
19 alphan <- alpha0EB + y
20 # Posterior parameters
21 S <- 100000
22 # Simulation draws from the Dirichlet distribution
23 thetas <- MCMCpack::rdirichlet(S, alphan)
24 colnames(thetas) <- c("Liverpool", "Draw", "Manchester")
25 # Predictive distribution based on simulations
26 y0 <- c(2, 0, 3)
27 # Liverpool two wins and Manchester city three wins in next
   five matches
28 Pred <- apply(thetas, 1, function(p) {rmultinom(1, size =
   sum(y0), prob = p)})
29 ProY0 <- sum(sapply(1:S, function(s){sum(Pred[,s]==y0)==3}))/
   S # Probability of y0
30 # Predictive distribution using analytical expression
31 PredY0 <- function(y0){
32   n <- sum(y0)
33   Res1 <- sum(sapply(1:length(y), function(l){lgamma(alphan[
   1]+y0[1]) - lgamma(alphan[1])-lfactorial(y0[1]))})
34   Res <- lfactorial(n) + lgamma(sum(alphan)) - lgamma(sum(
   alphan)+n) + Res1
35   return(exp(Res))
36 }
```

We observe that the Bayes factor provides evidence in favor of the

hyperparameters estimated via empirical Bayes, as these hyperparameters are specifically chosen to maximize the marginal likelihood.

Using the hyperparameters obtained from empirical Bayes, we calculate that the probability of Liverpool winning two out of the next five games, while Manchester City wins three, is 1.2%. The result obtained from the predictive distribution via simulations is similar to the probability derived using the exact predictive distribution.

## 2. Likelihood functions from continuous distributions

### The normal-normal/inverse-gamma model

Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from a normal distribution, then the conjugate prior density has the form

$$\begin{aligned}\pi(\mu, \sigma^2) &\propto \exp \left\{ b_0 \left( -\frac{\mu^2}{2\sigma^2} - \frac{\log \sigma^2}{2} \right) \right\} \exp \left\{ a_{01} \frac{\mu}{\sigma^2} - a_{02} \frac{1}{\sigma^2} \right\} \\ &= \exp \left\{ b_0 \left( -\frac{\mu^2}{2\sigma^2} - \frac{\log \sigma^2}{2} \right) \right\} \exp \left\{ a_{01} \frac{\mu}{\sigma^2} - a_{02} \frac{1}{\sigma^2} \right\} \\ &\quad \times \exp \left\{ -\frac{a_{01}^2}{2\sigma^2 b_0} \right\} \exp \left\{ \frac{a_{01}^2}{2\sigma^2 b_0} \right\} \\ &= \exp \left\{ -\frac{b_0}{2\sigma^2} \left( \mu - \frac{a_{01}}{b_0} \right)^2 \right\} \left( \frac{1}{\sigma^2} \right)^{\frac{b_0+1-1}{2}} \\ &\quad \times \exp \left\{ \frac{1}{\sigma^2} \frac{-2b_0 a_{02} + a_{01}^2}{2b_0} \right\} \\ &= \underbrace{\left( \frac{1}{\sigma^2} \right)^{\frac{1}{2}} \exp \left\{ -\frac{b_0}{2\sigma^2} \left( \mu - \frac{a_{01}}{b_0} \right)^2 \right\}}_1 \\ &\quad \times \underbrace{\left( \frac{1}{\sigma^2} \right)^{\frac{b_0-1}{2}} \exp \left\{ -\frac{1}{\sigma^2} \frac{2b_0 a_{02} - a_{01}^2}{2b_0} \right\}}_2.\end{aligned}$$

The first part is the kernel of a normal density with mean  $\mu_0 = \frac{a_{01}}{\beta_0}$  and variance  $\frac{\sigma^2}{\beta_0}$ , where  $\beta_0 = b_0$ . That is,  $\mu \mid \sigma^2 \sim N \left( \mu_0, \frac{\sigma^2}{\beta_0} \right)$ . The second part is the kernel of an inverse gamma density with shape parameter  $\frac{\alpha_0}{2} = \frac{\beta_0-3}{2}$  and scale parameter  $\frac{\delta_0}{2} = \frac{2\beta_0 a_{02} - a_{01}^2}{2\beta_0}$ , so  $\sigma^2 \sim IG \left( \frac{\alpha_0}{2}, \frac{\delta_0}{2} \right)$ .

Observe that  $b_0 = \beta_0$  represents the hypothetical sample size, and  $a_{01}$  is the hypothetical sum of prior observations. Therefore, it makes sense that  $\frac{a_{01}}{\beta_0}$  and  $\frac{\sigma^2}{\beta_0}$  represent the prior mean and variance, respectively.

Therefore, the posterior distribution is also a normal-inverse gamma distribution,

$$\begin{aligned}
\pi(\mu, \sigma^2 | \mathbf{y}) &\propto \left(\frac{1}{\sigma^2}\right)^{1/2} \exp\left\{-\frac{\beta_0}{2\sigma^2}(\mu - \mu_0)^2\right\} \left(\frac{1}{\sigma^2}\right)^{\alpha_0/2+1} \\
&\quad \times \exp\left\{-\frac{\delta_0}{2\sigma^2}\right\} (\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2\right\} \\
&= \left(\frac{1}{\sigma^2}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^2} \left(\beta_0(\mu - \mu_0)^2 + \sum_{i=1}^N (y_i - \bar{y})^2 + \delta_0\right.\right. \\
&\quad \left.\left.+ N(\mu - \bar{y})^2 + \frac{(\beta_0\mu_0 + N\bar{y})^2}{\beta_0 + N} - \frac{(\beta_0\mu_0 + N\bar{y})^2}{\beta_0 + N}\right)\right\} \\
&\quad \times \left(\frac{1}{\sigma^2}\right)^{\frac{\alpha_0+N}{2}+1} \\
&= \exp\left\{-\frac{1}{2\sigma^2} \left((\beta_0 + N) \left(\mu - \left(\frac{\beta_0\mu_0 + N\bar{y}}{\beta_0 + N}\right)\right)^2\right)\right\} \\
&\quad \times \left(\frac{1}{\sigma^2}\right)^{1/2} \\
&\quad \times \exp\left\{-\frac{1}{2\sigma^2} \left(\sum_{i=1}^N (y_i - \bar{y})^2 + \delta_0 + \frac{\beta_0 N}{\beta_0 + N}(\bar{y} - \mu_0)^2\right)\right\} \\
&\quad \times \left(\frac{1}{\sigma^2}\right)^{\frac{\alpha_0+N}{2}+1}.
\end{aligned}$$

The first term is the kernel of a normal density,  $\mu | \sigma^2, \mathbf{y} \sim N(\mu_n, \sigma_n^2)$ , where  $\mu_n = \frac{\beta_0\mu_0 + N\bar{y}}{\beta_0 + N}$  and  $\sigma_n^2 = \frac{\sigma^2}{\beta_n}$ ,  $\beta_n = \beta_0 + N$ . The second term is the kernel of an inverse gamma density,  $\sigma^2 | \mathbf{y} \sim IG(\alpha_n/2, \delta_n/2)$  where  $\alpha_n = \alpha_0 + N$  and  $\delta_n = \sum_{i=1}^N (y_i - \bar{y})^2 + \delta_0 + \frac{\beta_0 N}{\beta_0 + N}(\bar{y} - \mu_0)^2$ . Observe that the posterior mean is a weighted average between prior and sample information. The weights depends on the sample sizes ( $\beta_0$  and  $N$ ).

The marginal posterior for  $\sigma^2$  is inverse gamma with shape and scale parameters  $\alpha_n/2$  and  $\delta_n/2$ , respectively. The marginal posterior of  $\mu$  is

$$\begin{aligned}
\pi(\mu | \mathbf{y}) &\propto \int_0^\infty \left\{ \left(\frac{1}{\sigma^2}\right)^{\frac{\alpha_n+1}{2}+1} \exp\left\{-\frac{1}{2\sigma^2}(\beta_n(\mu - \mu_n)^2 + \delta_n)\right\} \right\} d\sigma^2 \\
&= \frac{\Gamma\left(\frac{\alpha_n+1}{2}\right)}{\left[\frac{\beta_n(\mu - \mu_n)^2 + \delta_n}{2}\right]^{\frac{\alpha_n+1}{2}}}
\end{aligned}$$

$$\begin{aligned} &\propto \left[ \frac{\beta_n(\mu - \mu_n)^2 + \delta_n}{2} \right]^{-\frac{\alpha_n+1}{2}} \left( \frac{\delta_n}{\delta_n} \right)^{-\frac{\alpha_n+1}{2}} \\ &\propto \left[ \frac{\alpha_n \beta_n (\mu - \mu_n)^2}{\alpha_n \delta_n} + 1 \right]^{-\frac{\alpha_n+1}{2}}, \end{aligned}$$

The second line follows from having the kernel of an inverse gamma density with parameters  $\frac{\alpha_n+1}{2}$  and  $\frac{1}{2}(\beta_n(\mu - \mu_n)^2 + \delta_n)$ .

This corresponds to the kernel of a Student's  $t$ -distribution:

$$\mu \mid \mathbf{y} \sim t \left( \mu_n, \frac{\delta_n}{\beta_n \alpha_n}, \alpha_n \right),$$

where  $\mathbb{E}[\mu \mid \mathbf{y}] = \mu_n$  and

$$\text{Var}[\mu \mid \mathbf{y}] = \frac{\alpha_n}{\alpha_n - 2} \left( \frac{\delta_n}{\beta_n \alpha_n} \right) = \frac{\delta_n}{(\alpha_n - 2)\beta_n}, \quad \alpha_n > 2.$$

Observe that the marginal posterior distribution for  $\mu$  has heavier tails than the conditional posterior distribution due to the incorporation of uncertainty regarding  $\sigma^2$ .

The marginal likelihood is

$$\begin{aligned} p(\mathbf{y}) &= \int_{-\infty}^{\infty} \int_0^{\infty} \left\{ (2\pi\sigma^2/\beta_0)^{-1/2} \exp \left\{ -\frac{1}{2\sigma^2/\beta_0} (\mu - \mu_0)^2 \right\} \right. \\ &\quad \times \frac{(\delta_0/2)^{\alpha_0/2}}{\Gamma(\alpha_0/2)} \left( \frac{1}{\sigma^2} \right)^{\alpha_0/2+1} \exp \left\{ -\frac{\delta_0}{2\sigma^2} \right\} (2\pi\sigma^2)^{-N/2} \\ &\quad \times \left. \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right\} \right\} d\sigma^2 d\mu \\ &= \frac{(\delta_0/2)^{\alpha_0/2}}{\Gamma(\alpha_0/2)} (2\pi)^{-(\frac{N+1}{2})} \beta_0^{1/2} \int_{-\infty}^{\infty} \int_0^{\infty} \left\{ \left( \frac{1}{\sigma^2} \right)^{\frac{\alpha_0+N+1}{2}+1} \right. \\ &\quad \times \left. \exp \left\{ -\frac{1}{2\sigma^2} (\beta_0(\mu - \mu_0)^2 + \sum_{i=1}^N (y_i - \mu)^2 + \delta_0) \right\} \right\} d\sigma^2 d\mu \\ &= \frac{(\delta_0/2)^{\alpha_0/2}}{\Gamma(\alpha_0/2)} (2\pi)^{-(\frac{N+1}{2})} \beta_0^{1/2} \Gamma \left( \frac{N+1+\alpha_0}{2} \right) \\ &\quad \times \int_{-\infty}^{\infty} \left[ \frac{\beta_0(\mu - \mu_0)^2 + \sum_{i=1}^N (y_i - \mu)^2 + \delta_0}{2} \right]^{-\frac{\alpha_0+N+1}{2}} d\mu \\ &= \frac{(\delta_0/2)^{\alpha_0/2}}{\Gamma(\alpha_0/2)} (2\pi)^{-(\frac{N+1}{2})} \beta_0^{1/2} \Gamma \left( \frac{N+1+\alpha_0}{2} \right) \\ &\quad \times \int_{-\infty}^{\infty} \left[ \frac{\beta_n(\mu - \mu_n)^2 + \delta_n}{2} \right]^{-\frac{\alpha_n+1}{2}} d\mu \left( \frac{\delta_n/2}{\delta_n/2} \right)^{-\frac{\alpha_n+1}{2}} \end{aligned}$$

$$\begin{aligned}
&= \frac{(\delta_0/2)^{\alpha_0/2}}{\Gamma(\alpha_0/2)} (2\pi)^{-\left(\frac{N+1}{2}\right)} \beta_0^{1/2} \Gamma\left(\frac{\alpha_n+1}{2}\right) \left(\frac{\delta_n}{2}\right)^{-\frac{\alpha_n+1}{2}} \\
&\times \frac{\left(\frac{\delta_n\pi}{\beta_n}\right)^{1/2} \Gamma\left(\frac{\alpha_n}{2}\right)}{\Gamma\left(\frac{\alpha_n+1}{2}\right)} \\
&= \frac{\Gamma\left(\frac{\alpha_n}{2}\right)}{\Gamma\left(\frac{\alpha_0}{2}\right)} \frac{(\delta_0/2)^{\alpha_0/2}}{(\delta_n/2)^{\alpha_n/2}} \left(\frac{\beta_0}{\beta_n}\right)^{1/2} (\pi)^{-N/2},
\end{aligned}$$

where we take into account that  $\int_{-\infty}^{\infty} \left[ \frac{\beta_n(\mu - \mu_n)^2 + \delta_n}{2} \right]^{-\frac{\alpha_n+1}{2}} d\mu \left( \frac{\delta_n/2}{\delta_n/2} \right)^{-\frac{\alpha_n+1}{2}} = \int_{-\infty}^{\infty} \left[ \frac{\beta_n\alpha_n(\mu - \mu_n)^2}{\delta_n\alpha_n} + 1 \right]^{-\frac{\alpha_n+1}{2}} d\mu \left( \frac{\delta_n}{2} \right)^{-\frac{\alpha_n+1}{2}}$ . The term in the integral is the kernel of a Student's t density, this means that the integral is equal to  $\frac{\left(\frac{\delta_n\pi}{\beta_n}\right)^{1/2} \Gamma\left(\frac{\alpha_n}{2}\right)}{\Gamma\left(\frac{\alpha_n+1}{2}\right)}$ .

The predictive density is

$$\begin{aligned}
\pi(y_0 | \mathbf{y}) &\propto \int_{-\infty}^{\infty} \int_0^{\infty} \left\{ \left( \frac{1}{\sigma^2} \right)^{1/2} \exp \left\{ -\frac{1}{2\sigma^2} (y_0 - \mu)^2 \right\} \left( \frac{1}{\sigma^2} \right)^{1/2} \right. \\
&\quad \times \left. \exp \left\{ -\frac{\beta_n}{2\sigma^2} (\mu - \mu_n)^2 \right\} \left( \frac{1}{\sigma^2} \right)^{\alpha_n/2+1} \exp \left\{ -\frac{\delta_n}{2\sigma^2} \right\} \right\} d\sigma^2 d\mu \\
&= \int_{-\infty}^{\infty} \int_0^{\infty} \left\{ \exp \left\{ -\frac{1}{2\sigma^2} ((y_0 - \mu)^2 + \beta_n(\mu - \mu_n)^2 + \delta_n) \right\} \right. \\
&\quad \times \left. \left( \frac{1}{\sigma^2} \right)^{\frac{\alpha_n+1}{2}+1} \right\} d\sigma^2 d\mu \\
&\propto \int_{-\infty}^{\infty} [\beta_n(\mu - \mu_n)^2 + (y_0 - \mu)^2 + \delta_n]^{-\left(\frac{\alpha_n}{2}+1\right)} d\mu \\
&= \int_{-\infty}^{\infty} \left[ (\beta_n + 1) \left( \mu - \left( \frac{\beta_n\mu_n + y_0}{\beta_n + 1} \right) \right)^2 + \frac{\beta_n(y_0 - \mu_n)^2}{\beta_n + 1} \right. \\
&\quad \left. + \delta_n \right]^{-\left(\frac{\alpha_n}{2}+1\right)} d\mu \\
&= \int_{-\infty}^{\infty} \left[ 1 + \frac{(\alpha_n + 1)(\beta_n + 1)^2 \left( \mu - \left( \frac{\beta_n\mu_n + y_0}{\beta_n + 1} \right) \right)^2}{(\alpha_n + 1)(\beta_n(y_0 - \mu_n)^2 + (\beta_n + 1)\delta_n)} \right]^{-\left(\frac{\alpha_n}{2}+1\right)} d\mu \\
&\times \left( \frac{\beta_n(y_0 - \mu_n)^2 + (\beta_n + 1)\delta_n}{\beta_n + 1} \right)^{-\left(\frac{\alpha_n}{2}+1\right)} \\
&\propto \left( \pi \frac{\beta_n(y_0 - \mu_n)^2 + (\beta_n + 1)\delta_n}{(\beta_n + 1)^2} \right)^{\frac{1}{2}}
\end{aligned}$$

$$\begin{aligned}
& \times \left( \frac{\beta_n(y_0 - \mu_n)^2 + (\beta_n + 1)\delta_n}{\beta_n + 1} \right)^{-\left(\frac{\alpha_n}{2}+1\right)} \\
& \propto (\beta_n(y_0 - \mu_n)^2 + (\beta_n + 1)\delta_n)^{-\left(\frac{\alpha_n+1}{2}\right)} \\
& \propto \left[ 1 + \frac{\beta_n\alpha_n}{(\beta_n + 1)\delta_n\alpha_n}(y_0 - \mu_n)^2 \right]^{-\left(\frac{\alpha_n+1}{2}\right)},
\end{aligned}$$

where we have that  $\left[ 1 + \frac{(\alpha_n+1)(\beta_n+1)^2(\mu - (\frac{\beta_n\mu_n+y_0}{\beta_n+1}))^2}{(\alpha_n+1)(\beta_n(y_0-\mu_n)^2+(\beta_n+1)\delta_n)} \right]^{-\left(\frac{\alpha_n}{2}+1\right)}$  is the kernel of a Student's t density with degrees of freedom  $\alpha_n + 1$  and scale  $\frac{\beta_n(y_0-\mu_n)^2+(\beta_n+1)\delta_n}{(\beta_n+1)^2(\alpha_n+1)}$ .

The last expression is the kernel of a Student's t density, that is,  $Y_0 | \mathbf{y} \sim t\left(\mu_n, \frac{(\beta_n+1)\delta_n}{\beta_n\alpha_n}, \alpha_n\right)$ .

### The multivariate normal-normal/inverse-Wishart model

We show in subsection 3.1 that the multivariate normal distribution is in the exponential family where

$$\begin{aligned}
C(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \exp \left\{ -\frac{1}{2} \left( \text{tr}(\boldsymbol{\mu}\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}) + \log(|\boldsymbol{\Sigma}|) \right) \right\}, \\
\eta(\boldsymbol{\mu}, \boldsymbol{\Sigma})^\top &= \left[ \left( \text{vec}(\boldsymbol{\Sigma}^{-1}) \right)^\top \quad \left( \text{vec}(\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}) \right)^\top \right], \\
T(\mathbf{y}) &= \left[ -\frac{1}{2} \left( \text{vec}(\mathbf{S})^\top + N \text{vec}(\hat{\boldsymbol{\mu}}\hat{\boldsymbol{\mu}}^\top)^\top \right) \quad -N\hat{\boldsymbol{\mu}}^\top \right]^\top
\end{aligned}$$

and

$$h(\mathbf{y}) = (2\pi)^{-pN/2}.$$

Then, its conjugate prior distribution should have the form

$$\begin{aligned}
\pi(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &\propto \exp \left\{ -\frac{b_0}{2} \left( \text{tr}(\boldsymbol{\mu}\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}) + \log(|\boldsymbol{\Sigma}|) \right) \right\} \\
&\times \exp \left\{ \mathbf{a}_{01}^\top \text{vec}(\boldsymbol{\Sigma}^{-1}) + \mathbf{a}_{02}^\top \text{vec}(\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1}) \right\} \\
&= |\boldsymbol{\Sigma}|^{-b_0/2} \exp \left\{ -\frac{b_0}{2} \left( \text{tr}(\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) + \text{tr}(\mathbf{a}_{02}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right) \right\} \\
&\times \exp \left\{ \mathbf{a}_{01}^\top \text{vec}(\boldsymbol{\Sigma}^{-1}) + \frac{\mathbf{a}_{02}^\top \boldsymbol{\Sigma}^{-1} \mathbf{a}_{02}}{2b_0} - \frac{\mathbf{a}_{02}^\top \boldsymbol{\Sigma}^{-1} \mathbf{a}_{02}}{2b_0} \right\} \\
&= |\boldsymbol{\Sigma}|^{-b_0/2} \exp \left\{ -\frac{b_0}{2} \left( \boldsymbol{\mu} - \frac{\mathbf{a}_{02}}{b_0} \right)^\top \boldsymbol{\Sigma}^{-1} \left( \boldsymbol{\mu} - \frac{\mathbf{a}_{02}}{b_0} \right) \right\} \\
&\times \exp \left\{ -\frac{1}{2} \text{tr} \left( \left( \mathbf{A}_{01} - \frac{\mathbf{a}_{02}\mathbf{a}_{02}^\top}{b_0} \right) \boldsymbol{\Sigma}^{-1} \right) \right\}
\end{aligned}$$

$$\begin{aligned}
&= \underbrace{|\Sigma|^{-1/2} \exp \left\{ -\frac{b_0}{2} \left( \boldsymbol{\mu} - \frac{\mathbf{a}_{02}}{b_0} \right)^\top \boldsymbol{\Sigma}^{-1} \left( \boldsymbol{\mu} - \frac{\mathbf{a}_{02}}{b_0} \right) \right\}}_1 \\
&\quad \times \underbrace{|\Sigma|^{-(\alpha_0+p+1)/2} \exp \left\{ -\frac{1}{2} \text{tr} \left( \left( \mathbf{A}_{01} - \frac{\mathbf{a}_{02}\mathbf{a}_{02}^\top}{b_0} \right) \boldsymbol{\Sigma}^{-1} \right) \right\}}_2.
\end{aligned}$$

Here, we set  $b_0 = 1 + \alpha_0 + p + 1$ , and represents the hypothetical sample size,  $\mathbf{a}_{01}$  and  $\mathbf{a}_{02}$  are  $p^2$ -dimensional and  $p$ -dimensional vectors of prior sufficient statistics, respectively. Specifically,  $\mathbf{a}_{01} = -\frac{1}{2}\text{vec}(\mathbf{A}_{01})$ , where  $\mathbf{A}_{01}$  is a  $p \times p$  positive semi-definite matrix.

We observe that the first part of the last expression is the kernel of a multivariate normal density with mean  $\boldsymbol{\mu}_0 = \frac{\mathbf{a}_{02}}{b_0}$  and covariance  $\frac{\boldsymbol{\Sigma}}{b_0}$ , i.e.,

$$\boldsymbol{\mu} | \boldsymbol{\Sigma} \sim N_p \left( \boldsymbol{\mu}_0, \frac{\boldsymbol{\Sigma}}{\beta_0} \right),$$

where  $b_0 = \beta_0$ . This choice of hyperparameters is intuitive because  $\mathbf{a}_{02}$  represents the hypothetical sum of prior observations, and  $b_0$  represents the hypothetical prior sample size.

Additionally, the second part of the last expression corresponds to the kernel of an inverse Wishart distribution with scale matrix  $\boldsymbol{\Psi}_0 = \left( \mathbf{A}_{01} - \frac{\mathbf{a}_{02}\mathbf{a}_{02}^\top}{b_0} \right)$  and  $\alpha_0$  degrees of freedom, i.e.,

$$\boldsymbol{\Sigma} \sim IW_p(\boldsymbol{\Psi}_0, \alpha_0).$$

Observe that  $\boldsymbol{\Psi}_0$  has the same structure as the first part of the sufficient statistics in  $T(\mathbf{y})$ , except that it should be understood as arising from prior hypothetical observations.

Therefore, the prior distribution in this setting is normal/inverse-Wishart, and, due to conjugacy, the posterior distribution belongs to the same family.

$$\begin{aligned}
\pi(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{y}) &\propto (2\pi)^{-pN/2} |\boldsymbol{\Sigma}|^{-N/2} \\
&\quad \times \exp \left\{ -\frac{1}{2} \text{tr} \left[ \left( \mathbf{S} + N(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})^\top \right) \boldsymbol{\Sigma}^{-1} \right] \right\} \\
&\quad \times |\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -\frac{\beta_0}{2} \text{tr} [(\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}^{-1}] \right\} \\
&\quad \times |\boldsymbol{\Sigma}|^{-(\alpha_0+p+1)/2} \exp \left\{ -\frac{1}{2} \text{tr} (\boldsymbol{\Psi}_0 \boldsymbol{\Sigma}^{-1}) \right\}.
\end{aligned}$$

Taking into account that  $N(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})^\top + \beta_0(\boldsymbol{\mu} - \boldsymbol{\mu}_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^\top$  is equal to

$$(N + \beta_0)(\boldsymbol{\mu} - \boldsymbol{\mu}_n)(\boldsymbol{\mu} - \boldsymbol{\mu}_n)^\top + \frac{N\beta_0}{N + \beta_0}(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_0)(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_0)^\top,$$

where  $\boldsymbol{\mu}_n = \frac{N}{N + \beta_0}\hat{\boldsymbol{\mu}} + \frac{\beta_0}{N + \beta_0}\boldsymbol{\mu}_0$  is the posterior mean. We have

$$\begin{aligned} \pi(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{y}) &\propto |\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -\frac{N + \beta_0}{2} \text{tr} \left[ \left( (\boldsymbol{\mu} - \boldsymbol{\mu}_n)(\boldsymbol{\mu} - \boldsymbol{\mu}_n)^\top \right) \boldsymbol{\Sigma}^{-1} \right] \right\} \\ &\quad \times |\boldsymbol{\Sigma}|^{-(N + \alpha_0 + p + 1)/2} \\ &\quad \times \exp \left\{ -\frac{1}{2} \text{tr} \left[ \left( \frac{N\beta_0}{N + \beta_0}(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_0)(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_0)^\top \right. \right. \right. \\ &\quad \left. \left. \left. + \boldsymbol{\Psi}_0 + \mathbf{S} \right) \boldsymbol{\Sigma}^{-1} \right] \right\}. \end{aligned}$$

Then,  $\boldsymbol{\mu} | \boldsymbol{\Sigma}, \mathbf{y} \sim N_p \left( \boldsymbol{\mu}_n, \frac{1}{\beta_n} \boldsymbol{\Sigma} \right)$ , and  $\boldsymbol{\Sigma} | \mathbf{y} \sim IW(\boldsymbol{\Psi}_n, \alpha_n)$  where  $\beta_n = N + \beta_0$ ,  $\alpha_n = N + \alpha_0$  and  $\boldsymbol{\Psi}_n = \boldsymbol{\Psi}_0 + \mathbf{S} + \frac{N\beta_0}{N + \beta_0}(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_0)(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_0)^\top$ .

The marginal posterior of  $\boldsymbol{\mu}$  is given by  $\int_{\mathcal{S}} \pi(\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\boldsymbol{\Sigma}$  where  $\mathcal{S}$  is the space of positive semi-definite matrices. Then,

$$\begin{aligned} \pi(\boldsymbol{\mu} | \mathbf{y}) &\propto \int_{\mathcal{S}} \left\{ |\boldsymbol{\Sigma}|^{-(\alpha_n + p + 2)/2} \right. \\ &\quad \left. \exp \left\{ -\frac{1}{2} \text{tr} \left[ \left( \beta_n(\boldsymbol{\mu} - \boldsymbol{\mu}_n)(\boldsymbol{\mu} - \boldsymbol{\mu}_n)^\top + \boldsymbol{\Psi}_n \right) \boldsymbol{\Sigma}^{-1} \right] \right\} \right\} d\boldsymbol{\Sigma} \\ &\propto |\left( \beta_n(\boldsymbol{\mu} - \boldsymbol{\mu}_n)(\boldsymbol{\mu} - \boldsymbol{\mu}_n)^\top + \boldsymbol{\Psi}_n \right)|^{-(\alpha_n + 1)/2} \\ &= \left[ |\boldsymbol{\Psi}_n| \times |1 + \beta_n(\boldsymbol{\mu} - \boldsymbol{\mu}_n)^\top \boldsymbol{\Psi}_n^{-1}(\boldsymbol{\mu} - \boldsymbol{\mu}_n)| \right]^{-(\alpha_n + 1)/2} \\ &\propto \left( 1 + \frac{1}{\alpha_n + 1 - p}(\boldsymbol{\mu} - \boldsymbol{\mu}_n)^\top \left( \frac{\boldsymbol{\Psi}_n}{(\alpha_n + 1 - p)\beta_n} \right)^{-1} \right. \\ &\quad \left. (\boldsymbol{\mu} - \boldsymbol{\mu}_n) \right)^{-(\alpha_n + 1 - p + p)/2}, \end{aligned}$$

where the second line uses properties of the inverse Wishart distribution, and the third line uses a particular case of the Sylvester's determinant theorem.<sup>4</sup>

We observe that the last line is the kernel of a multivariate t distribution, that is,  $\boldsymbol{\mu} | \mathbf{y} \sim t_p(v_n, \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$  where  $v_n = \alpha_n + 1 - p$  and  $\boldsymbol{\Sigma}_n = \frac{\boldsymbol{\Psi}_n}{(\alpha_n + 1 - p)\beta_n}$ .

<sup>4</sup> $\det(\mathbf{X} + \mathbf{AB}) = \det(\mathbf{X})\det(\mathbf{I} + \mathbf{B}\mathbf{X}^{-1}\mathbf{A})$ .

The marginal likelihood is given by

$$p(\mathbf{y}) = \frac{\Gamma_p\left(\frac{v_n}{2}\right)}{\Gamma_p\left(\frac{\alpha_0}{2}\right)} \frac{|\Psi_0|^{\alpha_0/2}}{|\Psi_n|^{\alpha_n/2}} \left(\frac{\beta_0}{\beta_n}\right)^{p/2} (2\pi)^{-Np/2},$$

where  $\Gamma_p$  is the multivariate gamma function (see Exercise 5).

The posterior predictive distribution is  $\mathbf{Y}_0 \mid \mathbf{y} \sim t_p(v_n, \boldsymbol{\mu}_n, (\beta_n + 1)\boldsymbol{\Sigma}_n)$  (see Exercise 6).

### Example: Tangency portfolio of US tech stocks

The tangency portfolio is the portfolio that maximizes the Sharpe ratio, which is defined as the excess return of a portfolio standardized by its risk.

We aim to find the portfolio weights  $\mathbf{w}$  that maximize the Sharpe ratio, where  $\mu_{i,T+\kappa} = \mathbb{E}(R_{i,T+\kappa} - R_{f,T+\kappa} \mid \mathcal{I}_T)$ , with  $R_{i,T+\kappa}$  and  $R_{f,T+\kappa}$  representing the returns of stock  $i$  and the risk-free asset, respectively. Here,  $\mu_{i,T+\kappa}$  is the expected value of the excess return at period  $T + \kappa$ , conditional on information available up to time  $T$  ( $\mathcal{I}_T$ ), and  $\boldsymbol{\Sigma}_{T+\kappa}$  is the covariance matrix of the excess returns, which quantifies the risk.

$$\arg \max_{\mathbf{w} \in \mathbb{R}^p} \frac{\mathbf{w}^\top \boldsymbol{\mu}_{T+\kappa}}{\sqrt{\mathbf{w}^\top \boldsymbol{\Sigma}_{T+\kappa} \mathbf{w}}}; \quad \text{s.t.} \quad \mathbf{w}^\top \mathbf{1} = 1,$$

where the solution is

$$\mathbf{w}^* = \frac{\boldsymbol{\Sigma}_{T+\kappa}^{-1} \boldsymbol{\mu}_{T+\kappa}}{\mathbf{1}^\top \boldsymbol{\Sigma}_{T+\kappa}^{-1} \boldsymbol{\mu}_{T+\kappa}}.$$

If we want to find the optimal portfolio for the next period under the assumption that the excess returns follow a multivariate normal distribution –a common assumption in these applications– we can set  $\kappa = 1$  and use the predictive distribution of the excess returns. In this case,  $\boldsymbol{\mu}_{T+1} = \boldsymbol{\mu}_n$  and  $\boldsymbol{\Sigma}_{T+1} = \frac{v_n}{v_n-2}(\beta_n + 1)\boldsymbol{\Sigma}_n$ , based on the previous predictive result.

We apply this framework to ten tech stocks of the US market between January first, 2021, and September ninth, 2022. In particular, we use information from Yahoo Finance for Apple (AAPL), Netflix (NFLX), Amazon (AMZN), Microsoft (MSFT), Google (GOOG), Meta (META), Tesla (TSLA), NVIDIA Corporation (NVDA), Intel (INTC), and PayPal (PYPL).

We use non-informative hyperparameters,  $\boldsymbol{\mu}_0 = \mathbf{0}$ ,  $\beta_0 = 1$ ,  $\boldsymbol{\Psi}_0 = 100\mathbf{I}$  and  $\alpha_0 = p + 2$ , where  $p = 10$  is the number of stocks.

### R code. Optimal tangency portfolio: Tech shares

```

1 library(quantmod); library(xts); library(ggplot2)
2 library(gridExtra); library(purrr); library(dplyr)
3 # Define date range
4 start_date <- as.Date("2021-01-01")
5 end_date <- as.Date("2022-09-30")
6 dates <- seq(start_date, end_date, by = "day")
7 # Tickers of interest
8 tickers <- c("AAPL", "NFLX", "AMZN", "GOOG", "INTC", "META",
9             "MSFT", "TSLA", "NVDA", "PYPL")
10 p <- length(tickers)
11 # Download adjusted closing prices
12 getSymbols(tickers, from = start_date, to = end_date, auto.assign = TRUE)
13 prices <- map(tickers, ~ Ad(get(.x))) %>%
14   reduce(merge) %>%
15   `colnames<-`((tickers)) %>%
16   as.data.frame()
17 # Calculate daily log returns
18 returns <- apply(prices, 2, function(x) diff(log(x))) %>%
19   as.data.frame()
20 # Download 10-year Treasury yield from FRED
21 t10yr_xts <- getSymbols("^TNX", src = "yahoo", from = start_date,
22                         to = end_date, auto.assign = FALSE)
23 # Annual yield in decimal
24 t10yr_annual <- Cl(t10yr_xts) / 1000 # e.g. 45 -> 0.045
25 # Convert to daily rate (using 365 days, consistent with
# your code)
26 t10yr_daily <- (1 + t10yr_annual)^(1 / 365) - 1
27 # Align risk-free series with stock returns
28 t10yr_daily <- t10yr_daily[rownames(returns), ]
29 # Compute excess returns
30 excess_returns <- as.matrix(returns) -
31   kronecker(t(rep(1, p)), as.matrix(t10yr_daily))
32 # Convert to data frame with dates
33 df <- as.data.frame(excess_returns)
34 df$Date <- as.Date(rownames(df))
35 df$Month <- months(df$Date)
36 df$Year <- format(df$Date, "%y")
37 # Aggregate monthly means
38 monthly_means <- map(1:p, function(i) {
39   aggregate(df[[i]] ~ Month + Year, data = df, FUN = mean)
40 })
41 # Extract values into matrix
42 data_excess <- matrix(0, nrow = nrow(monthly_means[[1]]),
43                         ncol = p)
44 for (i in 1:p) {
45   data_excess[, i] <- as.numeric(monthly_means[[i]][, 3])
46 }
47 colnames(data_excess) <- tickers

```

**R code. Optimal tangency portfolio: Tech shares**

```

1 # Hyperparameters #
2 N <- dim(DataExcRtn)[1]
3 mu0 <- rep(0, p)
4 beta0 <- 1
5 Psi0 <- 100 * diag(p)
6 alpha0 <- p + 2
7 # Posterior parameters #
8 alphan <- N + alpha0
9 vn <- alphan + 1 - p
10 muhat <- colMeans(DataExcRtn)
11 mun <- N/(N + beta0) * muhat + beta0/(N + beta0) * mu0
12 S <- t(DataExcRtn - rep(1, N) %*% t(muhat)) %*% (DataExcRtn -
   rep(1, N) %*% t(muhat))
13 Psin <- Psi0 + S + N*beta0/(N + beta0)*(muhat - mu0) %*% t(
   muhat - mu0)
14 betan <- N + beta0
15 Sigman <- Psin/((alphan + 1 - p)*betan)
16 Covarn <- (Sigman * (1 + betan)) * vn / (vn - 2)
17 Covari <- solve(Covarn)
18 OptShare <- t(Covari %*% mun/as.numeric((t(rep(1, p)) %*% Covari
   %*% mun)))
19 colnames(OptShare) <- tickers
20 OptShare

```

We find that the optimal tangency portfolio is composed by 24.9%, 10.6%, 17.6%, 23.6%, 2.9% and 31.1% weights of Netflix, Amazon, Intel, Meta, NVIDIA and PayPal, and -2.3%, -3.8%, -2.8% and -1.8% weights of Apple, Google, Microsoft and Tesla. A negative weight means being short in financial jargon, that is, borrowing a stock to sell it.

### 3.3 Linear regression: The conjugate normal-normal/inverse gamma model

In this setting, we analyze the conjugate normal-normal/inverse gamma model, which is a cornerstone in econometrics. In this model, the dependent variable  $y_i$  is related to a set of regressors  $\mathbf{x}_i = [x_{i1} \ x_{i2} \ \dots \ x_{iK}]^\top$  in a linear way, that is:

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_K x_{iK} + \mu_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \mu_i,$$

where  $\boldsymbol{\beta} = [\beta_1 \ \beta_2 \ \dots \ \beta_K]^\top$  and  $\mu_i \stackrel{iid}{\sim} N(0, \sigma^2)$  is a stochastic error such that  $\mathbb{E}[\mu_i | \mathbf{x}_i] = 0$ .

Defining the vectors and matrices:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1K} \\ x_{21} & x_{22} & \dots & x_{2K} \\ \vdots & \vdots & \vdots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NK} \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_N \end{bmatrix},$$

we can write the model in matrix form as:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu},$$

where  $\boldsymbol{\mu} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ . This implies that:

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}).^5$$

Thus, the likelihood function is:

$$\begin{aligned} p(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbf{X}) &= (2\pi\sigma^2)^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\} \\ &\propto (\sigma^2)^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}. \end{aligned}$$

The conjugate priors for the parameters are

$$\begin{aligned} \boldsymbol{\beta} | \sigma^2 &\sim N(\boldsymbol{\beta}_0, \sigma^2 \mathbf{B}_0), \\ \sigma^2 &\sim IG(\alpha_0/2, \delta_0/2). \end{aligned}$$

Then, the posterior distribution is

$$\begin{aligned} \pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &\propto (\sigma^2)^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\} \\ &\times (\sigma^2)^{-\frac{K}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \right\} \\ &\times \frac{(\delta_0/2)^{\alpha_0/2}}{\Gamma(\alpha_0/2)} \left( \frac{1}{\sigma^2} \right)^{\alpha_0/2+1} \exp \left\{ -\frac{\delta_0}{2\sigma^2} \right\} \\ &\propto \exp \left\{ -\frac{1}{2\sigma^2} [\boldsymbol{\beta}^\top (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X}) \boldsymbol{\beta} - 2\boldsymbol{\beta}^\top (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}})] \right\} \\ &\times (\sigma^2)^{-\frac{K}{2}} \left( \frac{1}{\sigma^2} \right)^{(\alpha_0+N)/2+1} \exp \left\{ -\frac{\delta_0 + \mathbf{y}^\top \mathbf{y} + \boldsymbol{\beta}_0^\top \mathbf{B}_0^{-1} \boldsymbol{\beta}_0}{2\sigma^2} \right\}, \end{aligned}$$

---

<sup>5</sup>In regression analysis, to simplify notation, we depart from the conventional statistical notation, which defines lowercase letters as realizations of random variables, typically denoted by uppercase letters. We hope it is clear from the context when we refer to random vectors and matrices, and their realizations. Thus, we use bold lowercase letters for vectors and bold uppercase letters for matrices. This applies to the rest of the book.

where  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  is the maximum likelihood estimator.

Adding and subtracting  $\beta_n^\top \mathbf{B}_n^{-1} \beta_n$  to complete the square, where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}$  and  $\beta_n = \mathbf{B}_n(\mathbf{B}_0^{-1} \beta_0 + \mathbf{X}^\top \mathbf{X} \hat{\beta})$ ,

$$\pi(\beta, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto \underbrace{(\sigma^2)^{-\frac{K}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\beta - \beta_n)^\top \mathbf{B}_n^{-1} (\beta - \beta_n) \right\}}_1 \\ \times \underbrace{(\sigma^2)^{-\left(\frac{\alpha_n}{2} + 1\right)} \exp \left\{ -\frac{\delta_n}{2\sigma^2} \right\}}_2.$$

The first expression is the kernel of a normal density function,  $\beta | \sigma^2, \mathbf{y}, \mathbf{X} \sim N(\beta_n, \sigma^2 \mathbf{B}_n)$ . The second expression is the kernel of a inverse gamma density,  $\sigma^2 | \mathbf{y}, \mathbf{X} \sim IG(\alpha_n/2, \delta_n/2)$ , where  $\alpha_n = \alpha_0 + N$  and  $\delta_n = \delta_0 + \mathbf{y}^\top \mathbf{y} + \beta_0^\top \mathbf{B}_0^{-1} \beta_0 - \beta_n^\top \mathbf{B}_n^{-1} \beta_n$ .

Taking into account that

$$\begin{aligned} \beta_n &= (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{B}_0^{-1} \beta_0 + \mathbf{X}^\top \mathbf{X} \hat{\beta}) \\ &= (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{B}_0^{-1} \beta_0 + (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \hat{\beta}, \end{aligned}$$

where  $(\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{B}_0^{-1} = \mathbf{I}_K - (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}$  [415]. Setting  $\mathbf{W} = (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}$  we have  $\beta_n = (\mathbf{I}_K - \mathbf{W}) \beta_0 + \mathbf{W} \hat{\beta}$ , that is, the posterior mean of  $\beta$  is a weighted average between the sample and prior information, where the weights depend on the precision of each piece of information. Observe that when the prior covariance matrix is highly vague (non-informative), such that  $\mathbf{B}_0^{-1} \rightarrow \mathbf{0}_K$ , we obtain  $\mathbf{W} \rightarrow \mathbf{I}_K$ , such that  $\beta_n \rightarrow \hat{\beta}$ , that is, the posterior mean location parameter converges to the maximum likelihood estimator.

In addition, we know that the posterior conditional covariance matrix of the location parameters  $\sigma^2 (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} - \sigma^2 ((\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{B}_0 + (\mathbf{X}^\top \mathbf{X})^{-1})^{-1} (\mathbf{X}^\top \mathbf{X})^{-1})$  is positive semi-definite.<sup>6</sup> Given that  $\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$  is the covariance matrix of the maximum likelihood estimator, we observe that prior information reduces estimation uncertainty.

Another way to see this model is by considering that both  $\mathbf{y}$  and  $\beta$  are treated as random variables under the Bayesian framework. Thus, we can express the joint distribution of these two vectors as follows:

$$\begin{bmatrix} \beta \\ \mathbf{y} \end{bmatrix} \sim N \left[ \begin{pmatrix} \beta_0 \\ \mathbf{X} \beta_0 \end{pmatrix}, \sigma^2 \begin{pmatrix} \mathbf{B}_0 & \mathbf{B}_0 \mathbf{X}^\top \\ \mathbf{X} \mathbf{B}_0^\top & \mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N \end{pmatrix} \right],$$

where we use that  $Cov[\beta, \mathbf{y}] = \mathbb{E}[\beta \mathbf{y}^\top] - \mathbb{E}[\beta] \mathbb{E}[\mathbf{y}^\top] = \mathbb{E}[\beta (\mathbf{X} \beta + \mu)^\top] -$

<sup>6</sup>A particular case of the Woodbury matrix identity,  $(\mathbf{A} + \mathbf{U} \mathbf{C} \mathbf{V})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{C}^{-1} + \mathbf{V} \mathbf{A}^{-1} \mathbf{U})^{-1} \mathbf{V} \mathbf{A}^{-1}$ .

$$\mathbb{E}[\beta]\mathbb{E}[\mathbf{y}^\top] = [Var[\beta] + \mathbb{E}[\beta]\mathbb{E}[\beta^\top]]\mathbf{X}^\top - \mathbb{E}[\beta]\mathbb{E}[\mathbf{y}^\top] = \sigma^2 \mathbf{B}_0 \mathbf{X}^\top + \beta_0 \beta_0^\top \mathbf{X}^\top - \beta_0 \beta_0^\top \mathbf{X}^\top = \sigma^2 \mathbf{B}_0 \mathbf{X}^\top.$$

Then, we can obtain the conditional distribution of  $\beta \mid \mathbf{y}$  using the properties of the multivariate normal distribution. This distribution is normal with mean equal to

$$\beta_0 + \mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1} (\mathbf{y} - \mathbf{X} \beta_0), \quad (3.5)$$

and covariance matrix

$$\sigma^2 \left( \mathbf{B}_0 - \mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1} \mathbf{X} \mathbf{B}_0^\top \right).$$

Observe that in this representation, the posterior mean is equal to the prior mean plus a correction term that takes into account the deviation between the observations and the prior expected value ( $\mathbf{X} \beta_0$ ). The weight of this correction is given by the matrix  $\mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1}$ .

This form of expressing the posterior distribution is relevant for gaining some intuition on Bayesian inference in time series models within the *Gaussian linear state-space representation* in Chapter 8, also known as the Kalman filter in time series literature.

We can show that both conditional posterior distributions are the same. In particular, the posterior mean in Equation 3.5 is  $[\mathbf{I}_K - \mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1} \mathbf{X}] \beta_0 + \mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1} \mathbf{y}$ , where

$$\begin{aligned} \mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1} &= \mathbf{B}_0 \mathbf{X}^\top [\mathbf{I}_N - \mathbf{I}_N \mathbf{X} (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{I}_N \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{I}_N] \\ &= \mathbf{B}_0 [\mathbf{I}_K - \mathbf{X}^\top \mathbf{I}_N \mathbf{X} (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{I}_N \mathbf{X})^{-1}] \mathbf{X}^\top \\ &= \mathbf{B}_0 [\mathbf{I}_K - [\mathbf{I}_K - \mathbf{B}_0^{-1} (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{I}_N \mathbf{X})^{-1}]] \mathbf{X}^\top \\ &= (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top, \end{aligned}$$

where the first equality uses the Woodbury matrix identity (matrix inversion lemma), and the third equality uses  $\mathbf{D}(\mathbf{D} + \mathbf{E})^{-1} = \mathbf{I} - \mathbf{E}(\mathbf{D} + \mathbf{E})^{-1}$ .

Thus,  $[\mathbf{I}_K - \mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1} \mathbf{X}] \beta_0 + \mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1} \mathbf{y} = [\mathbf{I}_K - (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}] \beta_0 + (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = [\mathbf{I}_K - (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}] \beta_0 + (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \hat{\beta}$ . Again, we see that the posterior mean is a weighted average between the prior mean, and the maximum likelihood estimator.

The equality of variances of both approaches is as follows:

$$\begin{aligned} Var[\beta \mid \mathbf{y}] &= \sigma^2 (\mathbf{B}_0 - \mathbf{B}_0 \mathbf{X}^\top (\mathbf{X} \mathbf{B}_0 \mathbf{X}^\top + \mathbf{I}_N)^{-1} \mathbf{X} \mathbf{B}_0) \\ &= \sigma^2 (\mathbf{B}_0 - \mathbf{B}_0 \mathbf{X}^\top (\mathbf{I}_N - \mathbf{I}_N \mathbf{X} (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{I}_N \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{I}_N) \mathbf{X} \mathbf{B}_0) \\ &= \sigma^2 (\mathbf{B}_0 - \mathbf{B}_0 \mathbf{X}^\top \mathbf{X} \mathbf{B}_0 + \mathbf{B}_0 \mathbf{X}^\top \mathbf{X} (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{B}_0) \\ &= \sigma^2 (\mathbf{B}_0 - \mathbf{B}_0 \mathbf{X}^\top \mathbf{X} \mathbf{B}_0 + \mathbf{B}_0 \mathbf{X}^\top \mathbf{X} [\mathbf{I}_K - (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{B}_0^{-1}] \mathbf{B}_0) \\ &= \sigma^2 (\mathbf{B}_0 - \mathbf{B}_0 \mathbf{X}^\top \mathbf{X} (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}) \end{aligned}$$

$$\begin{aligned}
&= \sigma^2(\mathbf{B}_0[\mathbf{I}_K - \mathbf{X}^\top \mathbf{X}(\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}]) \\
&= \sigma^2(\mathbf{B}_0[\mathbf{I}_K - (\mathbf{I}_K - \mathbf{B}_0^{-1}(\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1})]) \\
&= \sigma^2(\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1},
\end{aligned}$$

where the second equality uses the Woodbury matrix identity, the fourth equality uses  $(\mathbf{D} + \mathbf{E})^{-1}\mathbf{D} = \mathbf{I} - (\mathbf{D} + \mathbf{E})^{-1}\mathbf{E}$ , and the seventh equality uses  $\mathbf{D}(\mathbf{D} + \mathbf{E})^{-1} = \mathbf{I} - \mathbf{E}(\mathbf{D} + \mathbf{E})^{-1}$ .

Now, we calculate the posterior marginal distribution of  $\boldsymbol{\beta}$  following the standard approach,

$$\begin{aligned}
\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) &= \int_0^\infty \pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) d\sigma^2 \\
&= \int_0^\infty \left( \frac{1}{\sigma^2} \right)^{\frac{\alpha_n+K}{2}+1} \exp \left\{ -\frac{s}{2\sigma^2} \right\} d\sigma^2,
\end{aligned}$$

where  $s = \delta_n + (\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \mathbf{B}_n^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_n)$ . Then we can write

$$\begin{aligned}
\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) &= \int_0^\infty \left( \frac{1}{\sigma^2} \right)^{\frac{\alpha_n+K}{2}+1} \exp \left\{ -\frac{s}{2\sigma^2} \right\} d\sigma^2 \\
&= \frac{\Gamma((\alpha_n+K)/2)}{(s/2)^{(\alpha_n+K)/2}} \int_0^\infty \frac{(s/2)^{(\alpha_n+K)/2}}{\Gamma((\alpha_n+K)/2)} (\sigma^2)^{-(\alpha_n+K)/2-1} \exp \left\{ -\frac{s}{2\sigma^2} \right\} d\sigma^2.
\end{aligned}$$

The right term is the integral of the probability density function of an inverse gamma distribution with parameters  $\nu = (\alpha_n + K)/2$  and  $\tau = s/2$ . Since we are integrating over the whole support of  $\sigma^2$ , the integral is equal to 1, and therefore

$$\begin{aligned}
\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) &= \frac{\Gamma((\alpha_n+K)/2)}{(s/2)^{(\alpha_n+K)/2}} \\
&\propto s^{-(\alpha_n+K)/2} \\
&= [\delta_n + (\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \mathbf{B}_n^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_n)]^{-(\alpha_n+K)/2} \\
&= \left[ 1 + \frac{(\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \left( \frac{\delta_n}{\alpha_n} \mathbf{B}_n \right)^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_n)}{\alpha_n} \right]^{-(\alpha_n+K)/2} (\delta_n)^{-(\alpha_n+K)/2} \\
&\propto \left[ 1 + \frac{(\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \mathbf{H}_n^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_n)}{\alpha_n} \right]^{-(\alpha_n+K)/2},
\end{aligned}$$

where  $\mathbf{H}_n = \frac{\delta_n}{\alpha_n} \mathbf{B}_n$ . This last expression is a multivariate t distribution, that is,  $\boldsymbol{\beta} | \mathbf{y}, \mathbf{X} \sim t_K(\alpha_n, \boldsymbol{\beta}_n, \mathbf{H}_n)$ .

Observe that as we have incorporated the uncertainty of the variance, the posterior for  $\boldsymbol{\beta}$  changes from a normal to a t distribution, which has heavier tails, indicating more uncertainty.

The marginal likelihood of this model is

$$p(\mathbf{y}) = \int_0^\infty \int_{R^K} \pi(\boldsymbol{\beta} | \sigma^2, \mathbf{B}_0, \boldsymbol{\beta}_0) \pi(\sigma^2 | \alpha_0/2, \delta_0/2) p(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbf{X}) d\sigma^2 d\boldsymbol{\beta}.$$

Taking into account that  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_0) = (\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \mathbf{B}_n^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_n) + m$ , where  $m = \mathbf{y}^\top \mathbf{y} + \boldsymbol{\beta}_0^\top \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 - \boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n$ , we have that

$$\begin{aligned} p(\mathbf{y}) &= \int_0^\infty \int_{R^K} \pi(\boldsymbol{\beta} | \sigma^2) \pi(\sigma^2) p(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbf{X}) d\sigma^2 d\boldsymbol{\beta} \\ &= \int_0^\infty \pi(\sigma^2) \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2}m\right\} \frac{1}{(2\pi\sigma^2)^{K/2} |\mathbf{B}_0|^{1/2}} \\ &\quad \times \int_{R^K} \exp\left\{-\frac{1}{2\sigma^2}(\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \mathbf{B}_n^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_n)\right\} d\sigma^2 d\boldsymbol{\beta} \\ &= \int_0^\infty \pi(\sigma^2) \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2}m\right\} \frac{|\mathbf{B}_n|^{1/2}}{|\mathbf{B}_0|^{1/2}} d\sigma^2 \\ &= \int_0^\infty \frac{(\delta_0/2)^{\alpha_0/2}}{\Gamma(\alpha_0/2)} \left(\frac{1}{\sigma^2}\right)^{\alpha_0/2+1} \exp\left\{\left(-\frac{\delta_0}{2\sigma^2}\right)\right\} \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2}m\right\} \\ &\quad \times \frac{|\mathbf{B}_n|^{1/2}}{|\mathbf{B}_0|^{1/2}} d\sigma^2 \\ &= \frac{1}{(2\pi)^{N/2}} \frac{(\delta_0/2)^{\alpha_0/2}}{\Gamma(\alpha_0/2)} \frac{|\mathbf{B}_n|^{1/2}}{|\mathbf{B}_0|^{1/2}} \int_0^\infty \left(\frac{1}{\sigma^2}\right)^{\frac{\alpha_0+N}{2}+1} \exp\left\{\left(-\frac{\delta_0+m}{2\sigma^2}\right)\right\} d\sigma^2 \\ &= \frac{1}{\pi^{N/2}} \frac{\delta_0^{\alpha_0/2}}{\delta_n^{\alpha_n/2}} \frac{|\mathbf{B}_n|^{1/2}}{|\mathbf{B}_0|^{1/2}} \frac{\Gamma(\alpha_n/2)}{\Gamma(\alpha_0/2)}. \end{aligned}$$

We can show that  $\delta_n = \delta_0 + \mathbf{y}^\top \mathbf{y} + \boldsymbol{\beta}_0^\top \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 - \boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n = \delta_0 + (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top ((\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{B}_0)^{-1}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$  (see Exercise 7). Therefore, if we want to compare two models under this setting, the Bayes factor is

$$\begin{aligned} BF_{12} &= \frac{p(\mathbf{y} | \mathcal{M}_1)}{p(\mathbf{y} | \mathcal{M}_2)} \\ &= \frac{\frac{\delta_{10}^{\alpha_{10}/2}}{\delta_{1n}^{\alpha_{1n}/2}} \frac{|\mathbf{B}_{1n}|^{1/2}}{|\mathbf{B}_{10}|^{1/2}} \frac{\Gamma(\alpha_{1n}/2)}{\Gamma(\alpha_{10}/2)}}{\frac{\delta_{20}^{\alpha_{20}/2}}{\delta_{2n}^{\alpha_{2n}/2}} \frac{|\mathbf{B}_{2n}|^{1/2}}{|\mathbf{B}_{20}|^{1/2}} \frac{\Gamma(\alpha_{2n}/2)}{\Gamma(\alpha_{20}/2)}}, \end{aligned}$$

where subscripts 1 and 2 refer to each model.

Observe that, *ceteris paribus*, the model with better fit, coherence between sample and prior information regarding location parameters, higher prior to posterior precision, and fewer parameters is favored by the Bayes factor. The

Bayes factor rewards model fit, as the sum of squared errors appears in  $\delta_n$ ; the better the fit (i.e., the lower the sum of squared errors), the better the Bayes factor. In addition, a weighted distance between sample and prior location parameters also appears in  $\delta_n$ . The greater this distance, the worse the model support. The ratio of determinants between posterior and prior covariance matrices is also present; the higher this ratio, the better the Bayes factor supports a model due to information gains.

To see the effect of a model's parsimony, let's consider the common situation in applications where  $\mathbf{B}_{j0} = c\mathbf{I}_{K_j}$ , then  $|\mathbf{B}_{j0}| = c^{K_j}$ . Hence,

$$\left(\frac{|\mathbf{B}_{20}|}{|\mathbf{B}_{10}|}\right)^{1/2} = \left(\frac{c^{K_2/2}}{c^{K_1/2}}\right),$$

if  $\frac{K_2}{K_1} > 1$  and  $c \rightarrow \infty$  (the latter implying a non-informative prior), then  $BF_{12} \rightarrow \infty$ . This means infinite evidence supporting the parsimonious model, no matter what the sample information says.

Comparing models having the same number of regressors ( $K_1 = K_2$ ) is not a safe ground, as  $|\mathbf{B}_0|$  depends on the measurement units of the regressors. Conclusions regarding model selection depend on this, which is not a nice property. This prevents using non-informative priors when performing model selection in the Bayesian framework. However, this is not the case when  $\alpha_0 \rightarrow 0$  and  $\delta_0 \rightarrow 0$ , which implies a non-informative prior for the variance parameter.<sup>7</sup>

We observe that  $\Gamma(\alpha_{j0})$  cancels out, as  $\alpha_0 \rightarrow 0$  implies  $\alpha_{jn} \rightarrow N$ , and

$$\delta_{jn} \rightarrow (\mathbf{y} - \mathbf{X}_j \hat{\beta}_j)^\top (\mathbf{y} - \mathbf{X}_j \hat{\beta}_j) + (\hat{\beta}_j - \boldsymbol{\beta}_{j0})^\top ((\mathbf{X}_j^\top \mathbf{X}_j)^{-1} + \mathbf{B}_{j0})^{-1} (\hat{\beta}_j - \boldsymbol{\beta}_{j0}),$$

when  $\delta_0 \rightarrow 0$ , therefore, there is no effect. This is due to  $\sigma^2$  being a common parameter in both models.

In general, we can use non-informative priors for common parameters across all models, but we cannot use non-informative priors for non-common parameters when performing model selection using the Bayes factor. This issue raises the question of how to set informative priors. On one hand, we have those who advocate for *subjective* priors [362, 112, 396, 269]; on the other hand, those who prefer *objective* priors [32, 256, 223, 38].

Regarding the former, eliciting *subjective* priors, i.e., "formulating a person's knowledge and beliefs about one or more uncertain quantities into a (joint) probability distribution for those quantities" [155], is a very difficult task due to human beings' heuristics and biases associated with representativeness, information availability, conservatism, overconfidence, and anchoring-and-adjustment issues [438]. However, there have been good efforts using predictive and structural elicitation procedures [233, 234].

Regarding the latter, there are *reference priors* that are designed to have minimal impact on the posterior distribution and to be invariant to different

---

<sup>7</sup>See [161] for advice against this common practice.

parametrizations of the model [42, Chap. 5]. A remarkable example of *reference priors* is the *Jeffreys' prior* [224], which originated from the critique of *non-informative priors* that were not invariant to transformations of the parameter space. In particular, the *Jeffreys' prior* is given by:

$$\pi(\boldsymbol{\theta}) \propto |I(\boldsymbol{\theta})|^{1/2},$$

where  $I(\boldsymbol{\theta}) = \mathbb{E} \left( -\frac{\partial^2 \log p(\mathbf{y}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right)$ , i.e.,  $I(\boldsymbol{\theta})$  is the Fisher information matrix. However, the *Jeffreys' prior* is often improper, meaning it does not work well for model selection.

Thus, a standard *objective* approach is to use *intrinsic priors* [40], where a *minimal training* dataset is used with a *reference prior* to obtain a proper posterior distribution. This proper distribution is then used as a prior, and the standard Bayesian procedures are followed using the remaining dataset. In this way, we end up with meaningful Bayes factors for model selection.

Regardless of using a *subjective* or *objective* approach to define a prior distribution, it is always a good idea to assess the sensitivity of the posterior results to the prior assumptions. This is commonly done using local or pointwise assessments, such as partial derivatives [178, 221, 186] or, more often, in terms of multiple evaluations (*scenario analysis*) [368, 243, 8]. Recently, [220] extend these approaches to perform sensitivity analysis in high-dimensional hyperparameter settings.

Returning to the linear model, the posterior predictive is equal to

$$\begin{aligned} \pi(\mathbf{y}_0 \mid \mathbf{y}) &= \int_0^\infty \int_{R^K} p(\mathbf{Y}_0 \mid \boldsymbol{\beta}, \sigma^2, \mathbf{y}) \pi(\boldsymbol{\beta} \mid \sigma^2, \mathbf{y}) \pi(\sigma^2 \mid \mathbf{y}) d\boldsymbol{\beta} d\sigma^2 \\ &= \int_0^\infty \int_{R^K} p(\mathbf{Y}_0 \mid \boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta} \mid \sigma^2, \mathbf{y}) \pi(\sigma^2 \mid \mathbf{y}) d\boldsymbol{\beta} d\sigma^2, \end{aligned}$$

where we take into account independence between  $\mathbf{y}_0$  and  $\mathbf{y}$ . Given  $\mathbf{X}_0$ , which is the  $N_0 \times K$  matrix of regressors associated with  $\mathbf{y}_0$ , Then,

$$\begin{aligned} \pi(\mathbf{y}_0 \mid \mathbf{y}) &= \int_0^\infty \int_{R^K} \left\{ (2\pi\sigma^2)^{-\frac{N_0}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y}_0 - \mathbf{X}_0 \boldsymbol{\beta})^\top (\mathbf{y}_0 - \mathbf{X}_0 \boldsymbol{\beta}) \right\} \right. \\ &\quad \times (2\pi\sigma^2)^{-\frac{K}{2}} |\mathbf{B}_n|^{-1/2} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \mathbf{B}_n^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_n) \right\} \\ &\quad \left. \times \frac{(\delta_n/2)^{\alpha_n/2}}{\Gamma(\alpha_n/2)} \left( \frac{1}{\sigma^2} \right)^{\alpha_n/2+1} \exp \left\{ -\frac{\delta_n}{2\sigma^2} \right\} \right\} d\boldsymbol{\beta} d\sigma^2. \end{aligned}$$

Setting  $\mathbf{M} = (\mathbf{X}_0^\top \mathbf{X}_0 + \mathbf{B}_n^{-1})$  and  $\boldsymbol{\beta}_* = \mathbf{M}^{-1}(\mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \mathbf{X}_0^\top \mathbf{y}_0)$ , we have  $(\mathbf{y}_0 - \mathbf{X}_0 \boldsymbol{\beta})^\top (\mathbf{y}_0 - \mathbf{X}_0 \boldsymbol{\beta}) + (\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \mathbf{B}_n^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_n) = (\boldsymbol{\beta} - \boldsymbol{\beta}_*)^\top \mathbf{M} (\boldsymbol{\beta} - \boldsymbol{\beta}_*) + \boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \mathbf{y}_0^\top \mathbf{y}_0 - \boldsymbol{\beta}_*^\top \mathbf{M} \boldsymbol{\beta}_*$ . Thus,

$$\pi(\mathbf{y}_0 \mid \mathbf{y}) \propto \int_0^\infty \left\{ \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \mathbf{y}_0^\top \mathbf{y}_0 - \boldsymbol{\beta}_*^\top \mathbf{M} \boldsymbol{\beta}_* + \delta_n) \right\} \right\}$$

$$\begin{aligned} & \times \left( \frac{1}{\sigma^2} \right)^{-\frac{K+N_0+\alpha_n}{2}+1} \\ & \times \int_{R^K} \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta} - \boldsymbol{\beta}_*)^\top \mathbf{M} (\boldsymbol{\beta} - \boldsymbol{\beta}_*) \right\} d\boldsymbol{\beta} \} d\sigma^2, \end{aligned}$$

where the term in the second integral is the kernel of a multivariate normal density with mean  $\boldsymbol{\beta}_*$  and covariance matrix  $\sigma^2 \mathbf{M}^{-1}$ . Then,

$$\begin{aligned} \pi(\mathbf{y}_0 | \mathbf{y}) & \propto \int_0^\infty \left( \frac{1}{\sigma^2} \right)^{\frac{N_0+\alpha_n}{2}+1} \\ & \times \exp \left\{ -\frac{1}{2\sigma^2} (\boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \mathbf{y}_0^\top \mathbf{y}_0 - \boldsymbol{\beta}_*^\top \mathbf{M} \boldsymbol{\beta}_* + \delta_n) \right\} d\sigma^2, \end{aligned}$$

which is the kernel of an inverse gamma density. Thus,

$$\pi(\mathbf{y}_0 | \mathbf{y}) \propto \left[ \frac{\boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \mathbf{y}_0^\top \mathbf{y}_0 - \boldsymbol{\beta}_*^\top \mathbf{M} \boldsymbol{\beta}_* + \delta_n}{2} \right]^{-\frac{\alpha_n+N_0}{2}}.$$

Setting  $\mathbf{C}^{-1} = \mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{B}_n \mathbf{X}_0^\top$  such that  $\mathbf{C} = \mathbf{I}_{N_0} - \mathbf{X}_0 (\mathbf{B}_n^{-1} + \mathbf{X}_0^\top \mathbf{X}_0)^{-1} \mathbf{X}_0^\top = \mathbf{I}_{N_0} - \mathbf{X}_0 \mathbf{M}^{-1} \mathbf{X}_0^\top$ <sup>8</sup>, and  $\boldsymbol{\beta}_{**} = \mathbf{C}^{-1} \mathbf{X}_0 \mathbf{M}^{-1} \mathbf{B}_n^{-1} \boldsymbol{\beta}_n$ , then

$$\begin{aligned} \boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \mathbf{y}_0^\top \mathbf{y}_0 - \boldsymbol{\beta}_*^\top \mathbf{M} \boldsymbol{\beta}_* &= \boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \mathbf{y}_0^\top \mathbf{y}_0 - \\ & (\boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} + \mathbf{y}_0^\top \mathbf{X}_0) \mathbf{M}^{-1} (\mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \mathbf{X}_0^\top \mathbf{y}_0) \\ &= \boldsymbol{\beta}_n^\top (\mathbf{B}_n^{-1} - \mathbf{B}_n^{-1} \mathbf{M}^{-1} \mathbf{B}_n^{-1}) \boldsymbol{\beta}_n + \mathbf{y}_0^\top \mathbf{C} \mathbf{y}_0 \\ &- 2\mathbf{y}_0^\top \mathbf{C} \mathbf{C}^{-1} \mathbf{X}_0 \mathbf{M}^{-1} \mathbf{B}_n^{-1} \boldsymbol{\beta}_n + \boldsymbol{\beta}_{**}^\top \mathbf{C} \boldsymbol{\beta}_{**} - \boldsymbol{\beta}_{**}^\top \mathbf{C} \boldsymbol{\beta}_{**} \\ &= \boldsymbol{\beta}_n^\top (\mathbf{B}_n^{-1} - \mathbf{B}_n^{-1} \mathbf{M}^{-1} \mathbf{B}_n^{-1}) \boldsymbol{\beta}_n + \\ & + (\mathbf{y}_0 - \boldsymbol{\beta}_{**})^\top \mathbf{C} (\mathbf{y}_0 - \boldsymbol{\beta}_{**}) - \boldsymbol{\beta}_{**}^\top \mathbf{C} \boldsymbol{\beta}_{**}, \end{aligned}$$

where  $\boldsymbol{\beta}_n^\top (\mathbf{B}_n^{-1} - \mathbf{B}_n^{-1} \mathbf{M}^{-1} \mathbf{B}_n^{-1}) \boldsymbol{\beta}_n = \boldsymbol{\beta}_{**}^\top \mathbf{C} \boldsymbol{\beta}_{**}$  and  $\boldsymbol{\beta}_{**} = \mathbf{X}_0 \boldsymbol{\beta}_n$  (see Exercise 8).

Then,

$$\begin{aligned} \pi(\mathbf{y}_0 | \mathbf{y}) & \propto \left[ \frac{(\mathbf{y}_0 - \mathbf{X}_0 \boldsymbol{\beta}_n)^\top \mathbf{C} (\mathbf{y}_0 - \mathbf{X}_0 \boldsymbol{\beta}_n) + \delta_n}{2} \right]^{-\frac{\alpha_n+N_0}{2}} \\ & \propto \left[ \frac{(\mathbf{y}_0 - \mathbf{X}_0 \boldsymbol{\beta}_n)^\top \left( \frac{\mathbf{C} \boldsymbol{\alpha}_n}{\delta_n} \right) (\mathbf{y}_0 - \mathbf{X}_0 \boldsymbol{\beta}_n)}{\alpha_n} + 1 \right]^{-\frac{\alpha_n+N_0}{2}}. \end{aligned}$$

<sup>8</sup>Using  $(\mathbf{A} + \mathbf{B} \mathbf{D} \mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{D}^{-1} + \mathbf{C} \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{C} \mathbf{A}^{-1}$

The posterior predictive is a multivariate t distribution,  $\mathbf{y}_0 \mid \mathbf{y} \sim t\left(\mathbf{X}_0\boldsymbol{\beta}_n, \frac{\delta_n(\mathbf{I}_{N_0} + \mathbf{X}_0\mathbf{B}_n\mathbf{X}_0^\top)}{\alpha_n}, \alpha_n\right)$ .

### Example: Demand of electricity

We study in this example the determinants of the monthly demand for electricity by Colombian households. The data consists of information from 2103 households, including the following variables: the average price (USD/kWh), indicators of the socioeconomic conditions of the neighborhood where the household is located (with `IndSocio1` being the lowest and `IndSocio3` being the highest), an indicator for whether the household is located in a municipality that is above 1000 meters above sea level, the number of rooms in the house, the number of members in the household, the presence of children in the household (where 1 indicates yes), and the monthly income (USD). The specification is as follows:

$$\begin{aligned}\log(\text{Electricity}_i) = & \beta_1 \log(\text{price}_i) + \beta_2 \text{IndSocio1}_i + \beta_3 \text{IndSocio2}_i + \beta_4 \text{Altitude}_i \\ & + \beta_5 \text{Nrooms}_i + \beta_6 \text{HouseholdMem}_i + \beta_7 \text{Children}_i \\ & + \beta_8 \log(\text{Income}_i) + \beta_9 + \mu_i.\end{aligned}$$

We use a non-informative vague prior setting such that  $\alpha_0 = \delta_0 = 0.001$ ,  $\boldsymbol{\beta}_0 = \mathbf{0}$  and  $\mathbf{B}_0 = c_0 \mathbf{I}_K$ , where  $c_0 = 1000$  and  $K$  is the number of regressors.

The results from the **R** code (see below) indicate that the posterior mean of the own-price elasticity of electricity demand is -1.09, and the 95% symmetric credible interval is (-1.47, -0.71). Households in neighborhoods with low socioeconomic conditions and those located in municipalities situated more than 1000 meters above sea level consume less electricity, with reductions of 32.7% and 19.7% on average, respectively. An additional room leads to an 8.7% increase in electricity consumption, and each additional household member increases consumption by 5.9% on average. The mean estimate for income elasticity is 0.074, meaning that a 10% increase in income results in a 0.74% increase in electricity demand.

We want to check the results of the Bayes factor comparing the previous specification (model 1) with other specification without considering the price of electricity (model 2), that is,

$$\begin{aligned}\log(\text{Electricity}_i) = & \beta_1 \text{IndSocio1}_i + \beta_2 \text{IndSocio2}_i + \beta_3 \text{Altitude}_i + \beta_4 \text{Nrooms}_i \\ & + \beta_5 \text{HouseholdMem}_i + \beta_6 \text{Children}_i + \beta_7 \log(\text{Income}_i) \\ & + \beta_8 + \mu_i.\end{aligned}$$

In particular, we examine what happens as  $c_0$  increases from  $10^0$  to  $10^{20}$ . We observe that when  $c_0 = 1$ ,  $BF_{12} = 8.68 \times 10^{+16}$ , which indicates very strong evidence in favor of the model including the price of electricity. However, as  $c_0$  increases, the Bayes factor decreases, which suggests evidence supporting model 2. For instance, when  $c_0 = 10^{20}$ ,  $BF_{12} = 3.11 \times 10^{-4}$ . This is an example

of the issue with using non-informative priors to calculate the Bayes factor: there is very strong evidence supporting the parsimonious model as  $c_0 \rightarrow \infty$ .

We can obtain the posterior predictive distribution of the monthly electricity demand for a household located in the lowest socioeconomic condition in a municipality situated below 1000 meters above sea level, with 2 rooms, 3 members (with children), a monthly income of USD 500, and an electricity price of USD 0.15/kWh. Figure 3.1 shows the histogram of the predictive posterior distribution. The highest posterior density credible interval at 95% is between 44.4 kWh and 373.9 kWh, and the posterior mean is 169.4 kWh.

### *R code. Demand of electricity, posterior predictive distribution*

```

1 rm(list = ls()); set.seed(010101)
2 # Electricity demand
3 DataUt <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/Utilities.csv"
  ", sep = ",", header = TRUE, quote = "") 
4 library(dplyr)
5 DataUtEst <- DataUt %>% filter(Electricity != 0)
6 attach(DataUtEst)
7 # Dependent variable: Monthly consumption (kWh) in log
8 Y <- log(Electricity)
9 # Regressors quantity including intercept
10 X <- cbind(LnPriceElect, IndSocio1, IndSocio2, Altitude,
  Nrooms, HouseholdMem, Children, Lnincome, 1)
11 # LnPriceElect: Price per kWh (USD) in log. IndSocio1,
  IndSocio2, IndSocio3: Indicators socio-economic
  condition (1) is the lowest and (3) the highest.
  Altitude: Indicator of household location (1 is more
  than 1000 meters above sea level). Nrooms: Number of
  rooms in house. HouseholdMem: Number of household
  members. Children: Indicator por presence of children in
  household (1). Lnincome: Monthly income (USD) in log

```

**R code. Demand of electricity, posterior predictive distribution**

---

```

1 k <- dim(X)[2]; N <- dim(X)[1]
2 # Hyperparameters
3 d0 <- 0.001; a0 <- 0.001
4 b0 <- rep(0, k); B0 <- 1000*diag(k)
5 # Posterior parameters
6 bhat <- solve(t(X) %*% X) %*% t(X) %*% Y
7 Bn <- as.matrix(Matrix::forceSymmetric(solve(solve(B0) + t(X)
8 ) %*% X))) # Force this matrix to be symmetric
9 bn <- Bn %*% (solve(B0) %*% b0 + t(X) %*% X %*% bhat)
10 dn <- as.numeric(d0 + t(Y) %*% Y + t(b0) %*% solve(B0) %*% b0 - t(bn) %
11     %*% solve(Bn) %*% bn)
12 an <- a0 + N
13 Hn <- Bn * dn / an
14 # Posterior draws
15 S <- 10000 # Number of draws from posterior distributions
16 sig2 <- MCMCpack::rinvgamma(S, an/2, dn/2)
17 summary(coda::mcmc(sig2))
18 # Log marginal function (multiply by -1 due to minimization)
19 LogMarLikLM <- function(X, c0){
20     k <- dim(X)[2]
21     N <- dim(X)[1]
22     # Hyperparameters
23     B0 <- c0*diag(k)
24     b0 <- rep(0, k)
25     # Posterior parameters
26     bhat <- solve(t(X) %*% X) %*% t(X) %*% Y
27     # Force this matrix to be symmetric
28     Bn <- as.matrix(Matrix::forceSymmetric(solve(solve(B0) + t(X) %*% X)))
29     bn <- Bn %*% (solve(B0) %*% b0 + t(X) %*% X %*% bhat)
30     dn <- as.numeric(d0 + t(Y) %*% Y + t(b0) %*% solve(B0) %*% b0 - t(bn) %
31         %*% solve(Bn) %*% bn)
32     an <- a0 + N
33     # Log marginal likelihood
34     logpy <- (N/2)*log(1/pi)+(a0/2)*log(d0)-(an/2)*log(dn) +
35         0.5*log(det(Bn)/det(B0)) + lgamma(an/2)-lgamma(a0/2)
36     return(-logpy)
37 }
38 cs <- c(10^0, 10^3, 10^6, 10^10, 10^12, 10^15, 10^20)
39 # Observe -1 to recover the right sign
40 LogML <- sapply(cs, function(c) {-LogMarLikLM(c0=c, X = X)})
41 # Regressor without price
42 Xnew <- cbind(IndSocio1, IndSocio2, Altitude, Nrooms,
43     HouseholdMem, Children, Lnincome, 1)
44 # Observe -1 to recover the right sign
45 LogMLnew <- sapply(cs, function(c) {-LogMarLikLM(c0=c, X =
46     Xnew)})
47 # Bayes factor
48 BF <- exp(LogML - LogMLnew)

```

---

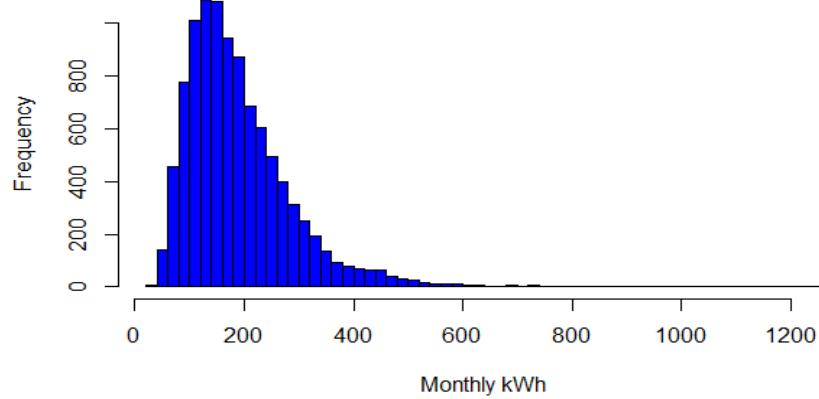
**R code. Demand of electricity, Bayes factor and predictive**

```

1 # Predictive distribution
2 Xpred <- c(log(0.15), 1, 0, 0, 2, 3, 1, log(500), 1)
3 Mean <- Xpred%*%bn
4 Hn <- dn*(1+t(Xpred)%*%Bn%*%Xpred)/an
5 ExpKwH <- exp(LaplacesDemon::rmvt(S, Mean, Hn, an))
6 summary(ExpKwH)
7 HDI <- HDInterval::hdi(ExpKwH, credMass = 0.95) # Highest
   posterior density credible interval
8 hist(ExpKwH, main = "Histogram: Monthly demand of
   electricity", xlab = "Monthly kWh", col = "blue", breaks
   = 50)

```

**Histogram: Monthly demand of electricity**



**FIGURE 3.1**

Histogram using the posterior predictive distribution of electricity demand

### 3.4 Multivariate linear regression: The normal-normal/inverse Wishart model

Let's study the multivariate regression setting where there are  $N$ -dimensional vectors  $\mathbf{y}_m$ , for  $m = 1, 2, \dots, M$ , such that  $\mathbf{y}_m = \mathbf{X}\boldsymbol{\beta}_m + \boldsymbol{\mu}_m$ . Here,  $\mathbf{X}$  represents the set of common regressors, and  $\boldsymbol{\mu}_m$  is the  $N$ -dimensional vector of stochastic errors for each equation. We assume that  $\mathbf{U} = [\boldsymbol{\mu}_1 \ \boldsymbol{\mu}_2 \ \dots \ \boldsymbol{\mu}_M] \sim MN_{N,M}(\mathbf{0}, \mathbf{I}_N, \boldsymbol{\Sigma})$ , which is a matrix variate normal distribution where  $\boldsymbol{\Sigma}$  is the covariance matrix of each  $i$ -th row of  $\mathbf{U}$ , for  $i = 1, 2, \dots, N$ , and we assume independence between the rows. Consequently, we have that  $\text{vec}(\mathbf{U}) \sim N_{N \times M}(\mathbf{0}, \boldsymbol{\Sigma} \otimes \mathbf{I}_N)$ .<sup>9</sup>

This framework can be written in matrix form

$$\underbrace{\begin{bmatrix} y_{11} & y_{12} & \dots & y_{1M} \\ y_{21} & y_{22} & \dots & y_{2M} \\ \vdots & \vdots & \dots & \vdots \\ y_{N1} & y_{N2} & \dots & y_{NM} \end{bmatrix}}_{\mathbf{Y}} = \underbrace{\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1K} \\ x_{21} & x_{22} & \dots & x_{2K} \\ \vdots & \vdots & \dots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NK} \end{bmatrix}}_{\mathbf{X}} \underbrace{\begin{bmatrix} \beta_{11} & \beta_{12} & \dots & \beta_{1M} \\ \beta_{21} & \beta_{22} & \dots & \beta_{2M} \\ \vdots & \vdots & \dots & \vdots \\ \beta_{K1} & \beta_{K2} & \dots & \beta_{KM} \end{bmatrix}}_{\boldsymbol{B}} + \underbrace{\begin{bmatrix} \mu_{11} & \mu_{12} & \dots & \mu_{1M} \\ \mu_{21} & \mu_{22} & \dots & \mu_{2M} \\ \vdots & \vdots & \dots & \vdots \\ \mu_{N1} & \mu_{N2} & \dots & \mu_{NM} \end{bmatrix}}_{\mathbf{U}}$$

Therefore,  $\mathbf{Y} \sim N_{N \times M}(\mathbf{X}\boldsymbol{B}, \boldsymbol{\Sigma} \otimes \mathbf{I}_N)$ ,<sup>10</sup>

$$\begin{aligned} p(\mathbf{Y} \mid \boldsymbol{B}, \boldsymbol{\Sigma}, \mathbf{X}) &\propto |\boldsymbol{\Sigma}|^{-N/2} \exp \left\{ -\frac{1}{2} \text{tr} [(\mathbf{Y} - \mathbf{X}\boldsymbol{B})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{B}) \boldsymbol{\Sigma}^{-1}] \right\} \\ &= |\boldsymbol{\Sigma}|^{-N/2} \exp \left\{ -\frac{1}{2} \text{tr} \left[ \left( \mathbf{S} + (\boldsymbol{B} - \widehat{\boldsymbol{B}})^\top \mathbf{X}^\top \mathbf{X} (\boldsymbol{B} - \widehat{\boldsymbol{B}}) \right) \boldsymbol{\Sigma}^{-1} \right] \right\}, \end{aligned}$$

where  $\mathbf{S} = (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{B}})^\top (\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{B}})$ ,  $\widehat{\boldsymbol{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  (see Exercise 9).

The conjugate prior for this model is  $\pi(\boldsymbol{B}, \boldsymbol{\Sigma}) = \pi(\boldsymbol{B} \mid \boldsymbol{\Sigma})\pi(\boldsymbol{\Sigma})$  where  $\boldsymbol{B} \mid \boldsymbol{\Sigma} \sim N_{K \times M}(\boldsymbol{B}_0, \mathbf{V}_0, \boldsymbol{\Sigma})$  and  $\boldsymbol{\Sigma} \sim IW(\boldsymbol{\Psi}_0, \alpha_0)$ , that is,

$$\pi(\boldsymbol{B}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-K/2} \exp \left\{ -\frac{1}{2} \text{tr} [(\boldsymbol{B} - \boldsymbol{B}_0)^\top \mathbf{V}_0^{-1} (\boldsymbol{B} - \boldsymbol{B}_0) \boldsymbol{\Sigma}^{-1}] \right\}$$

<sup>9</sup> $\text{vec}$  denotes the vectorization operation, and  $\otimes$  denotes the Kronecker product.

<sup>10</sup>We can write down the former expression in a more familiar way using vectorization properties,  $\underbrace{\text{vec}(\mathbf{Y})}_{\mathbf{y}} = \underbrace{(\mathbf{I}_M \otimes \mathbf{X})}_{\mathbf{Z}} \underbrace{\text{vec}(\boldsymbol{B})}_{\boldsymbol{\beta}} + \underbrace{\text{vec}(\mathbf{U})}_{\boldsymbol{\mu}}$ , where  $\mathbf{y} \sim N_{N \times M}(\mathbf{Z}\boldsymbol{\beta}, \boldsymbol{\Sigma} \otimes \mathbf{I}_N)$ .

$$\times |\Sigma|^{-(\alpha_0+M+1)/2} \exp \left\{ -\frac{1}{2} \text{tr} [\Psi_0 \Sigma^{-1}] \right\}.$$

The posterior distribution is given by

$$\begin{aligned} \pi(\mathbf{B}, \Sigma | \mathbf{Y}, \mathbf{X}) &\propto p(\mathbf{Y} | \mathbf{B}, \Sigma, \mathbf{X}) \pi(\mathbf{B} | \Sigma) \pi(\Sigma) \\ &\propto |\Sigma|^{-\frac{N+K+\alpha_0+M+1}{2}} \\ &\quad \times \exp \left\{ -\frac{1}{2} \text{tr} [(\Psi_0 + \mathbf{S} + (\mathbf{B} - \mathbf{B}_0)^\top \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) \right. \\ &\quad \left. + (\mathbf{B} - \widehat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}})) \Sigma^{-1}] \right\}. \end{aligned}$$

Completing the squares on  $\mathbf{B}$  and collecting the remaining terms in the bracket yields  $\Psi_0 + \mathbf{S} + (\mathbf{B} - \mathbf{B}_0)^\top \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) + (\mathbf{B} - \widehat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}})$  is equal to  $(\mathbf{B} - \mathbf{B}_n)^\top \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) + \Psi_n$ ,

where

$$\begin{aligned} \mathbf{B}_n &= (\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{V}_0^{-1} \mathbf{B}_0 + \mathbf{X}^\top \mathbf{Y}) \\ &= (\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{V}_0^{-1} \mathbf{B}_0 + \mathbf{X}^\top \mathbf{X} \widehat{\mathbf{B}}), \\ \mathbf{V}_n &= (\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}, \\ \Psi_n &= \Psi_0 + \mathbf{S} + \mathbf{B}_0^\top \mathbf{V}_0^{-1} \mathbf{B}_0 + \widehat{\mathbf{B}}^\top \mathbf{X}^\top \mathbf{X} \widehat{\mathbf{B}} - \mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n. \end{aligned}$$

Thus, the posterior distribution can be written as

$$\begin{aligned} \pi(\mathbf{B}, \Sigma | \mathbf{Y}, \mathbf{X}) &\propto |\Sigma|^{-K/2} \exp \left\{ -\frac{1}{2} \text{tr} [(\mathbf{B} - \mathbf{B}_n)^\top \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) \Sigma^{-1}] \right\} \\ &\quad \times |\Sigma|^{-\frac{N+\alpha_0+M+1}{2}} \exp \left\{ -\frac{1}{2} \text{tr} [\Psi_n \Sigma^{-1}] \right\}. \end{aligned}$$

That is  $\pi(\mathbf{B}, \Sigma | \mathbf{Y}, \mathbf{X}) = \pi(\mathbf{B} | \Sigma, \mathbf{Y}, \mathbf{X}) \pi(\Sigma | \mathbf{Y}, \mathbf{X})$  where  $\mathbf{B} | \Sigma, \mathbf{Y}, \mathbf{X} \sim N_{K \times M}(\mathbf{B}_n, \mathbf{V}_n, \Sigma)$  and  $\Sigma | \mathbf{Y}, \mathbf{X} \sim IW(\Psi_n, \alpha_n)$ ,  $\alpha_n = N + \alpha_0$ . Observe again that we can write down the posterior mean as a weighted average between prior and sample information such that  $\mathbf{V}_0 \rightarrow \infty$  implies  $\mathbf{B}_n \rightarrow \widehat{\mathbf{B}}$ , as we show in the univariate linear model.

The marginal posterior for  $\mathbf{B}$  is given by

$$\begin{aligned} \pi(\mathbf{B} | \mathbf{Y}, \mathbf{X}) &\propto \int_{\Sigma} |\Sigma|^{-(\alpha_n+K+M+1)/2} \\ &\quad \times \exp \left\{ -\frac{1}{2} \text{tr} \{ [(\mathbf{B} - \mathbf{B}_n)^\top \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) + \Psi_n] \Sigma^{-1} \} \right\} d\Sigma \\ &\propto |(\mathbf{B} - \mathbf{B}_n)^\top \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) + \Psi_n|^{-(K+\alpha_n)/2} \\ &= [|\Psi_n| \times |\mathbf{I}_K + \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) \Psi_n^{-1} (\mathbf{B} - \mathbf{B}_n)^\top|]^{-(\alpha_n+1-M+K+M-1)/2} \\ &\propto |\mathbf{I}_K + \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) \Psi_n^{-1} (\mathbf{B} - \mathbf{B}_n)^\top|^{-(\alpha_n+1-M+K+M-1)/2}. \end{aligned}$$

The second line uses the inverse Wishart distribution, the third line the Sylvester's theorem, and the last line is the kernel of a matrix  $t$ -distribution, that is,  $\mathbf{B} | \mathbf{Y}, \mathbf{X} \sim T_{K \times M}(\mathbf{B}_n, \mathbf{V}_n, \boldsymbol{\Psi}_n)$  with  $\alpha_n + 1 - M$  degrees of freedom.

Observe that  $\text{vec}(\mathbf{B})$  has mean  $\text{vec}(\mathbf{B}_n)$  and variance  $(\mathbf{V}_n \otimes \boldsymbol{\Psi}_n) / (\alpha_n - M - 1)$  based on its marginal distribution. On the other hand, the variance based on the conditional distribution is  $\mathbf{V}_n \otimes \boldsymbol{\Sigma}$ , where the mean of  $\boldsymbol{\Sigma}$  is  $\boldsymbol{\Psi}_n / (\alpha_n - M - 1)$ .

The marginal likelihood is the following,

$$\begin{aligned}
 p(\mathbf{Y}) &= \int_{\mathcal{B}} \int_{\mathcal{S}} \left\{ (2\pi)^{-NM/2} |\boldsymbol{\Sigma}|^{-N/2} \right. \\
 &\quad \exp \left\{ -\frac{1}{2} \text{tr} \left[ \mathbf{S} + (\mathbf{B} - \widehat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}}) \right] \boldsymbol{\Sigma}^{-1} \right\} \\
 &\quad \times (2\pi)^{-KM/2} |\mathbf{V}_0|^{-M/2} |\boldsymbol{\Sigma}|^{-K/2} \exp \left\{ -\frac{1}{2} \text{tr} \left[ (\mathbf{B} - \mathbf{B}_0)^\top \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) \boldsymbol{\Sigma}^{-1} \right] \right\} \\
 &\quad \times \frac{|\boldsymbol{\Psi}_0|^{\alpha_0/2}}{2^{\alpha_0 M/2} \Gamma_M(\alpha_0/2)} |\boldsymbol{\Sigma}|^{-(\alpha_0+M+1)/2} \exp \left\{ -\frac{1}{2} \text{tr} [\boldsymbol{\Psi}_0 \boldsymbol{\Sigma}^{-1}] \right\} d\boldsymbol{\Sigma} d\mathbf{B} \\
 &= (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\boldsymbol{\Psi}_0|^{\alpha_0/2}}{2^{\alpha_0 M/2} \Gamma_M(\alpha_0/2)} \\
 &\quad \times \int_{\mathcal{B}} \int_{\mathcal{S}} \left\{ |\boldsymbol{\Sigma}|^{-(\alpha_0+N+K+M+1)/2} \exp \left[ -\frac{1}{2} \text{tr} \left( \mathbf{S} + (\mathbf{B} - \widehat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}}) \right. \right. \right. \\
 &\quad \left. \left. \left. + (\mathbf{B} - \mathbf{B}_0)^\top \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) + \boldsymbol{\Psi}_0 \right) \boldsymbol{\Sigma}^{-1} \right] \right\} d\boldsymbol{\Sigma} d\mathbf{B} \\
 &= (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\boldsymbol{\Psi}_0|^{\alpha_0/2}}{2^{\alpha_0 M/2} \Gamma_M(\alpha_0/2)} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2) \\
 &\quad \times \int_{\mathcal{B}} \left| \mathbf{S} + (\mathbf{B} - \widehat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}}) \right. \\
 &\quad \left. + (\mathbf{B} - \mathbf{B}_0)^\top \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) + \boldsymbol{\Psi}_0 \right|^{-(\alpha_n+K)/2} d\mathbf{B} \\
 &= (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\boldsymbol{\Psi}_0|^{\alpha_0/2}}{2^{\alpha_0 M/2} \Gamma_M(\alpha_0/2)} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2) \\
 &\quad \times \int_{\mathcal{B}} \left| (\mathbf{B} - \widehat{\mathbf{B}}_n)^\top \mathbf{V}_n^{-1} (\mathbf{B} - \widehat{\mathbf{B}}_n) + \boldsymbol{\Psi}_n \right|^{-(\alpha_n+K)/2} d\mathbf{B} \\
 &= (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\boldsymbol{\Psi}_0|^{\alpha_0/2}}{2^{\alpha_0 M/2} \Gamma_M(\alpha_0/2)} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2) \\
 &\quad \times \int_{\mathcal{B}} \left[ |\boldsymbol{\Psi}_n| \times |\mathbf{I}_K + \mathbf{V}_n^{-1} (\mathbf{B} - \widehat{\mathbf{B}}_n) \boldsymbol{\Psi}_n^{-1} (\mathbf{B} - \widehat{\mathbf{B}}_n)^\top| \right]^{-(\alpha_n+K)/2} d\mathbf{B} \\
 &= |\boldsymbol{\Psi}_n|^{-(\alpha_n+K)/2} (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \\
 &\quad \times \frac{|\boldsymbol{\Psi}_0|^{\alpha_0/2} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2)}{2^{\alpha_0 M/2} \Gamma_M(\alpha_0/2)}
 \end{aligned}$$

$$\begin{aligned}
& \times \int_{\mathcal{B}} \left| \mathbf{I}_K + \mathbf{V}_n^{-1}(\mathbf{B} - \hat{\mathbf{B}}_n) \Psi_n^{-1} (\mathbf{B} - \hat{\mathbf{B}}_n)^\top \right|^{-(\alpha_n+1-M+K+M-1)/2} d\mathbf{B} \\
& = |\Psi_n|^{-(\alpha_n+K)/2} (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \\
& \quad \times \frac{|\Psi_0|^{\alpha_0/2} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2)}{2^{\alpha_0 M/2} \Gamma_M(\alpha_0/2)} \\
& \quad \times \pi^{MK/2} \frac{\Gamma_M((\alpha_n+1-M+M-1)/2)}{\Gamma_M((\alpha_n+1-M+K+M-1)/2)} |\Psi_n|^{K/2} |\mathbf{V}_n|^{M/2} \\
& = \frac{|\mathbf{V}_n|^{M/2} |\Psi_0|^{\alpha_0/2} \Gamma_M(\alpha_n/2)}{|\mathbf{V}_0|^{M/2} |\Psi_n|^{\alpha_n/2} \Gamma_M(\alpha_0/2)} \pi^{-MN/2}.
\end{aligned}$$

The third equality follows from the kernel of an inverse Wishart distribution, the fifth from Sylvester's theorem, and the seventh from the kernel of a matrix  $t$ -distribution.

Observe that this last expression is the multivariate case of the marginal likelihood of the univariate regression model. Taking into account that

$$\begin{aligned}
(\mathbf{A} + \mathbf{B})^{-1} &= \mathbf{A}^{-1} - (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \mathbf{A}^{-1} \\
&= \mathbf{B}^{-1} - (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \mathbf{B}^{-1} \\
&= \mathbf{A}^{-1} (\mathbf{A}^{-1} + \mathbf{B}^{-1}) \mathbf{B}^{-1},
\end{aligned}$$

we can show that  $\Psi_n = \Psi_0 + \mathbf{S} + (\hat{\mathbf{B}} - \mathbf{B}_0)^\top \mathbf{V}_n (\hat{\mathbf{B}} - \mathbf{B}_0)$  (see Exercise 7). Therefore, the marginal likelihood rewards fit (smaller sum of squares,  $\mathbf{S}$ ), similarity between prior and sample information regarding location parameters, and information gains in variability from  $\mathbf{V}_0$  to  $\mathbf{V}_n$ .

Given a matrix of regressors  $\mathbf{X}_0$  for  $N_0$  unobserved units, the predictive density of  $\mathbf{Y}_0$  given  $\mathbf{Y}$ ,  $\pi(\mathbf{Y}_0 | \mathbf{Y})$  is a matrix  $t$  distribution  $T_{N_0, M}(\alpha_n - M + 1, \mathbf{X}_0 \mathbf{B}_n, \mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{V}_n \mathbf{X}_0^\top, \Psi_n)$  (see Exercise 6). Observe that the prediction is centered at  $\mathbf{X}_0 \mathbf{B}_n$ , and the covariance matrix of  $\text{vec}(\mathbf{Y}_0)$  is  $\frac{(\mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{V}_n \mathbf{X}_0^\top) \otimes \Psi_n}{\alpha_n - M - 1}$ .

### 3.5 Summary

We introduce conjugate family models for both discrete and continuous data. These models form the foundation of the Bayesian framework due to their mathematical tractability, as they provide closed-form expressions for the posterior distributions, marginal likelihood, and predictive distribution. Additionally, we present the Bayesian linear regression frameworks for both univariate and multivariate cases under conjugate families. These frameworks are fundamental for performing regression analysis in the Bayesian setting.

---

### 3.6 Exercises

1. Write the distribution of the Bernoulli example in canonical form, and find the mean and variance of the sufficient statistic.
2. Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from  $N$  binomial experiments, each with known size  $n_i$  and the same unknown probability  $\theta$ , show that  $p(\mathbf{y} | \theta)$  is in the exponential family. Then, find the posterior distribution, the marginal likelihood, and the predictive distribution of the binomial-Beta model assuming the number of trials is known.
3. Given a random sample  $\mathbf{Y} = [Y_1 \ Y_2 \ \dots \ Y_N]^\top$  from an exponential distribution, show that  $p(\mathbf{y} | \lambda)$  is in the exponential family. Additionally, find the posterior distribution, the marginal likelihood, and the predictive distribution of the exponential-Gamma model.
4. Given that  $\mathbf{Y} \sim N_N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , that is, a multivariate normal distribution, show that  $p(\mathbf{y} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is in the exponential family.
5. Find the marginal likelihood in the normal/inverse-Wishart model.
6. Find the posterior predictive distribution in the normal/inverse-Wishart model, and show that  $\mathbf{Y}_0 \mid \mathbf{Y} \sim T_{N_0, M}(\alpha_n - M + 1, \mathbf{X}_0 \mathbf{B}_n, \mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{V}_n \mathbf{X}_0^\top, \boldsymbol{\Psi}_n)$ .
7. Show that  $\delta_n = \delta_0 + (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top((\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{B}_0)^{-1}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$  in the linear regression model, and that  $\boldsymbol{\Psi}_n = \boldsymbol{\Psi}_0 + \mathbf{S} + (\hat{\mathbf{B}} - \mathbf{B}_0)^\top \mathbf{V}_n (\hat{\mathbf{B}} - \mathbf{B}_0)$  in the linear multivariate regression model.
8. Show that in the linear regression model  $\boldsymbol{\beta}_n^\top(\mathbf{B}_n^{-1} - \mathbf{B}_n^{-1} \mathbf{M}^{-1} \mathbf{B}_n^{-1})\boldsymbol{\beta}_n = \boldsymbol{\beta}_{**}^\top \mathbf{C} \boldsymbol{\beta}_{**}$  and  $\boldsymbol{\beta}_{**} = \mathbf{X}_0 \boldsymbol{\beta}_n$ .
9. Show that  $(\mathbf{Y} - \mathbf{X}\mathbf{B})^\top(\mathbf{Y} - \mathbf{X}\mathbf{B}) = \mathbf{S} + (\mathbf{B} - \hat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \hat{\mathbf{B}})$  where  $\mathbf{S} = (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^\top(\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$ ,  $\hat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  in the multivariate regression model.
10. **What is the probability that the Sun will rise tomorrow?**  
This is the most famous example by Richard Price, developed in the Appendix of Bayes' theorem paper [33]. Here, we implicitly use Laplace's Rule of Succession to solve this problem. In particular, if we were a priori uncertain about the probability that the Sun will rise on a specified day, we can assume a uniform prior distribution over  $(0, 1)$ , that is, a Beta(1,1) distribution. Then, what is the probability that the Sun will rise?
11. Using information from Public Policy Polling in September 27th-28th for the 2016 presidential five-way race in USA, there are 411,

373 and 149 sampled people supporting Hillary Clinton, Donald Trump and other, respectively.

- Find the posterior probability of the percentage difference of people supporting Hillary versus Trump according to this data using a non-informative prior, that is,  $\alpha_0 = [1 \ 1 \ 1]$  in the multinomial-Dirichlet model. What is the probability of having more supports of Hillary vs Trump?
- What is the probability that sampling one hundred independent individuals 44, 40 and 16 support Hillary, Trump and other, respectively?

#### 12. Math test example continues

You have a random sample of math scores of size  $N = 50$  from a normal distribution,  $Y_i \sim N(\mu, \sigma^2)$ . The sample mean and variance are equal to 102 and 10, respectively. Using the normal-normal/inverse-gamma model where  $\mu_0 = 100$ ,  $\beta_0 = 1$ ,  $\alpha_0 = \delta_0 = 0.001$

- Get a 95% confidence and credible interval for  $\mu$ .
- What is the posterior probability that  $\mu > 103$ ?

#### 13. Demand of electricity example continues

Set  $c_0$  such that maximizes the marginal likelihood in the specifications with and without electricity price in the example of demand of electricity (empirical Bayes). Then, calculate the Bayes factor, and conclude if there is evidence supporting the inclusion of the price of electricity in the demand equation.

#### 14. Utility demand

Use the file *Utilities.csv* to estimate a multivariate linear regression model where  $\mathbf{Y}_i = [\log(\text{electricity}_i) \ \log(\text{water}_i) \ \log(\text{gas}_i)]$  as function of  $\log(\text{electricity price}_i)$ ,  $\log(\text{water price}_i)$ ,  $\log(\text{gas price}_i)$ ,  $\text{IndSocio1}_i$ ,  $\text{IndSocio2}_i$ ,  $\text{Altitude}_i$ ,  $\text{Nrooms}_i$ ,  $\text{HouseholdMem}_i$ ,  $\text{Children}_i$ , and  $\log(\text{Income}_i)$ , where electricity, water and gas are monthly consumption of electricity (kWh), water ( $\text{m}^3$ ) and gas ( $\text{m}^3$ ), and other definitions are given in the Example of Section 3.3. Omit households that do not consume any of the utilities in this exercise.

Set a non-informative prior framework,  $\mathbf{B}_0 = [0]_{11 \times 3}$ ,  $\mathbf{V}_0 = 1000\mathbf{I}_{11}$ ,  $\Psi_0 = 1000\mathbf{I}_3$  and  $\alpha_0 = 3$ , where we have  $K = 11$  (regressors plus intercept) and  $M = 3$  (equations) in this exercise.

- Find the posterior mean estimates and the highest posterior density intervals at 95% of  $\mathbf{B}$  and  $\Sigma$ . Use the marginal distribution and the conditional distribution to obtain the posterior estimates of  $\mathbf{B}$ , and compare the results.

- Find the Bayes factor comparing the baseline model in this exercise with the same specification but using the income in dollars. Now, calculate the Bayes factor using the income in thousand dollars. Is there any difference?
- Find the predictive distribution for the monthly demand of electricity, water and gas in the baseline specification of a household located in the lowest socioeconomic condition in a municipality located below 1000 meters above the sea level, 2 rooms, 3 members with children, a monthly income equal to USD 500, an electricity price equal to USD/kWh 0.15, a water price equal to USD/M<sup>3</sup> 0.70, and a gas price equal to USD/M<sup>3</sup> 0.75.

### 15. Ph.D. students sleeping hours [6, Chap. 2]

We are interested in learning about the proportion of Ph.D. students who sleep at least 6 hours per day. We have a sample of 52 students, where 15 report sleeping at least 6 hours, and the remaining 37 report not sleeping at least 6 hours. The prior distribution is a Beta distribution, with hyperparameters calibrated so that the prior probabilities of the proportion of students who sleep least than 6 hours being less than 0.4 and 0.75 are 0.6 and 0.95, respectively. Estimate the 95% posterior credible interval for the proportion of Ph.D. students who sleep at least 6 hours per day. Then, assume there is a group of experts whose beliefs about the proportion of Ph.D. students sleeping at least 6 hours are represented in the following table:

**TABLE 3.1**

Probability distribution: Ph.D students that sleep at least 6 hours per day.

$h$	0.05	0.10	0.15	0.20	0.25	0.30	0.35	0.40	0.45	0.50	0.55
$P(p = h)$	0.05	0.07	0.10	0.12	0.15	0.17	0.15	0.11	0.06	0.01	0.01

Use Table 3.1 as prior information, and find the posterior distribution of the proportion of students that sleep at least 6 hours.



# 4

---

## *Simulation methods*

---

In the previous chapters, we focused on conjugate families, where the posterior and predictive distributions have standard analytical forms (e.g., normal, Student's t, gamma, binomial, Poisson, etc.) and where the marginal likelihood has a closed-form analytical solution. However, realistic models are often more complex and lack such closed-form solutions.

To address this complexity, we rely on simulation (stochastic) methods to draw samples from posterior and predictive distributions. This chapter introduces posterior simulation, a cornerstone of Bayesian inference. We discuss Markov Chain Monte Carlo (MCMC) methods, including Gibbs sampling, Metropolis-Hastings, and Hamiltonian Monte Carlo, as well as other techniques like importance sampling and particle filtering (sequential Monte Carlo).

The simulation methods discussed in this chapter are specifically applied throughout this book. However, we do not delve into deterministic methods, such as numerical integration (quadrature), or other simulation methods, including discrete approximation, the probability integral transform, the method of composition, accept-reject sampling, and slice sampling algorithms. While these methods are also widely used, they are not as common as the approaches explicitly employed in this book.

For readers interested in these alternative methods, we recommend exploring [371, Chaps. 2 and 3], [372, Chaps. 2, 3, and 8], [184, Chap. 5], and [160, Chap. 10].

---

### 4.1 Markov chain Monte Carlo methods

Markov Chain Monte Carlo (MCMC) methods are algorithms used to approximate complex probability distributions by constructing a Markov chain. This chain is a sequence of random samples where each sample depends only on the previous one. The goal of MCMC methods is to obtain draws from the posterior distribution as the equilibrium distribution. The key point in MCMC methods is the transition kernel or density,  $q(\boldsymbol{\theta}^{(s)} \mid \boldsymbol{\theta}^{(s-1)})$ , which generates a draw  $\boldsymbol{\theta}^{(s)}$  at stage  $s$  that depends solely on  $\boldsymbol{\theta}^{(s-1)}$ . This transition distribution must be designed such that the Markov chain converges to a unique

stationary distribution, which, in our case, is the posterior distribution, that is,  $\pi(\boldsymbol{\theta}^{(s)} | \mathbf{y}) = \int_{\Theta} q(\boldsymbol{\theta}^{(s)} | \boldsymbol{\theta}^{(s-1)})\pi(\boldsymbol{\theta}^{(s-1)} | \mathbf{y})d\boldsymbol{\theta}^{(s-1)}$ .

Given that we start at an arbitrary point,  $\boldsymbol{\theta}^{(0)}$ , the algorithm requires that the Markov chain be *irreducible*, meaning that the process can reach any other state with positive probability. Additionally, the process must be *aperiodic*, meaning that for each state, the greatest common divisor of the number of steps it takes to return to the state is 1, ensuring that there are no cycles forcing the system to return to a state only after a fixed number of steps. Furthermore, the process must be *recurrent*, meaning that it will return to any state an infinite number of times with probability one. However, to ensure convergence to the stationary distribution, a stronger condition is required: the process must be *positive recurrent*, meaning that the expected return time to a state is finite. Given an *irreducible*, *aperiodic*, and *positive recurrent* transition density, the Markov chain algorithm will asymptotically converge to the stationary posterior distribution we are seeking. For more details, see [372, chap. 6].

#### 4.1.1 Gibbs sampler

The Gibbs sampler algorithm is one of the most widely used MCMC methods for sampling from non-standard distributions in Bayesian analysis. While it is a special case of the Metropolis-Hastings (MH) algorithm, it originated from a different theoretical background [165, 156]. The key requirement for implementing the Gibbs sampling algorithm is the availability of conditional posterior distributions. The algorithm works by cycling through the conditional posterior distributions corresponding to different blocks of the parameter space under inference.

Two simplify concepts let's focus on a parameter space composed by two blocks,  $\boldsymbol{\theta} = [\boldsymbol{\theta}_1 \ \boldsymbol{\theta}_2]^\top$ , the Gibbs sampling algorithm uses as transition kernel  $q(\boldsymbol{\theta}_1^{(s)}, \boldsymbol{\theta}_2^{(s)} | \boldsymbol{\theta}_1^{(s-1)}, \boldsymbol{\theta}_2^{(s-1)}) = \pi(\boldsymbol{\theta}_1^{(s)} | \boldsymbol{\theta}_2^{(s-1)}, \mathbf{y})\pi(\boldsymbol{\theta}_2^{(s)} | \boldsymbol{\theta}_1^{(s)}, \mathbf{y})$ . Thus,

$$\begin{aligned} \int_{\Theta} q(\boldsymbol{\theta}^{(s)} | \boldsymbol{\theta}^{(s-1)})\pi(\boldsymbol{\theta}^{(s-1)} | \mathbf{y})d\boldsymbol{\theta}^{(s-1)} &= \int_{\Theta_2} \int_{\Theta_1} \pi(\boldsymbol{\theta}_1^{(s)} | \boldsymbol{\theta}_2^{(s-1)}, \mathbf{y})\pi(\boldsymbol{\theta}_2^{(s)} | \boldsymbol{\theta}_1^{(s)}, \mathbf{y}) \\ &\quad \pi(\boldsymbol{\theta}_1^{(s-1)}, \boldsymbol{\theta}_2^{(s-1)} | \mathbf{y})d\boldsymbol{\theta}_1^{(s-1)}d\boldsymbol{\theta}_2^{(s-1)} \\ &= \pi(\boldsymbol{\theta}_2^{(s)} | \boldsymbol{\theta}_1^{(s)}, \mathbf{y}) \int_{\Theta_2} \int_{\Theta_1} \pi(\boldsymbol{\theta}_1^{(s)} | \boldsymbol{\theta}_2^{(s-1)}, \mathbf{y}) \\ &\quad \pi(\boldsymbol{\theta}_1^{(s-1)}, \boldsymbol{\theta}_2^{(s-1)} | \mathbf{y})d\boldsymbol{\theta}_1^{(s-1)}d\boldsymbol{\theta}_2^{(s-1)} \\ &= \pi(\boldsymbol{\theta}_2^{(s)} | \boldsymbol{\theta}_1^{(s)}, \mathbf{y}) \int_{\Theta_2} \pi(\boldsymbol{\theta}_1^{(s)} | \boldsymbol{\theta}_2^{(s-1)}, \mathbf{y}) \\ &\quad \pi(\boldsymbol{\theta}_2^{(s-1)} | \mathbf{y})d\boldsymbol{\theta}_2^{(s-1)} \\ &= \pi(\boldsymbol{\theta}_2^{(s)} | \boldsymbol{\theta}_1^{(s)}, \mathbf{y}) \int_{\Theta_2} \pi(\boldsymbol{\theta}_1^{(s)}, \boldsymbol{\theta}_2^{(s-1)} | \mathbf{y})d\boldsymbol{\theta}_2^{(s-1)} \end{aligned}$$

$$\begin{aligned}
&= \pi(\boldsymbol{\theta}_2^{(s)} \mid \boldsymbol{\theta}_1^{(s)}, \mathbf{y}) \pi(\boldsymbol{\theta}_1^{(s)} \mid \mathbf{y}) \\
&= \pi(\boldsymbol{\theta}_1^{(s)}, \boldsymbol{\theta}_2^{(s)} \mid \mathbf{y}).
\end{aligned}$$

Then,  $\pi(\boldsymbol{\theta} \mid \mathbf{y})$  is the stationary distribution for the Gibbs transition kernel.

A word of caution! Even if we have well-defined conditional posterior distributions  $\pi(\boldsymbol{\theta}_1^{(s)} \mid \boldsymbol{\theta}_2^{(s-1)}, \mathbf{y})$  and  $\pi(\boldsymbol{\theta}_2^{(s)} \mid \boldsymbol{\theta}_1^{(s)}, \mathbf{y})$ , and we can simulate from them, the joint posterior distribution  $\pi(\boldsymbol{\theta}_1^{(s)}, \boldsymbol{\theta}_2^{(s)} \mid \mathbf{y})$  may not correspond to any proper distribution. We should be mindful of this situation, especially when dealing with improper prior distributions (see [372, Chap. 10] for details).

Algorithm A1 demonstrates the implementation of a Gibbs sampler with  $d$  blocks. The number of iterations ( $S$ ) is chosen to ensure convergence to the stationary distribution. In Section 4.4, we review several convergence diagnostics to assess whether the posterior draws have reached convergence.

---

**Algorithm A1** Gibbs sampling

---

- 1: Set  $\boldsymbol{\theta}_2^{(0)}, \boldsymbol{\theta}_3^{(0)}, \dots, \boldsymbol{\theta}_d^{(0)}$
  - 2: **for**  $s = 1, \dots, S$  **do**
  - 3:     Draw  $\boldsymbol{\theta}_1^{(s)}$  from  $\pi(\boldsymbol{\theta}_1^{(s)} \mid \boldsymbol{\theta}_2^{(s-1)}, \dots, \boldsymbol{\theta}_d^{(s-1)}, \mathbf{y})$
  - 4:     Draw  $\boldsymbol{\theta}_2^{(s)}$  from  $\pi(\boldsymbol{\theta}_2^{(s)} \mid \boldsymbol{\theta}_1^{(s)}, \dots, \boldsymbol{\theta}_d^{(s-1)}, \mathbf{y})$
  - 5:     ⋮
  - 6:     Draw  $\boldsymbol{\theta}_d^{(s)}$  from  $\pi(\boldsymbol{\theta}_d^{(s)} \mid \boldsymbol{\theta}_1^{(s)}, \dots, \boldsymbol{\theta}_{d-1}^{(s)}, \mathbf{y})$
  - 7: **end for**
- 

**Example: Mining disaster change point**

Let's use the dataset *Mining.csv* provided by [68]. This dataset records the number of mining disasters per year from 1851 to 1962 in British coal mines.

We assume there is an unknown structural change point in the number of mining disasters, where the parameters of the Poisson distributions change. In particular,

$$p(y_t) = \begin{cases} \frac{\exp(-\lambda_1)\lambda_1^{y_t}}{y_t!}, & t = 1, 2, \dots, H \\ \frac{\exp(-\lambda_2)\lambda_2^{y_t}}{y_t!}, & t = H + 1, \dots, T \end{cases},$$

where  $H$  is the changing point.

We use conjugate families for  $\lambda_l$ ,  $l = 1, 2$ , where  $\lambda_l \sim G(\alpha_{l0}, \beta_{l0})$ , and set  $\pi(H) = 1/T$ , which corresponds to a discrete uniform distribution for the change point. This implies that, a priori, we assume equal probability for any time to be the change point.

The posterior distribution is

$$\pi(\lambda_1, \lambda_2, H \mid \mathbf{y}) \propto \prod_{t=1}^H \frac{\exp(-\lambda_1)\lambda_1^{y_t}}{y_t!} \prod_{t=H+1}^T \frac{\exp(-\lambda_2)\lambda_2^{y_t}}{y_t!}$$

$$\begin{aligned} & \times \exp(-\beta_{10}\lambda_1)\lambda_1^{\alpha_{10}-1} \exp(-\beta_{20}\lambda_2)\lambda_2^{\alpha_{20}-1} 1/T \\ & \propto \exp(-H\lambda_1)\lambda_1^{\sum_{t=1}^H y_t} \exp(-(T-H)\lambda_2)\lambda_2^{\sum_{t=H+1}^T y_t} \\ & \quad \times \exp(-\beta_{10}\lambda_1)\lambda_1^{\alpha_{10}-1} \exp(-\beta_{20}\lambda_2)\lambda_2^{\alpha_{20}-1}. \end{aligned}$$

Then, the conditional posterior distribution of  $\lambda_1 | \lambda_2, H, \mathbf{y}$  is

$$\pi(\lambda_1 | \lambda_2, H, \mathbf{y}) \propto \exp(-(H + \beta_{10})\lambda_1)\lambda_1^{\sum_{t=1}^H y_t + \alpha_{10}-1},$$

that is,  $\lambda_1 | \lambda_2, H, \mathbf{y} \sim G(\alpha_{1n}, \beta_{1n})$ ,  $\beta_{1n} = H + \beta_{10}$  and  $\alpha_{1n} = \sum_{t=1}^H y_t + \alpha_{10}$ .

The conditional posterior distribution of  $\lambda_2 | \lambda_1, H, \mathbf{y}$  is

$$\pi(\lambda_2 | \lambda_1, H, \mathbf{y}) \propto \exp(-((T-H) + \beta_{20})\lambda_2)\lambda_2^{\sum_{t=H+1}^T y_t + \alpha_{20}-1},$$

that is,  $\lambda_2 | \lambda_1, H, \mathbf{y} \sim G(\alpha_{2n}, \beta_{2n})$ ,  $\beta_{2n} = (T-H) + \beta_{20}$  and  $\alpha_{2n} = \sum_{t=H+1}^T y_t + \alpha_{20}$ .

The conditional posterior distribution of the change point is

$$\begin{aligned} \pi(H | \lambda_1, \lambda_2, \mathbf{y}) & \propto \exp(-H\lambda_1)\lambda_1^{\sum_{t=1}^H y_t} \exp(-(T-H)\lambda_2)\lambda_2^{\sum_{t=H+1}^T y_t} \\ & \propto \exp(-H(\lambda_1 - \lambda_2))\lambda_1^{\sum_{t=1}^H y_t} \lambda_2^{\sum_{t=H+1}^T y_t} \exp(-T\lambda_2) \frac{\lambda_2^{\sum_{t=1}^H y_t}}{\lambda_2^{\sum_{t=1}^H y_t}} \\ & \propto \exp(-H(\lambda_1 - \lambda_2)) \left( \frac{\lambda_1}{\lambda_2} \right)^{\sum_{t=1}^H y_t}. \end{aligned}$$

Thus, the conditional posterior distribution of  $H$  is

$$\pi(H | \lambda_1, \lambda_2, \mathbf{y}) = \frac{\exp(-H(\lambda_1 - \lambda_2)) \left( \frac{\lambda_1}{\lambda_2} \right)^{\sum_{t=1}^H y_t}}{\sum_{H=1}^T \exp(-H(\lambda_1 - \lambda_2)) \left( \frac{\lambda_1}{\lambda_2} \right)^{\sum_{t=1}^H y_t}}, \quad H = 1, 2, \dots, T.$$

The code shows how to do a Gibbs sampling algorithm to perform inference of this model using the hyperparameters suggested by [184, Chap. 7],  $\alpha_{l0} = 0.5$  and  $\beta_{l0} = 1$ ,  $l = 1, 2$ .

### **R code. Gibbs sampler: The mining disaster changepoint**

```

1 rm(list = ls()); set.seed(010101)
2 dataset<-read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/
  MiningDataCarlin.csv",header=T)
3 attach(dataset); str(dataset)
4 a10<-0.5; a20<-0.5; b10<-1; b20<-1
5 y<-Count; sumy<-sum(Count); N<-length(Count)
6 theta1<-NULL; theta2<-NULL
7 kk<-NULL; k<-60; S<-10000
8 for(i in 1:S){
9   a1<-a10+sum(y[1:k]); b1<-b10+k
10  theta11<-rgamma(1,a1,b1); theta1<-c(theta1,theta11)
11  a2<-a20+sum(y[(1+k):N]); b2<-b20+N-k
12  theta22<-rgamma(1,a2,b2); theta2<-c(theta2,theta22)
13  pp<-NULL
14  for(l in 1:N){
15    p<-exp(1*(theta22-theta11))*(theta11/theta22)^(sum(y[1:l]
16      ]))
17    pp<-c(pp,p)
18  }
19  prob<-pp/sum(pp); k<-sample(1:N,1,prob=prob)
20  kk<-c(kk,k)
21 library(coda); summary(mcmc(theta1)); summary(mcmc(theta2))
22 summary(mcmc(kk))
23 hist(kk, main = "Histogram: Posterior mean change point",
  xlab = "Posterior mean", col = "blue", breaks = 25)

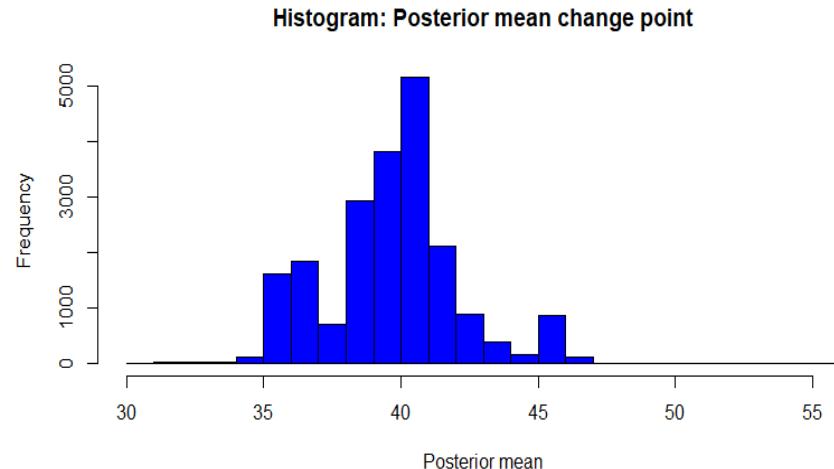
```

The posterior results indicate that the rate of disasters decrease from 3.1 to 0.92 per year in 1890.

Figure 4.1 shows the histogram of the posterior draws of the change point in mining disasters.

#### 4.1.2 Metropolis-Hastings

The Metropolis-Hastings (M-H) algorithm [308, 196] is a general MCMC method that does not require standard closed-form solutions for the conditional posterior distributions. The key idea is to use a transition kernel whose unique invariant distribution is  $\pi(\boldsymbol{\theta} | \mathbf{y})$ . This kernel must satisfy the *balancing condition*, meaning that, given a realization  $\boldsymbol{\theta}^{(s-1)}$  at stage  $s-1$  from the stationary distribution  $\pi(\boldsymbol{\theta} | \mathbf{y})$ , we generate a candidate draw  $\boldsymbol{\theta}^c$  from the

**FIGURE 4.1**

Histogram of posterior draws of change point: Mining disasters

proposal distribution  $q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)})$  at stage  $s$  such that:

$$q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)})\pi(\boldsymbol{\theta}^{(s-1)} \mid \mathbf{y}) = q(\boldsymbol{\theta}^{(s-1)} \mid \boldsymbol{\theta}^c)\pi(\boldsymbol{\theta}^c \mid \mathbf{y}),$$

which implies that the probability of moving from  $\boldsymbol{\theta}^{(s-1)}$  to  $\boldsymbol{\theta}^c$  is equal to the probability of moving from  $\boldsymbol{\theta}^c$  to  $\boldsymbol{\theta}^{(s-1)}$ .

In general, the *balancing condition* is not automatically satisfied, and we must introduce an *acceptance probability*  $\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c)$  to ensure that the condition holds:

$$q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)})\pi(\boldsymbol{\theta}^{(s-1)} \mid \mathbf{y})\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c) = q(\boldsymbol{\theta}^{(s-1)} \mid \boldsymbol{\theta}^c)\pi(\boldsymbol{\theta}^c \mid \mathbf{y}).$$

Thus, the acceptance probability is given by:

$$\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c) = \min \left\{ \frac{q(\boldsymbol{\theta}^{(s-1)} \mid \boldsymbol{\theta}^c)\pi(\boldsymbol{\theta}^c \mid \mathbf{y})}{q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)})\pi(\boldsymbol{\theta}^{(s-1)} \mid \mathbf{y})}, 1 \right\},$$

where  $q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)})$  and  $\pi(\boldsymbol{\theta}^{(s-1)} \mid \mathbf{y})$  must be nonzero, as transitioning from  $\boldsymbol{\theta}^{(s-1)}$  to  $\boldsymbol{\theta}^c$  is only possible under these conditions.

Algorithm A2 shows how to implement a Metropolis-Hastings algorithm. The number of iterations ( $S$ ) is chosen to ensure convergence to the stationary distribution.

Some remarks: First, we do not need to know the marginal likelihood to implement the M-H algorithm, as it cancels out when calculating the acceptance probability. Specifically, given that  $\pi(\boldsymbol{\theta} \mid \mathbf{y}) \propto \pi(\boldsymbol{\theta}) \times p(\mathbf{y} \mid \boldsymbol{\theta})$ , we can

**Algorithm A2** Metropolis-Hastings algorithm

---

```

1: Set  $\boldsymbol{\theta}^{(0)}$  in the support of  $\pi(\boldsymbol{\theta} \mid \mathbf{y})$ 
2: for  $s = 1, \dots, S$  do
3:   Draw  $\boldsymbol{\theta}^c$  from  $q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)})$ 
4:   Calculate  $\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c) = \min \left\{ \frac{q(\boldsymbol{\theta}^{(s-1)} \mid \boldsymbol{\theta}^c) \pi(\boldsymbol{\theta}^c \mid \mathbf{y})}{q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)}) \pi(\boldsymbol{\theta}^{(s-1)} \mid \mathbf{y})}, 1 \right\}$ 
5:   Draw  $U$  from  $U(0, 1)$ 
6:    $\boldsymbol{\theta}^{(s)} = \begin{cases} \boldsymbol{\theta}^c & \text{if } U \leq \alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c) \\ \boldsymbol{\theta}^{(s-1)} & \text{otherwise} \end{cases}$ 
7: end for

```

---

use the right-hand side expression to compute the acceptance probability. Second, the Gibbs sampling algorithm is a particular case of the M-H algorithm where the acceptance probability is equal to 1 ([164] and [372, Chap. 10], see Exercise 2). Third, we can combine the M-H and Gibbs sampling algorithms when dealing with relatively complex posterior distributions. Specifically, the Gibbs sampling algorithm can be used for blocks with conditional posterior distributions in standard closed forms, while the M-H algorithm is applied to sample from conditional posterior distributions that do not have standard forms. This approach is known as the M-H within Gibbs sampling algorithm. Fourth, we can note that the transition kernel in the M-H algorithm is a mixture of a continuous density ( $q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)})$ ) and a probability mass function ( $\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c)$ ) [92].

Fifth, a crucial point associated with the proposal densities is the acceptance probability. Low or high acceptance probabilities are not ideal. A low rate implies poor mixing, meaning the chain does not move effectively through the support of the posterior distribution. Conversely, a high acceptance rate implies that the chain will converge too slowly. A sensible value depends on the dimension of the parameter space. A rule of thumb is that if the dimension is less than or equal to 2, the acceptance rate should be around 0.50. If the dimension is greater than 2, the acceptance rate should be approximately 0.25 [373]. For technical details of the Metropolis-Hastings algorithm, see [372, Chap. 7].

Regarding the proposal density, it must be positive everywhere the posterior distribution is positive. This ensures that the Markov chain can explore the entire support of the posterior distribution. Additionally, the proposal density must allow the Markov chain to reach any region of the posterior distribution's support. There are three standard approaches for choosing the proposal density: the independent proposal, the random walk proposal, and the tailored proposal.

In the independent proposal,  $q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)}) = q(\boldsymbol{\theta}^c)$ , which implies that

$$\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c) = \min \left\{ \frac{q(\boldsymbol{\theta}^{(s-1)}) \pi(\boldsymbol{\theta}^c \mid \mathbf{y})}{q(\boldsymbol{\theta}^c) \pi(\boldsymbol{\theta}^{(s-1)} \mid \mathbf{y})}, 1 \right\}.$$

In this case, a move from  $\boldsymbol{\theta}^{(s-1)}$  to  $\boldsymbol{\theta}^c$  is always accepted if  $q(\boldsymbol{\theta}^{(s-1)})\pi(\boldsymbol{\theta}^c | \mathbf{y}) \geq q(\boldsymbol{\theta}^c)\pi(\boldsymbol{\theta}^{(s-1)} | \mathbf{y})$ .

In the random walk proposal,  $\boldsymbol{\theta}^c = \boldsymbol{\theta}^{(s-1)} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon}$  is a random perturbation. If  $p(\boldsymbol{\epsilon}) = p(-\boldsymbol{\epsilon})$ , meaning the distribution  $p(\boldsymbol{\epsilon})$  is symmetric around zero, then  $q(\boldsymbol{\theta}^c | \boldsymbol{\theta}^{(s-1)}) = q(\boldsymbol{\theta}^{(s-1)} | \boldsymbol{\theta}^c)$ . This was the original Metropolis algorithm [308]. Thus, the acceptance rate is

$$\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c) = \min \left\{ \frac{\pi(\boldsymbol{\theta}^c | \mathbf{y})}{\pi(\boldsymbol{\theta}^{(s-1)} | \mathbf{y})}, 1 \right\}.$$

In this case, a move from  $\boldsymbol{\theta}^{(s-1)}$  to  $\boldsymbol{\theta}^c$  is always accepted if  $\pi(\boldsymbol{\theta}^c | \mathbf{y}) \geq \pi(\boldsymbol{\theta}^{(s-1)} | \mathbf{y})$ .

In the tailored proposal, the density is designed to have fat tails, is centered at the mode of the posterior distribution, and its scale matrix is given by the negative inverse Hessian matrix evaluated at the mode. Specifically, for two blocks, the log posterior distribution is maximized with respect to  $\boldsymbol{\theta}_1$  given  $\boldsymbol{\theta}_2$ . This process is repeated at each iteration of the algorithm because  $\boldsymbol{\theta}_2$  changes at different stages. As a result, the algorithm can be slow since the optimization process is computationally demanding (see [184, Chap. 7 and 9] for examples).

A sensible recommendation when performing M-H algorithm is to use a random walk proposal such that  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, c^2 \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Sigma}$  is the negative inverse Hessian matrix evaluated at the mode, that is, maximize with respect to all parameters, and set  $c \approx 2.4/\sqrt{\dim\{\boldsymbol{\theta}\}}$ , which is the most efficient scale compared to independent sampling [160, Chap. 12] and [162]. After some iterations of the algorithm, adjust the scale matrix  $\boldsymbol{\Sigma}$  as before, and increase or decrease  $c$  if the acceptance rate of the simulations is too high or low, respectively. The objective is to bring the acceptance rate to the stated rule of thumb, that is, if the dimension is less than or equal to 2, the acceptance rate should be around 0.50, and if the dimension is greater than 2, the acceptance rate should be around 0.25. Once this is achieved, we should run the algorithm without modifications and use this part of the algorithm to perform inference.

#### **Example: Ph.D. students sleeping hours continues**

In the Ph.D. students sleeping hours exercise of Chapter 3 we get a posterior distribution that is Beta with parameters 16.55 and 39.57. We can sample from this posterior distribution using the function `rbeta` from **R**. However, we want to compare the performance of a M-H algorithm using as proposal density a  $U(0, 1)$  distribution.

The following code shows how to do a M-H algorithm to sample from the beta distribution using the uniform distribution.

**R code. Metropolis-Hastings algorithm: Ph.D. students sleeping hours**

---

```

1 rm(list = ls()); set.seed(010101)
2 an <- 16.55; bn <- 39.57
3 S <- 100000; p <- runif(S); accept <- rep(0, S)
4 for (s in 2:S){
5   pc <- runif(1) # Candidate
6   a <- dbeta(pc, an, bn)/dbeta(p[s-1], an, bn) # Acceptance
    rate
7   U <- runif(1)
8   if(U <= a){
9     p[s] <- pc
10    accept[s] <- 1
11  }else{
12    p[s] <- p[s-1]
13    accept[s] <- 0
14  }
15 }
16 mean(accept); mean(p); sd(p)
17 an/(an + bn); (an*bn/((an+bn)^2*(an+bn+1)))^0.5 # Population
  values
18 h <- hist(p, breaks=50, col="blue", xlab="Proportion Ph.D.
  students sleeping at least 6 hours", main="Beta draws
  from a Metropolis-Hastings algorithm")
19 pfit <- seq(min(p),max(p),length=50)
20 yfit<-dbeta(pfit, an, bn)
21 yfit <- yfit*diff(h$mids[1:2])*length(p)
22 lines(pfit, yfit, col="red", lwd=2)

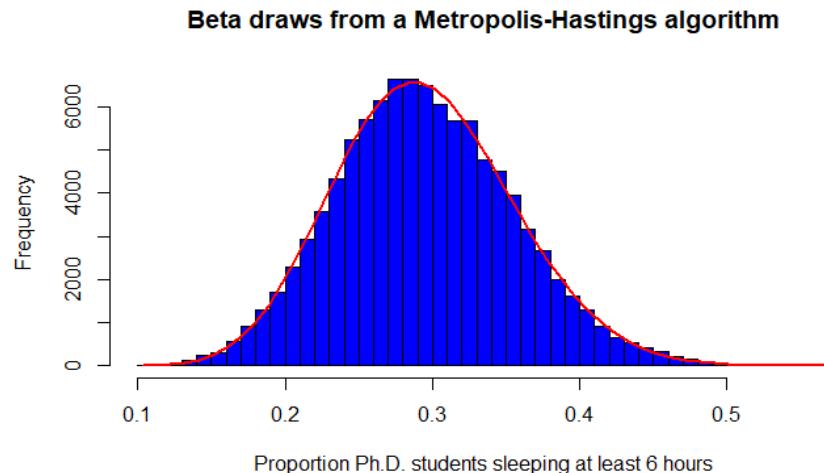
```

---

The results indicate that the mean and standard deviation obtained from the posterior draws are similar to the population values. Furthermore, Figure 4.2 presents the histogram of the posterior draws alongside the density of the beta distribution, demonstrating a good match between them.

#### 4.1.3 Hamiltonian Monte Carlo

Hamiltonian Monte Carlo (HMC) was proposed by [130] and later introduced to the statistical community by [316]. HMC extends the Metropolis algorithm to efficiently explore the parameter space by introducing *momentum variables*, which help overcome the random walk behavior of Gibbs sampling and the Metropolis-Hastings algorithm. Known also as hybrid Monte Carlo, HMC is particularly advantageous for high-dimensional posterior distributions, as it reduces the risk of getting stuck in local modes and significantly improves mixing [317].



**FIGURE 4.2**

Histogram of posterior draws of beta distribution and the density of the beta distribution.

However, HMC is designed to work with strictly positive target densities. Therefore, transformations are required to handle bounded parameters, such as variances and proportions. For example, logarithmic and logit transformations can be applied. These transformations necessitate the use of the change-of-variable theorem to compute the log posterior density and its gradient, which are essential for implementing the HMC algorithm.

HMC leverages concepts from physics, specifically Hamiltonian mechanics, to propose transitions in the Markov chain. In Hamiltonian mechanics, two key variables define the total energy of the system: the *position* ( $\theta$ ) and the *momentum* ( $\delta$ ). The Hamiltonian represents the total energy of the system, consisting of *potential energy* (energy due to position) and *kinetic energy* (energy associated with motion). The objective is to identify trajectories that preserve the system's total energy, meaning the Hamiltonian remains invariant, while avoiding trajectories that do not. This approach enhances the acceptance rate of proposed transitions.

To implement HMC, we solve the differential equations derived from the Hamiltonian, which involve derivatives with respect to position and momentum. However, these equations rarely have analytical solutions, requiring numerical methods for approximation. This necessitates discretizing Hamilton's equations, which introduces errors. To mitigate these errors, HMC uses the *leapfrog integrator*, a numerical method with smaller errors compared to simpler approaches like the Euler method.

HMC uses a *momentum variable* ( $\delta_k$ ) for each  $\theta_k$ , so that the transition

kernel of  $\boldsymbol{\theta}$  is determined by  $\boldsymbol{\delta}$ . Both vectors are updated using a Metropolis algorithm at each stage such that the distribution of  $\boldsymbol{\theta}$  remains invariant [317]. The joint density in HMC is given by  $p(\boldsymbol{\theta}, \boldsymbol{\delta} \mid \mathbf{y}) = \pi(\boldsymbol{\theta} \mid \mathbf{y}) \times p(\boldsymbol{\delta})$ , where  $\boldsymbol{\delta} \sim N(\mathbf{0}, \mathbf{M})$ , and  $\mathbf{M}$  is a diagonal matrix such that  $\delta_k \sim N(0, M_{kk})$ .

Algorithm A3 outlines the HMC implementation. The gradient vector  $\frac{d \log(\pi(\boldsymbol{\theta} \mid \mathbf{y}))}{d\boldsymbol{\theta}}$  must be computed analytically, as using finite differences can be computationally expensive. However, it is advisable to verify the analytical calculations by evaluating the gradient at the maximum posterior estimate, where the function should return values close to 0, or by comparing results with finite differences at a few points.

---

**Algorithm A3** Hamiltonian Monte Carlo

---

```

1: Initiate at  $\boldsymbol{\theta}^{(0)}$  in the support of  $\pi(\boldsymbol{\theta} \mid \mathbf{y})$ , and set step size  $\epsilon$ , number of
   leapfrog steps  $L$ , and total iterations  $S$ 
2: Draw  $\boldsymbol{\delta}^{(0)}$  from  $N(\mathbf{0}, \mathbf{M})$ 
3: for  $s = 1, \dots, S$  do
4:   for  $l = 1, \dots, L$  do
5:     if  $l = 1$  then
6:        $\boldsymbol{\delta}^c \leftarrow \boldsymbol{\delta}^{(s-1)} + \frac{1}{2}\epsilon \frac{d \log(\pi(\boldsymbol{\theta} \mid \mathbf{y}))}{d\boldsymbol{\theta}}$ 
7:        $\boldsymbol{\theta}^c \leftarrow \boldsymbol{\theta}^{(s-1)} + \epsilon \mathbf{M}^{-1} \boldsymbol{\delta}^c$ 
8:     else
9:       if  $l=2,\dots,L-1$  then
10:         $\boldsymbol{\delta}^c \leftarrow \boldsymbol{\delta}^c + \epsilon \frac{d \log(\pi(\boldsymbol{\theta} \mid \mathbf{y}))}{d\boldsymbol{\theta}}$ 
11:         $\boldsymbol{\theta}^c \leftarrow \boldsymbol{\theta}^c + \epsilon \mathbf{M}^{-1} \boldsymbol{\delta}^c$ 
12:      else
13:         $\boldsymbol{\delta}^c \leftarrow \boldsymbol{\delta}^c + \frac{1}{2}\epsilon \frac{d \log(\pi(\boldsymbol{\theta} \mid \mathbf{y}))}{d\boldsymbol{\theta}}$ 
14:         $\boldsymbol{\theta}^c \leftarrow \boldsymbol{\theta}^c + \epsilon \mathbf{M}^{-1} \boldsymbol{\delta}^c$ 
15:      end if
16:    end if
17:  end for
18:  Calculate  $\alpha([\boldsymbol{\theta} \; \boldsymbol{\delta}]^{(s-1)}, [\boldsymbol{\theta} \; \boldsymbol{\delta}]^c) = \min \left\{ \frac{p(\boldsymbol{\delta}^c) \pi(\boldsymbol{\theta}^c \mid \mathbf{y})}{p(\boldsymbol{\delta}^{(s-1)}) \pi(\boldsymbol{\theta}^{(s-1)} \mid \mathbf{y})}, 1 \right\}$ 
19:  Draw  $U$  from  $U(0, 1)$ 
20:   $\boldsymbol{\theta}^{(s)} = \begin{cases} \boldsymbol{\theta}^c & \text{if } U \leq \alpha(\cdot, \cdot) \\ \boldsymbol{\theta}^{(s-1)} & \text{otherwise} \end{cases}$ 
21: end for

```

---

Note that HMC does not require the marginal likelihood, as neither the gradient vector  $\frac{d \log(\pi(\boldsymbol{\theta} \mid \mathbf{y}))}{d\boldsymbol{\theta}}$  nor the acceptance rate depend on it. That is, we can use only  $\pi(\boldsymbol{\theta}) \times p(\mathbf{y} \mid \boldsymbol{\theta})$  to implement HMC. In addition, we do not retain  $\boldsymbol{\delta}$  after it is updated at the beginning of each iteration, as it is not required subsequently. To begin, the step size ( $\epsilon$ ) can be drawn randomly from a uniform distribution between 0 and  $2\epsilon_0$ , and the number of leapfrog steps ( $L$ ) is set as the largest integer near  $1/\epsilon$ , ensuring  $\epsilon \times L \approx 1$ . We need to

set  $\mathbf{M}$  to be the inverse of the posterior covariance matrix evaluated at the maximum a posteriori estimate under this setting.

The acceptance rate should be checked, with the optimal rate around 65% [160, Chap. 12]. If the acceptance rate is much higher than 65%, increase  $\epsilon_0$ ; if it is much lower, decrease it. This strategy may not always work, and alternative strategies can be tested, such as setting  $\mathbf{M} = \mathbf{I}$  and fine-tuning  $\epsilon$  and  $L$  to achieve an acceptance rate near 65%. Finally, the number of iterations ( $S$ ) is chosen to ensure convergence to the stationary distribution.

#### **Example: Sampling from a bi-variate Gaussian distribution**

As a toy example, let's compare the Gibbs sampling, M-H and HMC algorithms when the posterior distribution is a bi-variate Gaussian distribution with mean  $\mathbf{0}$  and covariance matrix  $\Sigma = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ . Let's set  $\rho = 0.98$ .

The Gibbs sampler requires the conditional posterior distributions, which in this case are  $\theta_1 | \theta_2 \sim N(\rho\theta_2, 1 - \rho^2)$  and  $\theta_2 | \theta_1 \sim N(\rho\theta_1, 1 - \rho^2)$ . We use the random walk proposal distribution for the M-H algorithm, where  $\boldsymbol{\theta}^c \sim N(\boldsymbol{\theta}^{(s-1)}, \text{diag}\{0.18^2\})$ . We set  $\epsilon = 0.05$ ,  $L = 20$  and  $\mathbf{M} = \mathbf{I}_2$  for the HMC algorithm, and given that  $\pi(\boldsymbol{\theta} | \mathbf{y}) \propto \exp\{-\frac{1}{2}\boldsymbol{\theta}^\top \Sigma^{-1} \boldsymbol{\theta}\}$ , then  $\frac{d\log(\pi(\boldsymbol{\theta} | \mathbf{y}))}{d\boldsymbol{\theta}} = -\Sigma^{-1}\boldsymbol{\theta}$ .

The following code shows how to implement the Gibbs sampler, the random walk M-H algorithm, and the HMC in this example such that the effective number of posterior draws is 400.

#### **R code. Gibbs, M-H and HMC: Bi-variate normal distribution**

```

1 rm(list = ls()); set.seed(010101)
2 # Gibbs sampler
3 Gibbs <- function(theta, rho){
4   thetal <- rnorm(1, mean = rho*theta, sd = (1- rho^2)^0.5)
5   return(thetal)
6 }
```

### R code. Gibbs, M-H and HMC: Bi-variate normal distribution

---

```

1 # Metropolis-Hastings
2 MH <- function(theta, rho, sig2){
3   SIGMA <- matrix(c(1, rho, rho, 1), 2, 2)
4   SIGMAd <- matrix(c(1, sig2, sig2, 1), 2, 2)
5   thetac <- MASS::mvrnorm(1, mu = theta, Sigma = SIGMAd)
6   a <- mvtnorm::dmvnorm(theta, c(0, 0), SIGMA)/mvtnorm::
7     dmvnorm(theta, c(0, 0), SIGMA)
8   U <- runif(1)
9   if(U <= a){
10     theta <- thetac
11     accept <- 1
12   }else{
13     theta <- theta
14     accept <- 0
15   }
16   return(list(theta = theta, accept = accept))
17 }
18 # Hamiltonian Monte Carlo
19 HMC <- function(theta, rho, epsilon, M){
20   SIGMA <- matrix(c(1, rho, rho, 1), 2, 2)
21   L <- ceiling(1/epsilon)
22   Minv <- solve(M); thetac <- theta
23   K <- length(theta)
24   mom <- t(mvtnorm::rmvnorm(1, rep(0, K), M))
25   logPost_Mom_t <- mvtnorm::dmvnorm(t(theta), rep(0, K),
26     SIGMA, log = TRUE) + mvtnorm::dmvnorm(t(mom), rep(0, K),
27     , M, log = TRUE)
28   for(l in 1:L){
29     if(l == 1 | l == L){
30       mom <- mom + 0.5*epsilon*(-solve(SIGMA)%*%theta)
31       theta <- theta + epsilon*Minv%*%mom
32     }else{
33       mom <- mom + epsilon*(-solve(SIGMA)%*%theta)
34       theta <- theta + epsilon*Minv%*%mom
35     }
36   }
37   logPost_Mom_star <- mvtnorm::dmvnorm(t(theta), rep(0, K),
38     SIGMA, log = TRUE) + mvtnorm::dmvnorm(t(mom), rep(0, K),
39     , M, log = TRUE)
40   alpha <- min(1, exp(logPost_Mom_star-logPost_Mom_t))
41   u <- runif(1)
42   if(u <= alpha){
43     thetaNew <- c(theta)
44   }else{
45     thetaNew <- thetac
46   }
47   rest <- list(theta = thetaNew, Prob = alpha)
48   return(rest)
49 }
```

---

### R code. Gibbs, M-H and HMC: Bi-variate normal distribution

```

1 # Hyperparameters
2 rho <- 0.98; sig2 <- 0.18^2
3 # Posterior draws Gibbs and M-H
4 S <- 8000; thin <- 20; K <- 2
5 thetaPostGibbs <- matrix(NA, S, K)
6 thetaPostMH <- matrix(NA, S, K)
7 AcceptMH <- rep(NA, S)
8 thetaGibbs <- c(-2, 3); thetaMH <- c(-2, 3)
9 for(s in 1:S){
10   theta1 <- Gibbs(thetaGibbs[2], rho)
11   theta2 <- Gibbs(theta1, rho)
12   thetaGibbs <- c(theta1, theta2)
13   ResMH <- MH(thetaMH, rho, sig2)
14   thetaMH <- ResMH$theta
15   thetaPostGibbs[s,] <- thetaGibbs
16   thetaPostMH[s,] <- thetaMH
17   AcceptMH[s] <- ResMH$accept
18 }
19 keep <- seq(0, S, thin)
20 mean(AcceptMH[keep[-1]])
21 thetaPostGibbsMCMC <- coda::mcmc(thetaPostGibbs[keep[-1],])
22 summary(thetaPostGibbsMCMC)
23 coda::autocorr.plot(thetaPostGibbsMCMC)
24 thetaPostMHMCMC <- coda::mcmc(thetaPostMH[keep[-1],])
25 plot(thetaPostMHMCMC)
26 coda::autocorr.plot(thetaPostMHMCMC)
27 # Posterior draws HMC
28 S <- 400; epsilon <- 0.05; L <- ceiling(1/epsilon); M <-
29   diag(2)
30 thetaPostHMC <- matrix(NA, S, K)
31 ProbAcceptHMC <- rep(NA, S)
32 thetaHMC <- c(-2, 3)
33 for(s in 1:S){
34   ResHMC <- HMC(theta = thetaHMC, rho, epsilon, M)
35   thetaHMC <- ResHMC$theta
36   thetaPostHMC[s,] <- thetaHMC
37   ProbAcceptHMC[s] <- ResHMC$Prob
38 }
39 thetaPostHMCMCMC <- coda::mcmc(thetaPostHMC)
40 plot(thetaPostHMCMCMC); coda::autocorr.plot(thetaPostHMCMCMC)
41 summary(ProbAcceptHMC)
42 #Figure
43 df <- as.data.frame(cbind(1:S, thetaPostHMC[,1], thetaPostMH
44   [keep[-1],1], thetaPostGibbs[keep[-1],1]))
45 colnames(df) <- c("Iter", "HMC", "MH", "Gibbs")

```

**R code. Gibbs, M-H and HMC: Bi-variate normal distribution**

```

1 library(latex2exp); library(ggpubr)
2 g1 <- ggplot(df, aes(x= Iter)) + geom_point(aes(y=HMC),
  colour="black") + labs(x = "Iteration", y = TeX("$\\theta_1$"),
  title = "HMC algorithm")
3 g2 <- ggplot(df, aes(x= Iter)) + geom_point(aes(y=MH),
  colour="black") + labs(x = "Iteration", y = TeX("$\\theta_1$"),
  title = "M-H algorithm")
4 g3 <- ggplot(df, aes(x= Iter)) + geom_point(aes(y=Gibbs),
  colour="black") + labs(x = "Iteration", y = TeX("$\\theta_1$"),
  title = "Gibbs sampling")
5 ggarrange(g3, g2, g1, labels = c("A", "B", "C"), ncol = 3,
  nrow = 1)

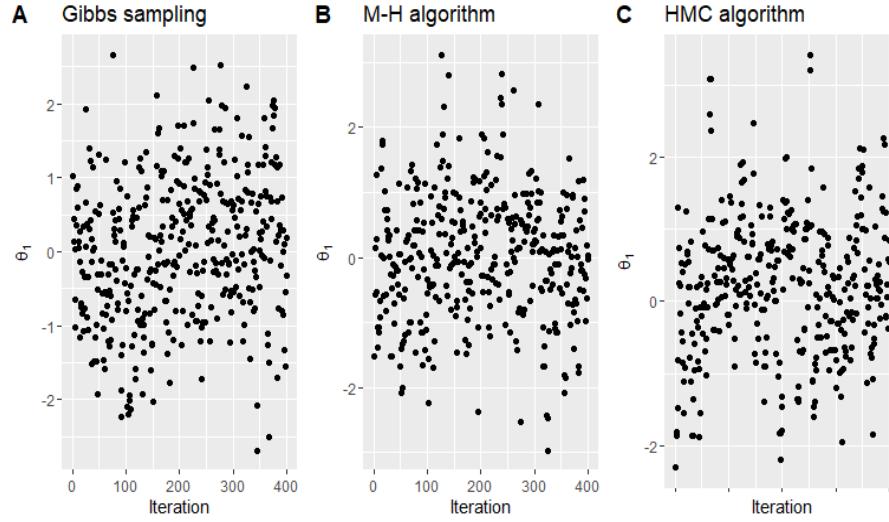
```

Figure 4.3 shows the posterior draws of  $\theta_1$  using the Gibbs sampler (Panel A, left), the Metropolis-Hastings algorithm (Panel B, middle), and the Hamiltonian Monte Carlo (Panel C, right). The convergence diagnostic plots (no shown) suggests that the three algorithms perform a good job. Although, the acceptance rate in HMC is higher than the M-H due to the HMC producing larger changes in  $\theta$  than a corresponding number of random-walk M-H iterations [317].

## 4.2 Importance sampling

Up to this section, we have introduced MCMC methods for sampling from the posterior distribution when it does not have a standard closed form. However, MCMC methods have some limitations. First, the samples are generated sequentially, which complicates parallel computing. Although multiple MCMC chains can be run simultaneously, this approach—often referred to as brute-force parallelization—does not fully address the sequential nature of individual chains. Second, consecutive samples are correlated, which reduces the effective sample size and complicates convergence diagnostics.

Thus, in this section, we introduce *importance sampling* (IS), a simulation method for drawing samples from the posterior distribution that avoids these limitations. Unlike MCMC, IS does not require satisfying the balancing condition, making it conceptually and mathematically simpler to implement in certain situations. Moreover, importance weights can be reused to analyze posterior quantities, compute marginal likelihoods, compare models, approx-

**FIGURE 4.3**

Posterior draws: Gibbs sampler (A), Metropolis-Hastings (B) and Hamiltonian Monte Carlo (C) in the bivariate normal example,  $\rho = 0.98$ .

imate new target distributions, and allow for straightforward parallelization in large-scale problems.

However, the critical challenge in IS lies in selecting an appropriate proposal distribution. This involves satisfying both support and stability conditions, which can be difficult to achieve, particularly in high-dimensional problems. In such cases, MCMC methods may be more suitable.

The starting point is evaluating the integral:

$$\mathbb{E}_\pi[h(\boldsymbol{\theta})] = \int_{\Theta} h(\boldsymbol{\theta})\pi(\boldsymbol{\theta} | \mathbf{y})d\boldsymbol{\theta}, \quad (4.1)$$

where  $\mathbb{E}_\pi$  denotes expected value under the posterior distribution. Thus, we can approximate Equation 4.1 by

$$\bar{h}(\boldsymbol{\theta})_S = \frac{1}{S} \sum_{s=1}^S h(\boldsymbol{\theta}^{(s)}), \quad (4.2)$$

where  $\boldsymbol{\theta}^{(s)}$  are draws from  $\pi(\boldsymbol{\theta} | \mathbf{y})$ . The *strong law of large numbers* shows that  $\bar{h}(\boldsymbol{\theta})_S$  converges (almost surely) to  $\mathbb{E}_\pi[h(\boldsymbol{\theta})]$  as  $S \rightarrow \infty$ .

The challenge arises when we do not know how to obtain samples from  $\pi(\boldsymbol{\theta} | \mathbf{y})$ . The ingenious idea is to express Equation 4.1 in a different way using the *importance sampling fundamental identity* [372, Chap. 3]:

$$\mathbb{E}_\pi[h(\boldsymbol{\theta})] = \int_{\Theta} h(\boldsymbol{\theta})\pi(\boldsymbol{\theta} | \mathbf{y})\frac{q(\boldsymbol{\theta})}{q(\boldsymbol{\theta})}d\boldsymbol{\theta}$$

$$= \mathbb{E}_q \left[ \frac{h(\boldsymbol{\theta})\pi(\boldsymbol{\theta} | \mathbf{y})}{q(\boldsymbol{\theta})} \right], \quad (4.3)$$

where  $q(\boldsymbol{\theta})$  is the proposal distribution.

Thus, we have

$$\frac{1}{S} \sum_{s=1}^S \left[ \frac{h(\boldsymbol{\theta}^{(s)})\pi(\boldsymbol{\theta}^{(s)} | \mathbf{y})}{q(\boldsymbol{\theta}^{(s)})} \right] = \frac{1}{S} \sum_{s=1}^S h(\boldsymbol{\theta}^{(s)})w(\boldsymbol{\theta}^{(s)}),$$

where  $w(\boldsymbol{\theta}^{(s)}) = \left[ \frac{\pi(\boldsymbol{\theta}^{(s)} | \mathbf{y})}{q(\boldsymbol{\theta}^{(s)})} \right]$  are called the *importance weights*, and  $\boldsymbol{\theta}^{(s)}$  are samples from the proposal distribution. This expression converges to  $\mathbb{E}_{\pi}[h(\boldsymbol{\theta})]$  given that the support of  $q(\boldsymbol{\theta})$  includes the support of  $\pi(\boldsymbol{\theta}^{(s)} | \mathbf{y})$ .

There are many proposal distributions that satisfy the support condition. However, the stability of the method depends heavily on the variability of the importance weights. In particular, the variance of

$$\frac{1}{S} \sum_{s=1}^S h(\boldsymbol{\theta}^{(s)})w(\boldsymbol{\theta}^{(s)})$$

can be large if the proposal distribution has lighter tails than the posterior distribution. In this case, the weights  $w(\boldsymbol{\theta}^{(s)})$  will vary widely, assigning too much importance to a few values of  $\boldsymbol{\theta}^{(s)}$ . Thus, it is important to use proposals that have thicker tails than the posterior distribution. In any case, we should check the adequacy of the proposal distribution by analyzing the behavior of the importance weights. If they are distributed more or less uniformly over the support, it is a good sign. Consider, for instance, the extreme case where  $q(\boldsymbol{\theta}) = \pi(\boldsymbol{\theta} | \mathbf{y})$ , then  $w(\boldsymbol{\theta}^{(s)}) = 1$  everywhere.

A natural choice in Bayesian inference is to use the prior distribution as the proposal, given that it is a proper density function. The prior distribution typically has heavier tails than the posterior by construction, and it is usually a distribution that allows for easy sampling.

The most relevant point for us is that importance sampling provides a way to simulate from the posterior distribution when there is no closed-form solution. The method generates samples  $\boldsymbol{\theta}^{(s)}$  from  $q(\boldsymbol{\theta})$  and computes the importance weights  $w(\boldsymbol{\theta}^{(s)})$ . Thus, if we *resample* with replacement from  $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(S)}$ , selecting  $\boldsymbol{\theta}^{(s)}$  with probability proportional to  $w(\boldsymbol{\theta}^{(s)})$ , we would get a sample  $\boldsymbol{\theta}^{*(1)}, \boldsymbol{\theta}^{*(2)}, \dots, \boldsymbol{\theta}^{*(L)}$  of size  $L$  from  $\pi(\boldsymbol{\theta} | \mathbf{y})$  [416, 388]. This is named *sampling/importance resampling* (SIR) algorithm. Observe that the number of times  $L^{(s)}$  each particular point  $\boldsymbol{\theta}^{(s)}$  is selected follows a binomial distribution with size  $L$ , and probabilities proportional to  $w^{(s)}$ . Consequently, the vector  $\mathbf{L}_{\boldsymbol{\theta}} = \{L_{\boldsymbol{\theta}^1}, L_{\boldsymbol{\theta}^2}, \dots, L_{\boldsymbol{\theta}^S}\}$  follows a multinomial distribution with  $L$  trials and probabilities proportional to  $w(\boldsymbol{\theta}^{(s)})$ ,  $s = 1, 2, \dots, S$  [66]. Therefore, the resampling step ensures that points in the first-stage sample with small importance weights are more likely to be discarded, while points with high weights are replicated in proportion to their importance weights. In most applications, it is typical to have  $S \gg L$ .

The intuition is that importance weights are scaling factors that correct for the bias introduced by drawing from  $q(\boldsymbol{\theta}^{(s)})$  instead of  $\pi(\boldsymbol{\theta}^{(s)} \mid \mathbf{y})$ ; thus, when combined, the samples and weights effectively recreate the posterior distribution, ensuring the resampled data set reflects the posterior. Let's proof this

$$\begin{aligned} P(\boldsymbol{\theta}^* \in A) &= \frac{1}{S} \sum_{s=1}^S w^{(s)} \mathbb{1}_A(\boldsymbol{\theta}^{(s)}) \\ &\rightarrow \mathbb{E}_q \left[ \mathbb{1}_{\in A}(\boldsymbol{\theta}) \frac{\pi(\boldsymbol{\theta} \mid \mathbf{y})}{q(\boldsymbol{\theta})} \right] \\ &= \int_A \left[ \frac{\pi(\boldsymbol{\theta} \mid \mathbf{y})}{q(\boldsymbol{\theta})} \right] q(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \int_A \pi(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}. \end{aligned}$$

Thus,  $\boldsymbol{\theta}^*$  is approximately distributed as an observation from  $\pi(\boldsymbol{\theta} \mid \mathbf{y})$ .

However, the weights  $\pi(\boldsymbol{\theta}^{(s)} \mid \mathbf{y})/(Sq(\boldsymbol{\theta}^{(s)}))$  do not sum up to 1, and we need to standardize them,

$$w^*(\boldsymbol{\theta}^{(s)}) = \frac{\frac{1}{S} w(\boldsymbol{\theta}^{(s)})}{\frac{1}{S} \sum_{s=1}^S w(\boldsymbol{\theta}^{(s)})}.$$

Note that we could alternatively arrive to these weights as follow

$$\begin{aligned} \mathbb{E}_\pi[h(\boldsymbol{\theta})] &= \int_\Theta \left[ \frac{h(\boldsymbol{\theta})\pi(\boldsymbol{\theta} \mid \mathbf{y})}{q(\boldsymbol{\theta})} \right] q(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \frac{\int_\Theta \left[ \frac{h(\boldsymbol{\theta})\pi(\boldsymbol{\theta} \mid \mathbf{y})}{q(\boldsymbol{\theta})} \right] q(\boldsymbol{\theta}) d\boldsymbol{\theta}}{\int_\Theta \left[ \frac{\pi(\boldsymbol{\theta} \mid \mathbf{y})}{q(\boldsymbol{\theta})} \right] q(\boldsymbol{\theta}) d\boldsymbol{\theta}}. \end{aligned}$$

Then,

$$\frac{\frac{1}{S} \sum_{s=1}^S h(\boldsymbol{\theta}^{(s)}) w(\boldsymbol{\theta}^{(s)})}{\frac{1}{S} \sum_{s=1}^S w(\boldsymbol{\theta}^{(s)})} = \sum_{s=1}^S h(\boldsymbol{\theta}^{(s)}) w^*(\boldsymbol{\theta}^{(s)}).$$

This alternative expression also converges (almost surely) to  $\mathbb{E}_\pi[h(\boldsymbol{\theta})]$ . In addition, this expression is very useful because if we do not have the marginal likelihood in the posterior distribution, this constant cancels out in  $w^*(\boldsymbol{\theta}^{(s)})$ . And, although this estimator is biased, the bias is small, and provides good gains in variance reduction compared with the non-standardized option [372, Chap. 3].

A nice by-product of implementing IS is that it easily allows the calculation of the marginal likelihood. In particular, we know from Bayes' rule that

$$p(\mathbf{y})^{-1} = \frac{\pi(\boldsymbol{\theta} \mid \mathbf{y})}{p(\mathbf{y} \mid \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})},$$

then,

$$\begin{aligned}\int_{\Theta} p(\boldsymbol{y})^{-1} q(\boldsymbol{\theta}) d\boldsymbol{\theta} &= \int_{\Theta} \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{y} | \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})} \pi(\boldsymbol{\theta} | \boldsymbol{y}) d\boldsymbol{\theta} \\ &= \mathbb{E}_{\pi} \left[ \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{y} | \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})} \right].\end{aligned}$$

Thus, an estimate of the marginal likelihood is  $\left[ \frac{1}{S} \sum_{s=1}^S \frac{q(\boldsymbol{\theta}^{*(s)})}{p(\boldsymbol{y} | \boldsymbol{\theta}^{*(s)}) \times \pi(\boldsymbol{\theta}^{*(s)})} \right]^{-1}$ . This is the Gelfand-Dey method to calculate the marginal likelihood [157] (see subsection 10.4.3).

### Example: Cauchy distribution

Let's assume that the posterior distribution is Cauchy with parameters 0 and 1. We perform an importance sampling algorithm using as proposals a standard normal distribution and a Student's t distribution with 3 degrees of freedom. The following code shows how to do this.

#### *R code. Importance sampling: Cauchy distribution*

```

1 rm(list = ls()); set.seed(010101)
2 S <- 20000 # Size proposal
3 # Importance sampling from standard normal proposal
4 thetaNs <- rnorm(S)
5 wNs <- dcauchy(thetaNs)/dnorm(thetaNs)
6 wNstars <- wNs/sum(wNs)
7 L <- 10000 # Size posterior
8 thetaCauchyN <- sample(thetaNs, L, replace = TRUE, prob =
  wNstars)
9 h <- hist(thetaCauchyN, breaks=50, col="blue", xlab="x",
  main="Cauchy draws from importance sampling: Normal
  standard proposal")
10 pfit <- seq(min(thetaCauchyN),max(thetaCauchyN),length=50)
11 yfit<-dcauchy(pfit)
12 yfit <- yfit*diff(h$mid[1:2])*length(thetaCauchyN)
13 lines(pfit, yfit, col="red", lwd=2)
14 # Importance sampling from Student's t proposal
15 df <- 3
16 thetaTs <- rt(S, df = df)
17 wTs <- dcauchy(thetaTs)/dt(thetaTs, df = df)
18 wTstars <- wTs/sum(wTs)
19 thetaCauchyT <- sample(thetaTs, L, replace = TRUE, prob =
  wTstars)
```

**R code. Importance sampling: Cauchy distribution**

```

1 h <- hist(thetaCauchyT, breaks=50, col="blue", xlab="x",
            main="Cauchy draws from importance sampling: Student's t
            proposal")
2 pfit <- seq(min(thetaCauchyT),max(thetaCauchyT),length=50)
3 yfit<-dcauchy(pfit)
4 yfit <- yfit*diff(h$mids[1:2])*length(thetaCauchyT)
5 lines(pfit, yfit, col="red", lwd=2)
6 plot(wNstars, main = "Importance sampling: Cauchy
        distribution", ylab = "Weights", xlab = "Iterations")
7 points(wTstars, col = "blue")
8 legend("topright", legend = c("Normal", "Student's t"), col
        = c("black", "blue"), pch = c(1, 1))
9

```

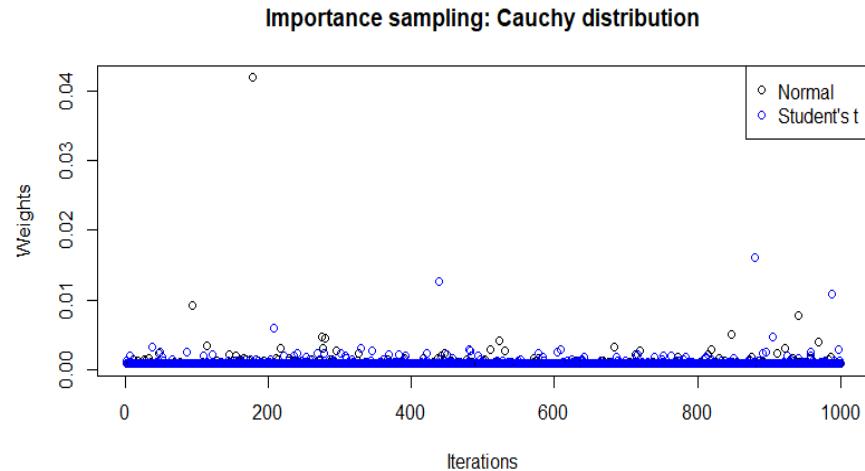
Figure 4.4 shows the weights using the standard normal distribution (black dots) and the Student's t-distribution with 3 degrees of freedom (blue dots) as proposals. We observe that a few draws carry too much weight when using the normal proposal; this occurs because the normal distribution has much lighter tails compared to the Cauchy distribution. In contrast, using the Student's t-distribution with 3 degrees of freedom improves this situation.

Figures 4.5 and 4.6 show the histograms of the posterior draws using the normal and Student's t-distributions, respectively, along with the density of the Cauchy distribution. The spike in the posterior draws from the standard normal proposal arises due to the lighter tails of the standard normal compared to the Cauchy distribution, consequently assigning too much weight to a specific draw from the normal distribution.

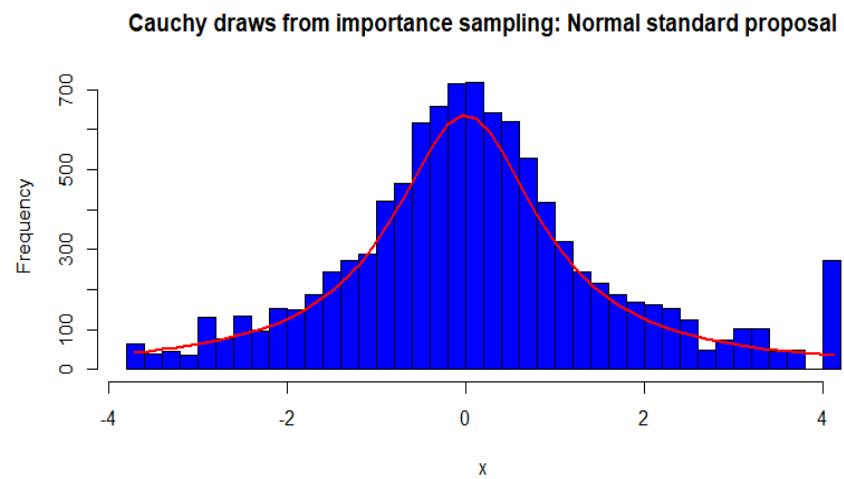
---

### 4.3 Particle filtering

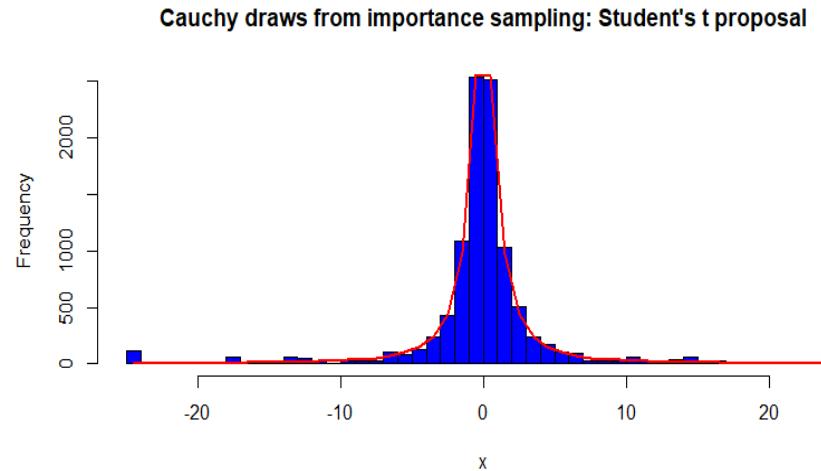
Now, we consider the scenario where we need to sample from a posterior distribution whose dimension increases over time,  $\pi(\boldsymbol{\theta}_{0:t} | \mathbf{y}_{0:t})$ , for  $t = 0, 1, \dots$ . The challenge arises from the fact that, even if this posterior distribution is known, the computational complexity of implementing a sampling scheme in this context increases linearly with  $t$ . This makes MCMC methods, which operate in batch mode and require a complete re-run whenever new information becomes available, less optimal. Consequently, we present sequential algorithms, which operate incrementally as new data becomes available, and are often a better alternative. These algorithms are typically faster and are

**FIGURE 4.4**

Importance sampling: 1000 weights in the Cauchy distribution example using the standard normal and student's t with 3 degrees of freedom as proposals.

**FIGURE 4.5**

Importance sampling: Draws from a Cauchy distribution using the standard normal as proposal.

**FIGURE 4.6**

Importance sampling: Draws of a Cauchy distribution using the Student's t with 3 degrees of freedom as proposal.

well-suited for scenarios requiring real-time updates, commonly referred to as online mode.

Specifically, we consider the dynamic system in the *state-space* representation. This is a system where there is an *unobservable state vector*  $\boldsymbol{\theta}_t \in \mathbb{R}^K$ , and an observed variable  $\mathbf{Y}_t$ ,  $t = 0, 1, \dots$  such that (i)  $\boldsymbol{\theta}_t$  is a *Markov process*, this is,  $\pi(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{1:t-1}) = \pi(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1})$ ,  $t = 1, 2, \dots$ , all the relevant information to define  $\boldsymbol{\theta}_t$  is in  $\boldsymbol{\theta}_{t-1}$ ,<sup>1</sup> and ii)  $\mathbf{Y}_t \perp \mathbf{Y}_s | \boldsymbol{\theta}_t$ ,  $s < t$ , there is independence between observable variables regarding their history conditional on the actual state vector. We can see in Figure 4.7 a graphical representation of the dynamic system.

Formally,

$$\begin{aligned}\boldsymbol{\theta}_t &= h(\boldsymbol{\theta}_{t-1}, \mathbf{w}_t) && \text{(State equations)} \\ Y_t &= f(\boldsymbol{\theta}_t, \mu_t) && \text{(Observation equation),}\end{aligned}$$

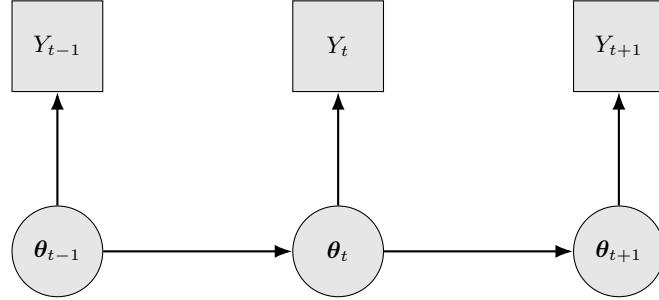
where  $\mathbf{w}_t$  and  $\mu_t$  are stochastic errors such that their probability distributions define the transition density  $\pi(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1})$  and observation density  $p(Y_t | \boldsymbol{\theta}_t)$ .

We present *particle filtering*, a specific case of *sequential Monte Carlo* (SMC), which is one of the most commonly used algorithms for scenarios requiring sequential updates of the posterior distribution as described by the *state-space* model.

---

<sup>1</sup> $\boldsymbol{\theta}_0$  comes from the given distribution  $\pi(\boldsymbol{\theta}_0)$ .

$$Y_t \sim p(Y_t | \theta_t)$$



$$\theta_t \sim \pi(\theta_t | \theta_{t-1})$$

**Notes:** The figure illustrates the structure of a *state-space* model where the latent states  $\theta_t$  evolve according to  $\pi(\theta_t | \theta_{t-1})$ , and the observations  $Y_t$  depend on the states via  $p(Y_t | \theta_t)$ .

**FIGURE 4.7**  
*State-space* model representation.

The starting point is *sequential importance sampling* (SIS), originally proposed by [191], which is a modification of IS to compute an estimate of  $\pi(\theta_{0:t} | y_{0:t})$  without altering the past trajectories  $\{\theta_{1:t-1}^{(s)}, s = 1, 2, \dots, S\}$ . The key idea is to use a proposal density that takes the form

$$\begin{aligned} q(\theta_{0:t} | y_{0:t}) &= q(\theta_{0:t-1} | y_{1:t-1})q(\theta_t | \theta_{t-1}, y_t) \\ &= q(\theta_0) \prod_{h=1}^t q(\theta_h | \theta_{h-1}, y_h). \end{aligned}$$

This proposal density allows calculating the weights sequentially,

$$\begin{aligned} w_t(\theta_{0:t}^{(s)}) &= \frac{\pi(\theta_{0:t}^{(s)} | y_{0:t})}{q(\theta_{0:t}^{(s)} | y_{0:t})} \\ &= \frac{p(y_{0:t} | \theta_{0:t}^{(s)})\pi(\theta_{0:t}^{(s)})}{p(y_{0:t})q(\theta_{0:t}^{(s)} | y_{0:t})} \\ &= \frac{p(y_t | \theta_t^{(s)})p(y_{1:t-1} | \theta_{0:t-1}^{(s)})\pi(\theta_t^{(s)} | \theta_{t-1}^{(s)})\pi(\theta_{0:t-1}^{(s)})}{p(y_{0:t})q(\theta_t^{(s)} | \theta_{t-1}^{(s)}, y_t)q(\theta_{0:t-1}^{(s)} | y_{1:t-1})} \\ &\propto \frac{p(y_{1:t-1} | \theta_{0:t-1}^{(s)})\pi(\theta_{0:t-1}^{(s)})}{q(\theta_{0:t-1}^{(s)} | y_{1:t-1})} \frac{p(y_t | \theta_t^{(s)})\pi(\theta_t^{(s)} | \theta_{t-1}^{(s)})}{q(\theta_t^{(s)} | \theta_{t-1}^{(s)}, y_t^{(s)})} \end{aligned}$$

$$\begin{aligned} & \propto w_{t-1}(\boldsymbol{\theta}_{0:t-1}^{(s)}) \frac{p(\mathbf{y}_t | \boldsymbol{\theta}_t^{(s)}) \pi(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1}^{(s)})}{q(\boldsymbol{\theta}_t^{(s)} | \boldsymbol{\theta}_{t-1}^{(s)}, \mathbf{y}_t)} \\ & \propto w_{t-1}^*(\boldsymbol{\theta}_{0:t-1}^{(s)}) \frac{p(\mathbf{y}_t | \boldsymbol{\theta}_t^{(s)}) \pi(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1}^{(s)})}{q(\boldsymbol{\theta}_t^{(s)} | \boldsymbol{\theta}_{t-1}^{(s)}, \mathbf{y}_t)}. \end{aligned}$$

Take into account that  $p(\mathbf{y}_{0:t})$  does not depend on  $\boldsymbol{\theta}_{0:t}^{(s)}$ . The term  $\alpha_t(\boldsymbol{\theta}_{0:t}^{(s)}) = \frac{p(\mathbf{y}_t | \boldsymbol{\theta}_t^{(s)}) \pi(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1}^{(s)})}{q(\boldsymbol{\theta}_t^{(s)} | \boldsymbol{\theta}_{t-1}^{(s)}, \mathbf{y}_t)}$  is called the *incremental importance weight*, and implies that

$$w_t(\boldsymbol{\theta}_{0:t}^{(s)}) = w_0(\boldsymbol{\theta}_0^{(s)}) \prod_{h=1}^t \alpha_h(\boldsymbol{\theta}_{1:h}^{(s)}).$$

This algorithm possesses the desirable property of maintaining fixed computational complexity. Consequently, we sequentially obtain draws  $\boldsymbol{\theta}_t^{(s)}$ , referred to as particles:  $\boldsymbol{\theta}_0^{(s)}$  is drawn from  $q(\boldsymbol{\theta}_0)$  at  $t = 0$ , and subsequently,  $\boldsymbol{\theta}_h^{(s)}$  is drawn from  $q(\boldsymbol{\theta}_h | \boldsymbol{\theta}_{h-1}, \mathbf{y}_h)$  at  $t = h$  [124, 66].

A relevant case is when the proposal distribution takes the form of the prior distribution, that is,  $q(\boldsymbol{\theta}_{0:t} | \mathbf{y}_{0:t}) = \pi(\boldsymbol{\theta}_{0:t}) = \pi(\boldsymbol{\theta}_0) \prod_{h=1}^t \pi(\boldsymbol{\theta}_h | \boldsymbol{\theta}_{h-1})$ . This implies that

$$w_t(\boldsymbol{\theta}_{0:t}^{(s)}) \propto w_{t-1}^*(\boldsymbol{\theta}_{0:t-1}^{(s)}) p(\mathbf{y}_t | \boldsymbol{\theta}_t^{(s)}),$$

this means that the *incremental importance weight* is given by  $p(\mathbf{y}_t | \boldsymbol{\theta}_t^{(s)})$ .

Algorithm A4 shows how to perform SIS [66]. We set  $w_t^{(s)} := w_t(\boldsymbol{\theta}_{0:t}^{(s)})$  to simplify notation.

---

**Algorithm A4** Sequential importance sampling algorithm

---

```

1: for  $s = 1, \dots, S$  do
2:   Sample  $\boldsymbol{\theta}_0^{(s)}$  from  $q(\boldsymbol{\theta}_0 | y_0)$ 
3:   Calculate the importance weights  $w_0^{(s)} \propto \frac{p(y_0 | \boldsymbol{\theta}_0^{(s)}) \pi(\boldsymbol{\theta}_0^{(s)})}{q(\boldsymbol{\theta}_0^{(s)} | y_0)}$ 
4: end for
5: for  $t = 1, \dots, T$  do
6:   for  $s = 1, \dots, S$  do
7:     Draw particles  $\boldsymbol{\theta}_t^{(s)}$  from  $q_t(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1}, \mathbf{y}_t)$ 
8:     Compute the weights  $w_t^{(s)} = w_{t-1}^{*(s)} \frac{p(\mathbf{y}_t | \boldsymbol{\theta}_t^{(s)}) \pi(\boldsymbol{\theta}_t^{(s)} | \boldsymbol{\theta}_{t-1}^{(s)})}{q(\boldsymbol{\theta}_t^{(s)} | \boldsymbol{\theta}_{t-1}^{(s)}, \mathbf{y}_t)}$ 
9:   end for
10:  Standardize the weights  $w_t^{*(s)} = \frac{w_t^{(s)}}{\sum_{h=1}^S w_t^{(h)}}, s = 1, 2, \dots, S$ 
11: end for

```

---

**Example: Dynamic linear model**

Let's assume that the state-space representation is

$$\begin{aligned}\theta_t &= \theta_{t-1} + w_t && \text{(State equation)} \\ Y_t &= \phi\theta_t + \mu_t && \text{(Observation equation),}\end{aligned}$$

where  $w_t \sim N(0, \sigma_w^2)$  and  $\mu_t \sim N(0, \sigma_\mu^2)$ ,  $t = 1, 2, \dots, 50$ . In addition, we use the proposal distribution  $q(\theta_t | y_t) = \pi(\theta_t)$ , which is normal with mean  $\theta_{t-1}$  and variance  $\sigma_w^2$ . Then, the weights are given by the recursion

$$w_t^{(s)} \propto w_{t-1}^{*(s)} p(y_t | \theta_t, \sigma_\mu^2),$$

where  $p(y_t | \theta_t, \sigma_\mu^2)$  is  $N(\phi\theta_t, \sigma_\mu^2)$ .

We can compute the mean and standard deviation of the state at each  $t$  using

$$\hat{\theta}_t = \sum_{s=1}^S w_t^{*(s)} \theta_t^{(s)}$$

and

$$\hat{\sigma}_\theta = \left( \sum_{s=1}^S w_t^{*(s)} \theta_t^{2(s)} - \hat{\theta}_t^2 \right)^{1/2}.$$

The following code demonstrates the implementation of this algorithm, setting  $\sigma_w^2 = \sigma_\mu^2 = 1$  and  $\phi = 0.5$ . First, we simulate the process, and then we implement the SIS algorithm.

### R code. Sequential importance sampling: Dynamic linear model

```

1 rm(list = ls()); set.seed(010101)
2 S <- 50000 # Number of particles
3 sigma_w <- 1 # State noise
4 sigma_mu <- 1 # Observation noise
5 phi <- 0.5 # Coefficient in observation equation
6 T <- 50 # Sample size
7 # Simulate true states and observations
8 theta_true <- numeric(T); y_obs <- numeric(T)
9 theta_true[1] <- rnorm(1, mean = 0, sd = sigma_w) # Initial
   state
10 for (t in 2:T) {
11   theta_true[t] <- rnorm(1, mean = theta_true[t-1], sd =
      sigma_w)
12 }
13 y_obs <- rnorm(T, mean = phi*theta_true, sd = sigma_mu)

```

### R code. Sequential importance sampling: Dynamic linear model

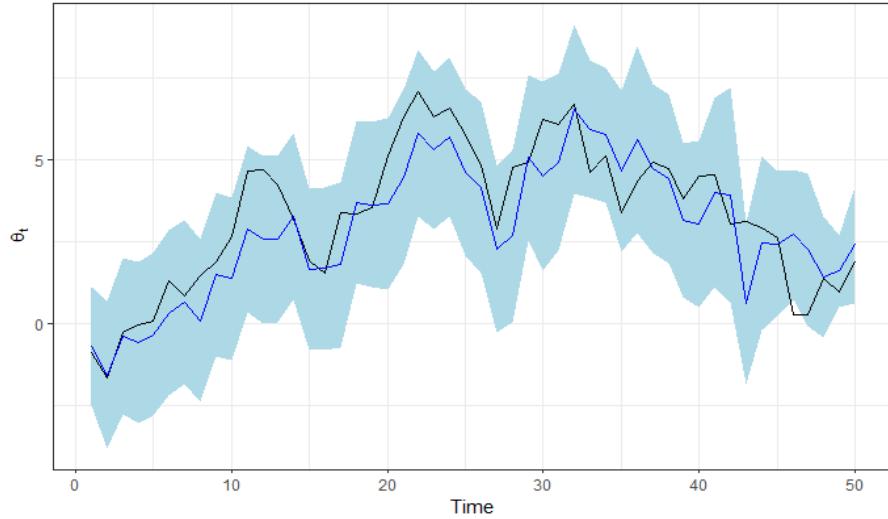
```

1 # Sequential Importance Sampling (SIS)
2 particles <- matrix(0, nrow = S, ncol = T)
3 weights <- matrix(0, nrow = S, ncol = T)
4 weightsSt <- matrix(0, nrow = S, ncol = T)
5 # Initialization
6 particles[, 1] <- rnorm(S, mean = 0, sd = sigma_w) # Sample
    initial particles
7 weights[, 1] <- dnorm(y_obs[1], mean = phi*particles[, 1],
    sd = sigma_mu) # Importance weights
8 weightsSt[, 1] <- weights[, 1] / sum(weights[, 1]) # Standardized weights
9 # Sequential updating
10 for (t in 2:T) {
11     # Propagate particles
12     particles[, t] <- rnorm(S, mean = particles[, t-1], sd =
        sigma_w)
13     # Compute weights
14     weights[, t] <- weightsSt[, t-1] * dnorm(y_obs[t], mean =
        phi*particles[, t], sd = sigma_mu) # Recursive weight
        update
15     weightsSt[, t] <- weights[, t] / sum(weights[, t]) # Normalize weights
16 }
17 # Estimate the states (weighted mean)
18 FilterDist <- colSums(particles * weightsSt)
19 SDFilterDist <- (colSums(particles^2 * weightsSt) -
    FilterDist^2)^0.5
20 library(dplyr); library(ggplot2); library(latex2exp)
21 ggplot2::theme_set(theme_bw())
22 df <- tibble(t = 1:T, mean = FilterDist, lower = FilterDist
    - 2*SDFilterDist, upper = FilterDist + 2*SDFilterDist,
    theta_true = theta_true)
23 # Function to plot
24 plot_filtering_estimates <- function(df) {
25     p <- ggplot(data = df, aes(x = t)) + geom_ribbon(aes(ymin =
        lower, ymax = upper), alpha = 1, fill = "lightblue") +
        geom_line(aes(y = theta_true), colour = "black", alpha
        = 1, linewidth = 0.5) + geom_line(aes(y = mean), colour
        = "blue", linewidth = 0.5) + ylab(TeX("\theta_{\{t\}}"))
        + xlab("Time")
26     print(p)
27 }
28 plot_filtering_estimates(df)
29

```

Figure 4.8 shows the trajectory of the true state vector (black line), the

posterior mean (blue line), and the area defined by  $\pm 2\hat{\sigma}_\theta$  (light blue shaded area).



**FIGURE 4.8**

State in linear state-space model: True and mean estimate using sequential importance sampling,  $T = 50$ .

Sequential importance sampling is effective for sampling from the posterior distribution in the short term. However, it is important to note that SIS is a particular case of IS and, consequently, inherits the drawbacks of importance sampling. In particular, the variance of the weights increases exponentially with  $t$  [246]. This implies that, as  $t$  increases, the importance weights tend to degenerate in the long run; that is, all probability mass concentrates on a few weights, a phenomenon known as sample impoverishment or weight degeneracy. This is because it is impossible to accurately represent a distribution on a space of arbitrarily high dimension with a sample of fixed, finite size. This phenomenon can be observed, for instance, in the dynamic linear model example, where the highest standardized weight at  $t = 50$  is 53%, and 7 out of 50,000 particles account for 87% of the total probability.

Given that, in practice, we are often interested in lower-dimensional marginal distributions, ideas from sampling/importance resampling can be employed. This strategy avoids the accumulation of errors due to resetting the system, although resampling introduces some additional Monte Carlo variation. [181] proposed the *Bootstrap filter*, where, at each time step, resampling is performed by drawing  $S$  particles from the current set using the standardized weights as probabilities of selection. This ensures that particles with small weights have a low probability of being selected. After resampling, the standardized weights are set equal to  $1/S$ . Note that the *Bootstrap filter* involves

multiple iterations of the SIR algorithm, which implies that the resampled trajectories are no longer independent. This multinomial resampling provides an unbiased approximation to the posterior distribution obtained by SIS [125].

Algorithm A5 shows how to perform the *particle filter*. We set  $w_t^{(s)} := w_t(\boldsymbol{\theta}_{0:t}^{(s)})$  to simplify notation. [125].

---

**Algorithm A5** The *particle filter* algorithm

---

```

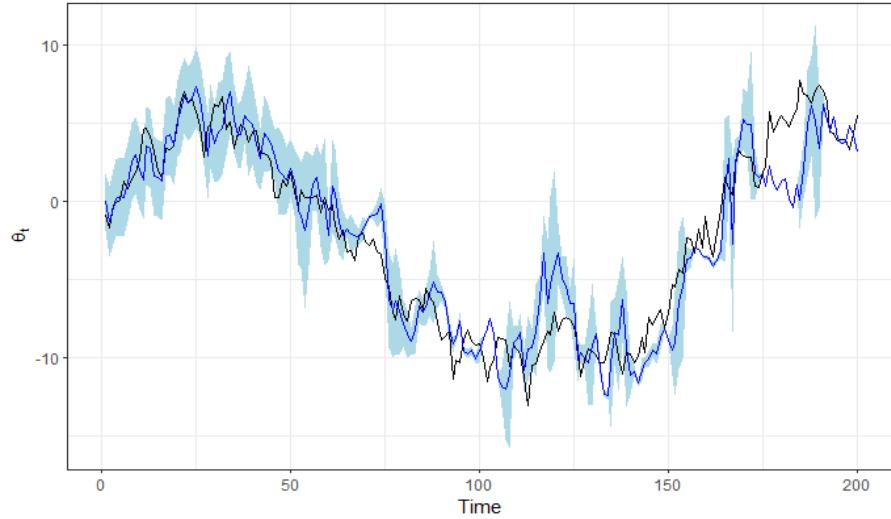
1: for  $s = 1, \dots, S$  do
2:   Sample  $\boldsymbol{\theta}_0^{(s)}$  from  $q(\boldsymbol{\theta}_0 | y_0)$ 
3:   Calculate the importance weights  $w_0^{(s)} \propto \frac{p(y_0 | \boldsymbol{\theta}_0^{(s)})\pi(\boldsymbol{\theta}_0^{(s)})}{q(\boldsymbol{\theta}_0^{(s)} | y_0)}$ 
4: end for
5: Standardize the weights  $w_0^{*(s)} = \frac{w_0^{(s)}}{\sum_{h=1}^S w_0^{(h)}}, s = 1, 2, \dots, S$ 
6: Select  $S$  particles from  $\{\boldsymbol{\theta}_0^{(s)}, w_0^{*(s)}\}$  to obtain  $\{\boldsymbol{\theta}_0^{r(s)}, 1/S\}$ 
7: for  $t = 1, \dots, T$  do
8:   for  $s = 1, \dots, S$  do
9:     Draw particles  $\boldsymbol{\theta}_t^{(s)}$  from  $q_t(\boldsymbol{\theta}_t | \boldsymbol{\theta}_{t-1}, \mathbf{y}_t)$ 
10:    Set  $\boldsymbol{\theta}_{1:t}^{(s)} \leftarrow (\boldsymbol{\theta}_{1:t-1}^{(s)}, \boldsymbol{\theta}_t^{(s)})$ 
11:    Compute the weights  $\alpha_t^{(s)} = \frac{p(\mathbf{y}_t | \boldsymbol{\theta}_t^{(s)})\pi(\boldsymbol{\theta}_t^{(s)} | \boldsymbol{\theta}_{t-1}^{(s)})}{q(\boldsymbol{\theta}_t^{(s)} | \boldsymbol{\theta}_{t-1}^{(s)}, \mathbf{y}_t)}$ 
12:   end for
13:   Standardize the weights  $w_t^{*(s)} = \frac{w_t^{(s)}}{\sum_{h=1}^S w_t^{(h)}}, s = 1, 2, \dots, S$ 
14:   Select  $S$  particles from  $\{\boldsymbol{\theta}_{1:t}^{(s)}, w_t^{*(s)}\}$  to obtain  $\{\boldsymbol{\theta}_{1:t}^{r(s)}, 1/S\}$ 
15: end for
```

---

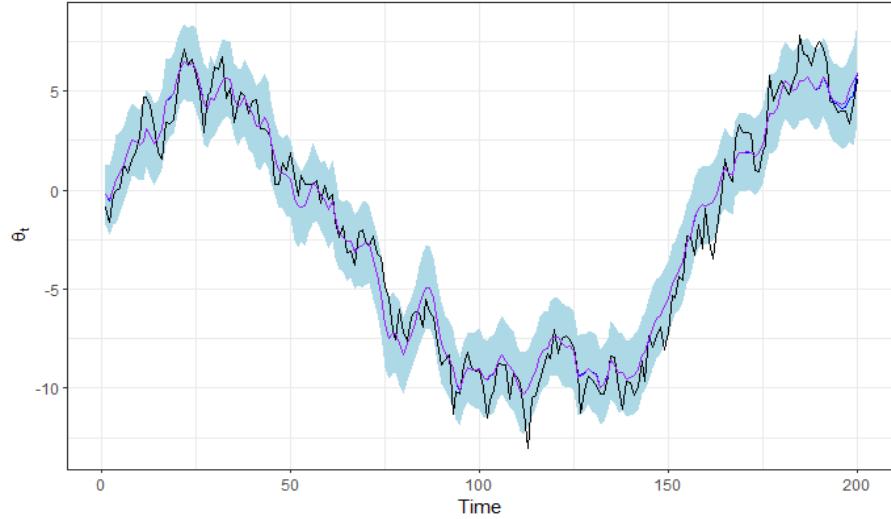
**Example: Dynamic linear model continues**

Let's apply the SIS algorithm to the dynamic linear model with a sample size of 200. Figure 4.9 illustrates the performance of sequential importance sampling. We observe that the algorithm's performance deteriorates as  $t$  increases. This is due to particle degeneration; at  $t = 200$ , a single particle holds a weight close to 100%.

Let's perform particle filtering in this example. The following code illustrate the procedure. Figure 4.10 shows the performance of particle filtering in this example. There is the true state vector (black line), the means based on  $\{\boldsymbol{\theta}_{1:t}^{(s)}, w_t^{*(s)}\}$  (blue line) and  $\{\boldsymbol{\theta}_{1:t}^{r(s)}, 1/S\}$  (purple line), and the area defined by  $\pm 2\hat{\sigma}_\theta$  based on the former (light blue shaded area). Note that the particle filtering algorithm has better performance than the SIS algorithm.

**FIGURE 4.9**

State in linear state-space model: True and mean estimate using sequential importance sampling,  $T = 200$ .

**FIGURE 4.10**

State in linear state-space model: True and mean estimate using particle filtering,  $T = 200$ .

**R code. Particle filtering: Dynamic linear model**

```

1 rm(list = ls()); set.seed(010101)
2 S <- 50000 # Number of particles
3 sigma_w <- 1; sigma_mu <- 1 # # State and observation
   noises
4 phi <- 0.5 # Coefficient in observation equation
5 T <- 200 # Sample size
6 # Simulate true states and observations
7 theta_true <- numeric(T); y_obs <- numeric(T)
8 theta_true[1] <- rnorm(1, mean = 0, sd = sigma_w)
9 for (t in 2:T) {
10   theta_true[t] <- rnorm(1, mean = theta_true[t-1], sd =
      sigma_w)
11 }
12 y_obs <- rnorm(T, mean = phi*theta_true, sd = sigma_mu)
13 # Particle filtering
14 particles <- matrix(0, nrow = S, ncol = T) # Store
   particles
15 particlesT <- matrix(0, nrow = S, ncol = T) # Store
   resampling particles
16 weights <- matrix(0, nrow = S, ncol = T) # Store weights
17 weightsSt <- matrix(0, nrow = S, ncol = T) # Store
   standardized weights
18 weightsSTT <- matrix(1/S, nrow = S, ncol = T) # Store
   standardized weights
19 logalphas <- matrix(0, nrow = S, ncol = T) # Store log
   incremental weights
20 particles[, 1] <- rnorm(S, mean = 0, sd = sigma_w)
21 weights[, 1] <- dnorm(y_obs[1], mean = phi*particles[, 1],
   sd = sigma_mu) # Importance weights
22 weightsSt[, 1] <- weights[, 1] / sum(weights[, 1]) #
   Normalize weights
23 ind <- sample(1:S, size = S, replace = TRUE, prob =
   weightsSt[, 1]) # Resample
24 particles[, 1] <- particles[ind, 1] # Resampled particles
25 particlesT[, 1] <- particles[, 1] # Resampled particles

```

### R code. Particle filtering: Dynamic linear model

```

1 # Sequential updating
2 pb <- txtProgressBar(min = 0, max = T, style = 3)
3 for (t in 2:T) {
4   particles[, t] <- rnorm(S, mean = particles[, t-1], sd =
5     sigma_w)
6   logalphas[, t] <- dnorm(y_obs[t], mean = phi*particles[, t],
7     sd = sigma_mu, log = TRUE)
8   weights[, t] <- exp(logalphas[, t])
9   weightsSt[, t] <- weights[, t] / sum(weights[, t])
10  if(t < T){
11    ind <- sample(1:S, size = S, replace = TRUE, prob =
12      weightsSt[, t])
13    particles[, 1:t] <- particles[ind, 1:t]
14  }else{
15    ind <- sample(1:S, size = S, replace = TRUE, prob =
16      weightsSt[, t])
17    particlesT[, 1:t] <- particles[ind, 1:t]
18  }
19  setTxtProgressBar(pb, t)
20 }
21 close(pb)
22 FilterDist <- colSums(particles * weightsSt)
23 SDFilterDist <- (colSums(particles^2 * weightsSt) -
24   FilterDist^2)^0.5
25 FilterDistT <- colSums(particlesT * weightsSTT)
26 SDFilterDistT <- (colSums(particlesT^2 * weightsSTT) -
27   FilterDistT^2)^0.5
28 MargLik <- colMeans(weights)
29 plot(MargLik, type = "l")
30 library(dplyr)
31 library(ggplot2)
32 require(latex2exp)
33 ggplot2::theme_set(theme_bw())
34 df <- tibble(t = 1:T, mean = FilterDist, lower = FilterDist -
35   2*SDFilterDist, upper = FilterDist + 2*SDFilterDist,
36   meanT = FilterDistT, lowerT = FilterDistT - 2*
37   SDFilterDistT,
38   upperT = FilterDistT + 2*SDFilterDistT, x_true = theta_true)
39 plot_filtering_estimates <- function(df) {
40   p <- ggplot(data = df, aes(x = t)) + geom_ribbon(aes(ymin =
41     lower, ymax = upper), alpha = 1, fill = "lightblue") +
42     geom_line(aes(y = x_true), colour = "black", alpha = 1,
43     linewidth = 0.5) + geom_line(aes(y = mean), colour =
44     "blue", linewidth = 0.5) +
45     geom_line(aes(y = meanT), colour = "purple", linewidth =
46     0.5) + ylab(TeX("\theta_{\{t\}}")) + xlab("Time")
47   print(p)
48 }
49 plot_filtering_estimates(df)

```

Algorithm A5 performs resampling at every time step. However, it is common to perform resampling only when the effective sample size of the particles ( $ESS = (\sum_{s=1}^S (w_t^{*(s)})^2)^{-1}$ ) falls below a specific threshold, such as 50% of the initial number of particles. Note that when  $w_t^{*(s)} = 1/S$ , the effective sample size is  $S$ , the total number of particles. Additionally, we should use  $\{\boldsymbol{\theta}_{1:t}^{(s)}, w_t^{*(s)}\}$  to estimate the posterior distribution, as it results in lower Monte Carlo error compared to calculations based on  $\{\boldsymbol{\theta}_{1:t}^{r(s)}, 1/S\}$  [66]. Finally, an estimate of the marginal likelihood can be obtained using

$$\hat{p}(y_t) = \frac{1}{S} \sum_{s=1}^S w_t^{(s)}.$$

*Particle filtering* offers several advantages, such as being quick and easy to implement, its modularity—allowing one to simply adjust the expressions for the importance distribution and weights when changing the problem—and its suitability for parallel algorithms. Moreover, it enables straightforward sequential inference for very complex models.

However, there are also disadvantages. The resampling step introduces extra Monte Carlo variability. Using the state transition (prior) density as the importance distribution often leads to poor performance, manifested in a lack of robustness with respect to the observed sequence. For instance, performance deteriorates when outliers occur in the data or when the variance of the observation noise is small. Furthermore, the procedure is not well suited for sampling from  $\pi(\boldsymbol{\theta}_{0:t} | y_{1:t})$  because most particles originate from the same ancestor.

Alternative resampling approaches, such as residual resampling [276] and systematic resampling [70], preserve unbiasedness while reducing variance. Additionally, auxiliary particle filtering [67] can help decrease Monte Carlo variability.

Lastly, estimating fixed parameters such as  $\sigma_w^2$ ,  $\sigma_\mu^2$ , and  $\phi$  in the dynamic linear model poses a challenge. Various methods exist to address this issue; see [236, 237] for a comprehensive review and [13] for a seminal work in *particle MCMC* methods.

#### 4.4 Convergence diagnostics

MCMC methods rely on *irreducibility*, *positive recurrence*, and *aperiodicity*, ensuring that, after a sufficient burn-in (warm-up) period, the posterior draws are sampled from the invariant stationary posterior distribution. This can be achieved by running multiple chains initiated at different points and then

mixing them, or by running a single longer chain. In most of the second part of this book, we follow the latter approach, as suggested by [175].

In this section, we present diagnostics to assess whether the sample draws come from the stationary posterior distribution. First, we calculate the numerical standard error associated with the MCMC algorithm. Next, we review the effective number of simulation draws and various convergence tests. Finally, we examine potential errors in the posterior simulator.

#### 4.4.1 Numerical standard error

Many times, the goal in Bayesian inference is to obtain a set of independent draws  $\boldsymbol{\theta}^{(s)}$ ,  $s = 1, 2, \dots, S$ , from the posterior distribution, such that a measure of interest can be estimated with reasonable precision. In particular, we approximate (4.1) using (4.2). By the central limit theorem, we know that

$$\frac{\bar{h}(\boldsymbol{\theta})_S - \mathbb{E}_\pi[h(\boldsymbol{\theta})]}{\sigma_h(\boldsymbol{\theta})/\sqrt{S}} \xrightarrow{d} N(0, 1), \quad (4.4)$$

where  $\sigma_h^2(\boldsymbol{\theta})$  is the variance of  $h(\boldsymbol{\theta})$ .

If we have independent draws, we can estimate  $\sigma_h^2(\boldsymbol{\theta})$  using the posterior draws as follows:

$$\hat{\sigma}_{Sh}^2(\boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S \left[ h(\boldsymbol{\theta}^{(s)}) \right]^2 - [\bar{h}(\boldsymbol{\theta})_S]^2.$$

However, if there are dependent draws, we have

$$\begin{aligned} \hat{\sigma}_{Sh}^{2*}(\boldsymbol{\theta}) &= \frac{1}{S} \left\{ \sum_{s=1}^S \left[ h(\boldsymbol{\theta}^{(s)}) - \bar{h}(\boldsymbol{\theta})_S \right]^2 \right. \\ &\quad \left. + 2 \sum_{l=k+1}^K (h(\boldsymbol{\theta}^{(l)}) - \bar{h}(\boldsymbol{\theta})) (h(\boldsymbol{\theta}^{(l-k)}) - \bar{h}(\boldsymbol{\theta})) \right\}. \end{aligned}$$

The *numerical standard error* is given by  $\sigma_h(\boldsymbol{\theta})/\sqrt{S}$  and serves as a measure of the approximation error in the Monte Carlo integration. Note that this error can be decreased by increasing  $S$ . For instance,  $S = 1000$  implies an error proportional to 3.2%, while  $S = 10000$  reduces the error to approximately 1%.

#### 4.4.2 Effective number of simulation draws

MCMC posterior draws are not independent; therefore, the effective sample size of the posterior chains is not equal to  $S$ . To assess the effective sample size of the posterior draws, we use the following measure:

$$S_{\text{ef}} = \frac{S}{1 + 2 \sum_{k=1}^{\infty} \rho_k(h)},$$

where  $\rho_k(h)$  is the autocorrelation of the sequence  $h(\boldsymbol{\theta})$  at lag  $k$ .

The sample counterpart of this expression is:

$$\hat{S}_{\text{ef}} = \frac{S}{1 + 2 \sum_{k=1}^K \hat{\rho}_k(h)},$$

where

$$\hat{\rho}_k(h) = \frac{\sum_{l=k+1}^K (h(\boldsymbol{\theta}^{(l)}) - \bar{h}(\boldsymbol{\theta}))(h(\boldsymbol{\theta}^{(l-k)}) - \bar{h}(\boldsymbol{\theta}))}{\sum_{s=1}^K (h(\boldsymbol{\theta}^{(s)}) - \bar{h}(\boldsymbol{\theta}))^2}.$$

If  $\hat{\rho}_k(h)$  declines to zero slowly as  $k$  increases, it indicates significant memory in the draws. Consequently, the effective sample size of the posterior draws is small, and it becomes necessary to either decrease the autocorrelation or increase the number of posterior draws.

Note that

$$\hat{\sigma}_{Sh}^{2*}(\boldsymbol{\theta}) = \hat{\sigma}_{Sh}^2(\boldsymbol{\theta}) (1 + 2 \sum_{k=1}^K \hat{\rho}_k(h)),$$

where  $\hat{\sigma}_{Sh}^{2*}(\boldsymbol{\theta})$  and  $\hat{\sigma}_{Sh}^2$  are the simulation variances using dependent and independent draws, and  $\hat{\kappa}(h) = (1 + 2 \sum_{k=1}^K \hat{\rho}_k(h))$  is called the *inefficiency factor*, which represents the inflation of the simulation variance due to autocorrelation in the draws. Values near one indicate draws with little correlation.

#### 4.4.3 Tests of convergence

Regarding convergence issues, there are several diagnostics to assess the adequacy of the posterior chains. In particular, graphical approaches such as trace plots and autocorrelation plots are widely used. Trace plots display the sampled values of a parameter (or multiple parameters) as a function of the iteration number, while autocorrelation plots graphically represent  $\hat{\rho}_k$ . The latter shows how correlated the values of  $\boldsymbol{\theta}$ , or functions of  $\boldsymbol{\theta}$ , are at different lags. Trace plots should fluctuate around a stable mean, exploring the entire parameter space without becoming stuck in any particular region. Autocorrelation plots, on the other hand, should exhibit values close to zero or diminish quickly as the lag increases.

Additionally, Geweke's test [170] provides a simple two-sample test of means. If the mean of the first window (10% of the chain) is not significantly different from the mean of the second window (50% of the chain), we do not reject the null hypothesis that the two segments of the chain are drawn from the same stationary distribution.

The Raftery and Lewis test [352] is designed to calculate the approximate number of iterations ( $S$ ), burn-in ( $b$ ), and thinning parameter ( $d$ ) required to estimate  $p[H(\boldsymbol{\theta}) \leq h]$ , where  $H(\boldsymbol{\theta}) : \mathcal{R}^K \rightarrow \mathcal{R}$ . This calculation is based on a specific quantile of interest ( $q$ ), precision ( $r$ ), and probability ( $p$ ). The diagnostic is based on the dependence factor,  $I = \frac{S+b}{S_{\text{Min}}}$ , where

$S_{\text{Min}} = \Phi^{-1} \left( \frac{1}{2}(p+1) \right)^2 q(1-q)/r^2$ , and  $\Phi(\cdot)$  is the standard normal cumulative distribution function. Values of  $I$  much greater than 5 indicate a high level of dependence.

Heidelberger and Welch's test [202] uses a Cramér-von Mises statistic to test the null hypothesis that the sampled values,  $\boldsymbol{\theta}^{(s)}$ , are drawn from a stationary distribution. The statistic is given by

$$\text{CVM}(B_S) = \int_0^1 B_S(t)^2 dt,$$

where  $B_S(t) = \frac{S_{[St]} - [St]\bar{\theta}^S}{\sqrt{Sp(0)}}$ ,  $S_S = \sum_{s=1}^S \boldsymbol{\theta}^{(s)}$ ,  $\bar{\theta}^S = S_S/S$ , and  $p(0)$  is the spectral density at 0, with  $0 \leq t \leq 1$ . Under the null hypothesis,  $B_S(t)$  converges in distribution to a Brownian bridge.

This test is recursively applied until either the null hypothesis is not rejected, or  $s = 50\%$  of the chain has been discarded. Subsequently, the half-width test calculates a 95% credible interval for the mean using the portion of the chain that passed the stationarity test. If the ratio of the half-width of this interval to the mean is less than 0.1, the test is considered passed. This indicates no evidence to reject the null hypothesis that the estimated mean is accurate and stable.

There are other diagnostics in Bayesian inference that we do not mention here, such as the Gelman and Rubin test [164]. This is because we focus on the available diagnostics in our Graphical User Interface (GUI).

#### 4.4.4 Checking for errors in the posterior simulator

In this book, we provide basic code templates to get posterior draws for performing inference under the Bayesian framework when there is no closed-form solution. We are prone to making mistakes and greatly appreciate your feedback to help improve our code and identify any other potential issues. One way to check if our code works correctly is to perform simulations where the population parameters are known. If the code is functioning properly, the posterior estimates should converge to these values as the sample size increases due to the Bayesian consistency. This is an informal approach to identifying potential mistakes.

[172] offers a more formal method for code validation. The starting point is the joint density  $p(\mathbf{y}, \boldsymbol{\theta}) = p(\mathbf{y} \mid \boldsymbol{\theta})\pi(\boldsymbol{\theta})$  and a test function  $h(\mathbf{y}, \boldsymbol{\theta})$  such that  $\sigma_h^2 = \text{Var}[h(\mathbf{y}, \boldsymbol{\theta})] < \infty$ .

Assume that there is a *marginal-conditional simulator* for the joint distribution of  $\mathbf{y}$  and  $\boldsymbol{\theta}$ :

$$\begin{aligned}\boldsymbol{\theta}^{(s)} &\sim \pi(\boldsymbol{\theta}) \\ \mathbf{y}^{(s)} &\sim p(\mathbf{y} \mid \boldsymbol{\theta}^{(s)}) \\ h^{(s)} &= h(\mathbf{y}^{(s)}, \boldsymbol{\theta}^{(s)}).\end{aligned}$$

The sequence  $\{\mathbf{y}^{(s)}, \boldsymbol{\theta}^{(s)}\}$  is i.i.d.,  $\bar{h}_S$  converges almost surely to  $\mathbb{E}[h(\mathbf{y}, \boldsymbol{\theta})]$ , and there is convergence in distribution when  $\bar{h}_S$  is well standardized (see Equation 4.4) and  $\hat{\sigma}_{Sh}^2(\boldsymbol{\theta})$  converges to  $\sigma_h^2(\boldsymbol{\theta})$  almost surely.

A posterior simulator produces draws  $\boldsymbol{\theta}^{(s)}$  given a particular realization  $\mathbf{y}_{\text{Obs}}$ , using the transition density  $q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(s-1)}, \mathbf{y}_{\text{Obs}})$ . Thus, a *successive-conditional simulator* consists of an initial draw  $\boldsymbol{\theta}^{(0)}$  from  $\pi(\boldsymbol{\theta})$  followed by:

$$\begin{aligned}\mathbf{y}^{(l)} &\sim p(\mathbf{y} \mid \boldsymbol{\theta}^{(l-1)}) \\ \boldsymbol{\theta}^{(l)} &\sim q(\boldsymbol{\theta} \mid \mathbf{y}^{(l)}, \boldsymbol{\theta}^{(l-1)}) \\ h^{(l)} &= h(\mathbf{y}^{(l)}, \boldsymbol{\theta}^{(l)}),\end{aligned}$$

where  $\bar{h}_L = L^{-1} \sum_{l=1}^L h(\mathbf{y}^{(l)}, \boldsymbol{\theta}^{(l)})$  converges almost surely to  $\mathbb{E}[h(\mathbf{y}, \boldsymbol{\theta})]$ , and there is convergence in distribution when  $\bar{h}_L$  is well standardized, and  $\hat{\sigma}_{Lh}^{*2}(\boldsymbol{\theta})$  converges to  $\sigma_h^2(\boldsymbol{\theta})$  almost surely, for  $l = 1, 2, \dots, L$ . Thus,

$$\frac{\bar{h}_S - \bar{h}_L}{(S^{-1}\hat{\sigma}_{Sh}^2(\boldsymbol{\theta}) + L^{-1}\hat{\sigma}_{Lh}^{*2}(\boldsymbol{\theta}))^{1/2}} \xrightarrow{d} N(0, 1).$$

Thus, we can test  $H_0: \bar{h}_S - \bar{h}_L = 0$  versus  $H_1: \bar{h}_S - \bar{h}_L \neq 0$ . Rejection of the null indicates potential errors in implementing the posterior simulator.

### Example: Mining disaster change point continues

Let's revisit the mining disaster change point example from subsection 4.1.1 and examine some convergence diagnostics for the posterior draws of the rate of disasters after the change point ( $\lambda_2$ ). The following code demonstrates how to perform these diagnostics using the **R** package *coda*. For clarity and replicability of the results, we present the Gibbs sampler again.

Figures 4.11 and 4.12 show the trace and autocorrelation plots. We observe that the posterior draws of  $\lambda_2$  appear stationary around their mean, and the autocorrelation decreases rapidly to zero.

The mean and standard deviation of the rate after the change point are 0.92 and 0.12, respectively. The naive and time series standard errors are 0.0008245 and 0.0008945, respectively. The naive standard error assumes iid posterior draws, whereas the time series standard error accounts for autocorrelation. Both standard errors are very similar, indicating a low level of autocorrelation, which is consistent with the results shown in Figure 4.12. The effective sample size of the posterior draws is 16,991, while the total number of posterior draws is 20,000 after a burn-in period of 1,000.

The Geweke test statistic is 1.43, which implies no statistical evidence to reject the null hypothesis of equal means in the two segments of the posterior draws. The Raftery and Lewis test yields a dependence factor near 1, indicating a low level of dependence. The Heidelberger and Welch test does not reject the null hypothesis of stationarity for the posterior draws and also confirms that the mean is accurate and stable.

In summary, all posterior diagnostics indicate that the posterior draws originate from an invariant stationary distribution.

The second part of the code implements the proposal by [172] to assess the reliability of the posterior simulator. The parameter vector is defined as  $\boldsymbol{\theta} = [\lambda_1 \ \lambda_2 \ H]$ , and the first moments of these parameters are used as test functions. We do not reject the null hypothesis of equal means across the three test functions, indicating that the posterior simulator is functioning correctly.

To evaluate the effectiveness of the test, we run the *marginal-conditional simulator* with prior parameters  $\alpha_{l0} = 0.5$  and  $\beta_{l0} = 1$ ,  $l = 1, 2$ . In contrast, for the *successive-conditional simulator*, we use prior parameters  $\alpha_{l0} = 1$  and  $\beta_{l0} = 0.5$ ,  $l = 1, 2$ . In this case, we reject the null hypothesis in two out of three test functions, suggesting that the test performs well in this example to detect issues.

### R code. Posterior diagnostics: The mining disaster changepoint

```

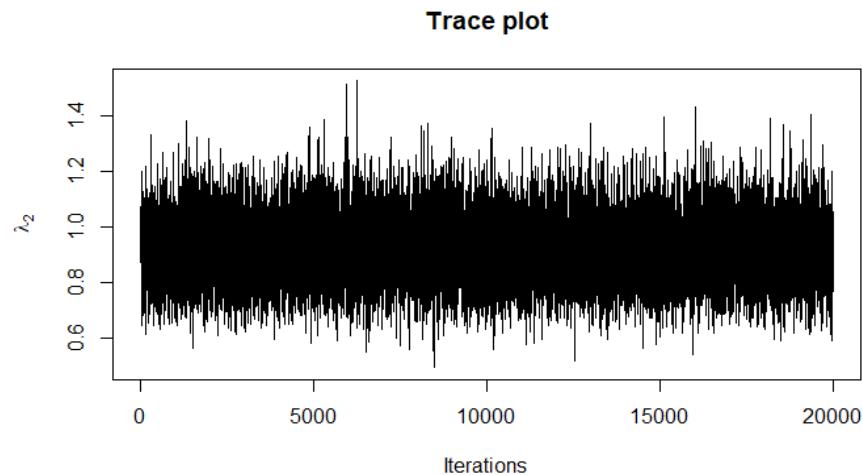
1 rm(list = ls()); set.seed(010101)
2 dataset<-read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/
  MiningDataCarlin.csv",header=T)
3 attach(dataset); str(dataset)
4 a10 <- 0.5; a20 <- 0.5; b10 <- 1; b20 <- 1
5 y <- Count; sumy <- sum(Count); T <- length(Count)
6 theta1 <- NULL; theta2 <- NULL
7 kk <- NULL; H <- 60
8 MCMC <- 20000; burnin <- 1000; S <- MCMC + burnin
9 keep <- (burnin+1):S
10 pb <- txtProgressBar(min = 0, max = S, style = 3)
11 for(s in 1:S){
12   a1 <- a10 + sum(y[1:H]); b1 <- b10+H
13   theta11 <- rgamma(1,a1,b1); theta1 <- c(theta1,theta11)
14   a2 <- a20 + sum(y[(1+H):T]); b2 <- b20 + T-H
15   theta22 <- rgamma(1,a2,b2); theta2 <- c(theta2,theta22)
16   pp<-NULL
17   for(l in 1:T){
18     p <- exp(1*(theta22-theta11))*(theta11/theta22)^(sum(y
      [1:l]))
19     pp <- c(pp,p)
20   }
21   prob <- pp/sum(pp)
22   H <- sample(1:T,1,prob=prob)
23   kk <- c(kk,H)
24   setTxtProgressBar(pb, s)
25 }
```

**R code. Posterior diagnostics: The mining disaster changepoint**

```

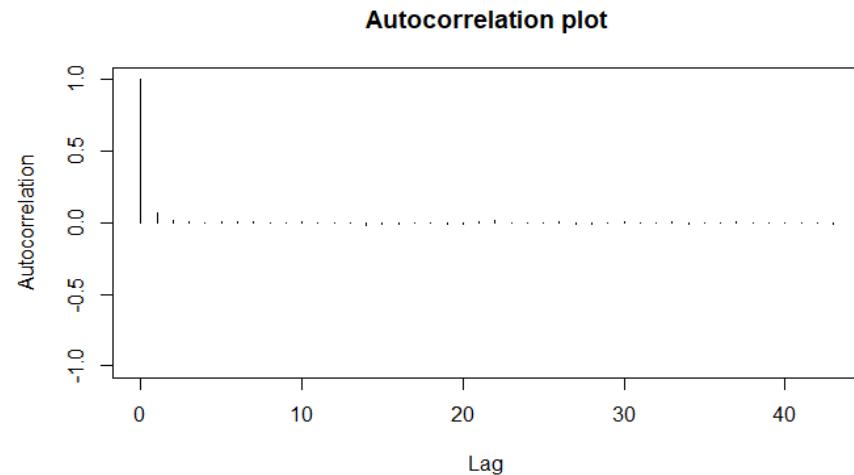
1 close(pb)
2 library(coda); library(latex2exp)
3 theta1Post <- mcmc(theta1[keep]); summary(theta1Post)
4 HPost <- mcmc(kk); summary(HPost)
5 theta2Post <- mcmc(theta2[keep]); summary(theta2Post)
6 plot(theta2Post, density = FALSE, main = "Trace plot", ylab
= TeX("$\\lambda_2$"))
7 autocorr.plot(theta2Post, main = "Autocorrelation plot")
8 raftery.diag(theta2Post, q = 0.025, r = 0.005, s = 0.95)
9 geweke.diag(theta2Post, frac1 = 0.1, frac2 = 0.5)
10 heidel.diag(theta2Post, eps = 0.1, pvalue = 0.05)
11 effectiveSize(theta2Post)

```



**FIGURE 4.11**

Mining disaster change point: Trace plot of  $\lambda_2$ .

**FIGURE 4.12**

Mining disaster change point: Autocorrelation plot of  $\lambda_2$ .

*R code. Posterior diagnostics: The mining disaster changepoint*

```

1 # Marginal-conditional simulator
2 Theta1Prior <- rgamma(MCMC,a10,b10); Theta2Prior <- rgamma(
  MCMC,a20,b20)
3 kPrior <- sample(1:T, MCMC, replace = TRUE, prob = rep(1/T,T
  ))
4 ytPrior <- function(par){
5   y1t <- rpois(par[3], par[1])
6   if(par[3] == T){y2t <- NULL
7   }else{y2t <- rpois(T-par[3], par[2])
8   }
9   yt <- c(y1t, y2t)
10  return(yt)
11 }
12 pars1 <- cbind(Theta1Prior, Theta2Prior, kPrior); Yt <-
  apply(pars1, 1, ytPrior)
13 parsMcmc1 <- coda::mcmc(pars1); Summ1 <- summary(parsMcmc1)

```

**R code. Posterior diagnostics: The mining disaster changepoint**

```

1 # Successive-conditional simulator
2 SucConSim <- function(a10, b10, a20, b20, par){
3   y <- ytPrior(par)
4   theta1 <- par[1]; theta2 <- par[2]; H <- par[3]
5   a1 <- a10 + sum(y[1:H]); b1 <- b10+H
6   theta11 <- rgamma(1,a1,b1)
7   if(H == T){ a2 <- a20
8 } else{ a2 <- a20 + sum(y[(1+H):T])
9 }
10 b2 <- b20 + T-H; theta22 <- rgamma(1,a2,b2); pp<-NULL
11 for(l in 1:T){
12   p <- l*(theta22-theta11) + (sum(y[1:l]))*log(theta11/
13     theta22)
14   pp <- c(pp,p)
15 }
16 pps <- exp(pp - max(pp)); prob <- pps/sum(pps)
17 H <- sample(1:T, 1, prob=prob)
18 parNew <- list(y = y, pars = c(theta11, theta22, H))
19 return(parNew)
20 }
21 a10 <- 0.5; b10 <- 1; a20 <- 0.5; b20 <- 1
22 # a10 <- 1; b10 <- 0.5; a20 <- 1; b20 <- 0.5
23 par1 <- rgamma(1,a10,b10); par2 <- rgamma(1,a20,b20)
24 par3 <- sample(1:T, 1, replace = TRUE, prob = rep(1/T,T))
25 pars2 <- matrix(NA, MCMC, 3); pars2[,1] <- c(par1, par2,
26   par3)
27 for(s in 2:MCMC){
28   Res <- SucConSim(a10 = a10, b10 = b10, a20 = a20, b20 =
29     b20, par = pars2[s-1,])
30   pars2[s, ] <- Res$pars
31 }
32 parsMcmc2 <- coda::mcmc(pars2); Summ2 <- summary(parsMcmc2)
33 TestGeweke <- function(j){
34   Test <- (Summ1[["statistics"]][j,1] - Summ2[["statistics"
35     ][j,1]])/(Summ1[["statistics"]][j,4]+Summ2[["statistics"
36     ][j,4]])^0.5
37   Reject <- abs(Test) > qnorm(0.975)
38   return(list(Test = Test, Reject = Reject))
39 }
40 TestGeweke(1); TestGeweke(2); TestGeweke(3)
41

```

---

## 4.5 Summary

In this chapter, we present the most popular methods for obtaining posterior draws when the posterior distribution does not have a standard closed-form solution. In particular, Markov chain Monte Carlo (MCMC) methods, such as Gibbs sampling and the Metropolis-Hastings algorithm, are the most commonly used approaches in this book. However, Hamiltonian Monte Carlo is gaining particular relevance in high-dimensional problems, while particle filtering (Sequential Monte Carlo) is widely applied in time series models.

Each problem requires careful consideration to determine the most appropriate method, and in many cases, a combination of methods is necessary. For instance, estimating fixed parameters in *state-space* models typically requires MCMC methods, while recursion of the state vector requires particle filtering. Additionally, convergence diagnostics are crucial because MCMC methods rely on technical assumptions that must be verified.

---

## 4.6 Exercises

### 1. Example: The normal model with independent priors

Let's recap the math test exercise in Chapter 3, this time assuming independent priors. Specifically, let  $Y_i \sim N(\mu, \sigma^2)$ , where  $\mu \sim N(\mu_0, \sigma_0^2)$  and  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ . The sample size is 50, and the mean and standard deviation of the math scores are 102 and 10, respectively. We set  $\mu_0 = 100$ ,  $\sigma_0^2 = 100$ , and  $\alpha_0 = \delta_0 = 0.001$ .

- Find the posterior distribution of  $\mu$  and  $\sigma^2$ .
  - Program a Gibbs sampler algorithm and plot the histogram of the posterior draws of  $\mu$
2. Show that the Gibbs sampler is a particular case of the Metropolis-Hastings where the acceptance probability is equal to 1.
  3. Implement a Metropolis-Hastings to sample from the Cauchy distribution,  $C(0, 1)$ , using as proposals a standard normal distribution and a Student's t distribution with 5 degrees of freedom.
  4. This exercise was proposed by Professor Hedibert Freitas Lopes, who cites [431] as a useful reference for an introduction to Hamiltonian Monte Carlo in **R** and the *hmclearn* package. The task is to obtain posterior draws using the Metropolis-Hastings and Hamiltonian Monte Carlo algorithms for the posterior distribution given

by

$$\pi(\theta_1, \theta_2 | \mathbf{y}) \propto \exp \left\{ -\frac{1}{2} (\theta_1^2 \theta_2^2 + \theta_1^2 + \theta_2^2 - 8\theta_1 - 8\theta_2) \right\}.$$

### 5. Ph.D. students sleeping hours continues

- (a) Use importance sampling based on a  $U(0, 1)$  proposal to obtain draws of  $\boldsymbol{\theta} | \mathbf{y} \sim B(16.55, 39.57)$  in the Ph.D. students' sleeping hours example in Chapter 3. Note that, based on Exercise 15 in Chapter 3,  $\alpha_0 = 1.44$  and  $\beta_0 = 2.57$ .
- (b) Compute the marginal likelihood in this context (Bernoulli-Beta model) and compare it to the result obtained using the Gelfand-Dey method.

### 6. Example 4.1 in [181] is

$$\begin{aligned}\theta_t &= 0.5\theta_{t-1} + 25 \frac{\theta_{t-1}}{1 + \theta_{t-1}^2} + 8\cos(1.2t) + w_t \\ y_t &= \frac{\theta_t^2}{20} + \mu_t,\end{aligned}$$

where  $\theta_0 \sim N(0, \sqrt{10})$ ,  $w_t \sim \mathcal{N}(0, \sqrt{10})$  and  $\mu_t \sim N(0, \sqrt{1})$ .

- Perform sequential importance sampling in this example.
- Perform particle (Bootstrap) filtering in this example.
- Estimate the marginal likelihood in this example.

### 7. Ph.D. students sleeping hours continues

- Perform the diagnostics of Section 4.4 in this example.
- Check if there are errors in the posterior simulator of the Metropolis-Hastings algorithm in this example using the Geweke approach using as test functions the first moments of  $p$  and  $p^2$ . Remember from Exercise 15 in Chapter 3 that the sample size is 52, and  $\alpha_0 = 1.22$  and  $\beta_0 = 2.57$ .
- Run the Geweke test using  $\alpha_0 = 2.57$  and  $\beta_0 = 1.22$ , and check the results.

— | — | —

## Part II

# Regression models: A GUIDed toolkit



# 5

---

## *Graphical user interface*

---

This chapter introduces our graphical user interface (GUI) for conducting Bayesian regression analysis in a user-friendly environment that requires no programming skills, using a simple drag-and-drop design. The GUI is implemented as an interactive web application built with *shiny* [77] and integrates packages such as *MCMCpack* [295] and *bayesm* [377] within the **R** software environment [346]. It is primarily designed for teaching and applied purposes at the introductory level. In the subsequent chapters of the second part of this book, we present several applications that illustrate the potential of our GUI for applied researchers and practitioners.

---

### 5.1 Introduction

Our GUI enables users to perform inference using Bayesian regression analysis without requiring programming skills. The lack of such skills is often a significant impediment to the widespread adoption of the Bayesian framework [454, 238].

Several other graphical user interfaces are available for Bayesian regression analysis. *ShinyStan* [422] is a highly flexible and open-source program; however, it requires users to have some programming skills, as it is based on the *Stan* software for Bayesian data analysis [69]. *BugsXLA* [454] is also open source but less flexible, though it does not require any programming experience. *Bayesian Regression: Nonparametric and Parametric Models* [238] is a user-friendly and flexible GUI based on the *MATLAB Compiler* for 64-bit Windows systems. It primarily focuses on Bayesian nonparametric regression and is designed for users already familiar with basic parametric models, such as those implemented in our GUI. Additionally, there are tools such as the *MATLAB Toolkit*, *Stata*, and *BayES*, but these are not open source.

We developed our GUI as an interactive web application using *shiny* [77] and various libraries in **R** [345]. The specific libraries and commands used in our GUI are listed in Table A.1. It includes ten univariate models, four multivariate models, four time-series models, three hierarchical longitudinal models, and seven Bayesian model averaging frameworks. In addition, it provides basic

summaries and diagnostics of posterior chains, along with visualizations such as trace plots, autocorrelation plots, and density plots.

In terms of flexibility and functionality, our GUI lies between *ShinyStan* and *BugsXLA*: it does not require programming skills but is not as advanced as the software described in [238]. However, unlike the latter, our GUI runs on any operating system. We call our GUI **BEsma**rter,<sup>1</sup> and it is freely available at <https://github.com/besmarter/BSTApp>, where users can access all source code and datasets.

Simulated and applied datasets are stored in the **DataSim** and **DataApp** folders of our **GitHub** repository (see Tables A.2 and A.3 for details). The **DataSim** folder contains files used to simulate various processes, providing access to population parameters. These files serve as valuable pedagogical tools for illustrating the statistical properties of the inferential frameworks available in our GUI. The **DataApp** folder includes the datasets used throughout this book, which users can adopt as templates when structuring their own data.

There are four ways to install our GUI. The simplest method, requiring only the installation of **R** and, optionally, an **R** code editor, is to install *shiny*, and then type:

### **R** code. How to display our graphical user interface

```
1 shiny::runGitHub("besmarter/BSTApp", launch.browser = T)
```

in the **R** package console or any **R** code editor, we strongly recommend typing this command manually rather than copying and pasting it, as quotation marks can sometimes cause issues.

The second option is to visit

<https://andres-ramirez-hassan.shinyapps.io/BSTApp/>

Please note: the free Posit Cloud tier sometimes runs out of memory, which can cause the app to stop. Sorry for the inconvenience.

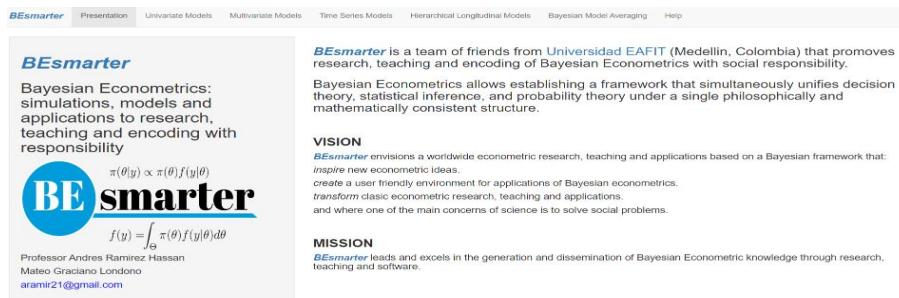
The third option is to visit <https://fly-besmarter.fly.dev/>. As with Posit Cloud, occasional memory limits on the free tier may affect performance.

The fourth approach, and our recommendation, is using a **Docker** image by running:

1. docker pull aramir21/besmartergui:latest
2. docker run -rm -p 3838:3838 aramir21/besmartergui

---

<sup>1</sup>Bayesian Econometrics: Simulations, Models, and Applications to Research, Teaching, and Encoding with Responsibility.



**FIGURE 5.1**  
Display of graphical user interface.

in your **Command Prompt**, this command creates an isolated environment for our GUI, ensuring consistent performance across different systems. Note that **Docker** must be installed to deploy our GUI using this method. Users can then access the app by navigating to `127.0.0.1:3838` or `http://localhost:3838/`.

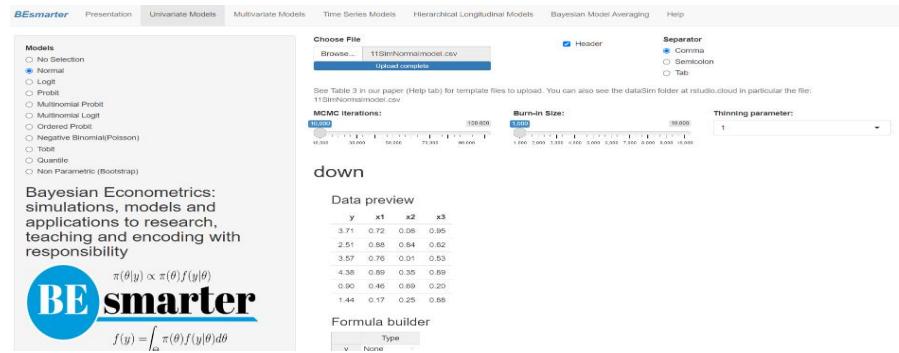
After using any of the four methods to run our GUI, users will see a new window displaying a presentation of our research team (see Figure 5.1). Additionally, the top panel in Figure 5.1 shows the categories of models that can be estimated in our GUI.

## 5.2 Univariate models

After deploying our GUI (see Figure 5.1), the user should select *Univariate Models* from the top panel. Then, Figure 5.2 is displayed, showing a radio button on the left-hand side that lists the specific models within this category. In particular, users can see that the normal model is selected from the univariate models class.

The right-hand panel then displays a widget for uploading the input dataset, which must be a *csv* file with headers in the first row. Users must also select the type of separator used in the input file —comma, semicolon, or tab— (see the **DataSim** and **DataApp** folders for input file templates). Once the dataset is uploaded, users can preview the data. Range sliders allow users to set the number of iterations for the Markov Chain Monte Carlo algorithm, specify the burn-in period, and adjust the thinning parameter (see the following chapters in this part of the book for technical details).

Next, users must specify the model equation. This can be done using the formula builder, where they select the dependent and independent variables

**FIGURE 5.2**

Univariate models: Specification.

and then click on the *Build Formula* tab. The equation is displayed in the *Main Equation* window, formatted according to **R** syntax, for example,  $y \sim x1 + x2 + x3$ . Users can modify this expression as needed to include higher-order terms, interaction effects, or other transformations. These modifications must follow standard formula syntax.<sup>2</sup>

By default, univariate models include an intercept, except for ordered probit models, where the specification must explicitly exclude it due to *identification* constraints (see details below).<sup>3</sup> Thus, users should specify this explicitly as follows:  $y \sim x1 + x2 + x3 - 1$ .

Finally, users must define the prior hyperparameters. For example, in the normal-inverse-gamma model, these include the mean vector, covariance matrix, shape parameter, and scale parameter (see Figure 5.3). However, our GUI uses *noninformative* hyperparameters by default across all modeling frameworks, so this step is optional.

After completing the specification process, users should click the *Go!* button to initiate estimation (see Figure 5.3). Once the process is complete, the GUI displays summary statistics and convergence diagnostics. In addition, widgets allow users to download the posterior chains (as *csv* files) and graphs (as *pdf* or *eps* files). Note that in the results —summary tables, posterior chains, and graphs— the coefficients are ordered such that location parameters appear first, followed by scale parameters.

For multinomial models (probit and logit), the dataset must be structured as follows: the first column should contain the dependent variable, followed by alternative-specific regressors (e.g., alternatives' price), and then non-alternative-specific regressors (e.g., income). The formula builder allows users to specify the dependent variable as well as both types of independent

<sup>2</sup>See <https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/formula>

<sup>3</sup>An *identification* issue arises when multiple sets of model parameters yield the same likelihood function value.

Main Equation  
y ~ x1+x2+x3

Introduce prior mean vector: location parameters

cte0	0
x1	0
x2	0
x3	0

Introduce prior covariances location parameters by row. It has to be symmetric:

cte0	x1	x2	x3
x1	1	0	0
x2	0	1	0
x3	0	0	1

Prior Shape Parameter: Variance Parameter  
0.001

Prior Scale Parameter: Variance Parameter  
0.001

Introduce Prior shape Parameter. Example: 0.001

Introduce prior scale parameter. Example: 0.001

Click the button (Go!) after importing the dataset and selecting the model to update the value displayed in the main panel.

Go!

Warning: Be patient this may take several minutes!!

Download Posterior Chains | Download Posterior Graphs

Working on it. Running MCMC sampling

**FIGURE 5.3**  
Univariate models: Priors.

variables (see technical details in the next chapter). Users must also define the base category, the number of alternatives (which is also required for ordered probit models), the number of alternative-specific regressors, and the number of non-alternative-specific regressors (see Figure 5.4).

For multinomial logit models, users can additionally specify a tuning parameter corresponding to the degrees of freedom of the Metropolis–Hastings algorithm (see technical details in the next chapter). This tuning option is available in our GUI given that estimation relies on the Metropolis–Hastings algorithm.

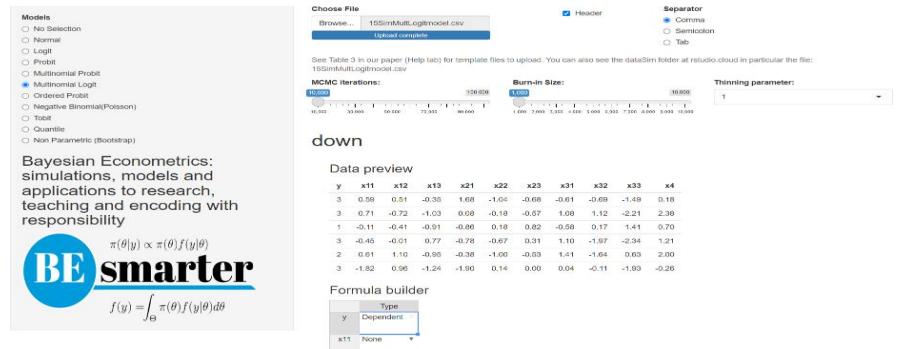
In the results of these models, the coefficients are ordered as follows:

1. Intercepts ( $cte_l$  in the summary display, where  $l$  represents the alternative).
2. Non-alternative-specific regressors ( $NAS_{jl}$  in the summary display, where  $l$  represents the alternative and  $j$  the non-alternative regressor).
3. Alternative-specific regressors ( $AS_j$  in the summary display, where  $j$  represents the alternative-specific regressor).

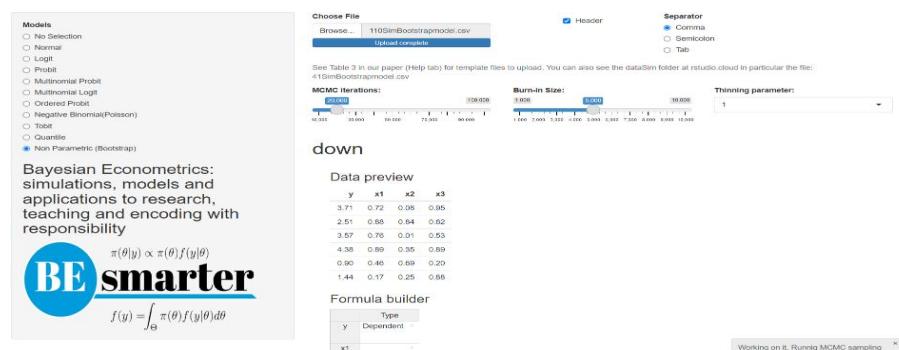
Note that the non-alternative-specific regressors associated with the base category are set to zero and do not appear in the results. Additionally, due to identification constraints in multinomial and multivariate probit models, some coefficients in the main diagonal of the covariance matrix remain constant.

For the negative binomial model, users must specify a dispersion parameter (see the next chapter for details). Similarly, for Tobit and quantile models, users need to define the censorship points and quantiles, respectively.

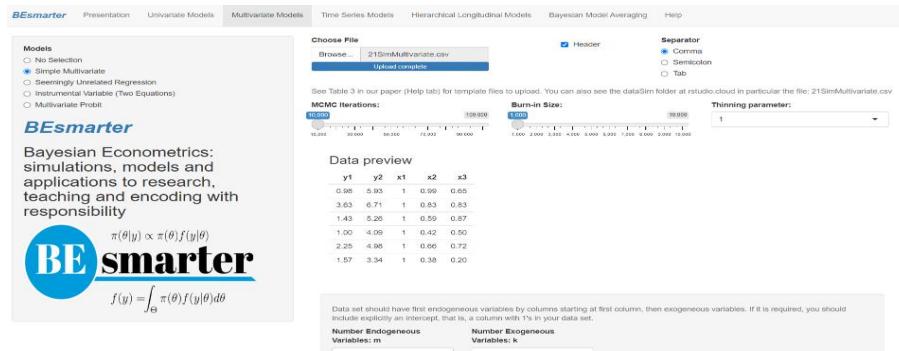
The Bayesian bootstrap method only requires uploading a dataset, specifying the number of MCMC iterations, burn-in size, and defining the equation (see Figure 5.5). The input file should follow the same structure as the one used for the univariate normal model.



**FIGURE 5.4**  
Univariate models: Multinomial.



**FIGURE 5.5**  
Univariate models: Bootstrap.



**FIGURE 5.6**  
Multivariate models: Simple multivariate.

### 5.3 Multivariate models

After deploying our GUI (see Figure 5.1), the user should select *Multivariate Models* from the top panel. Figure 5.6 will then appear, displaying a radio button on the left-hand side that lists the specific models within this category.

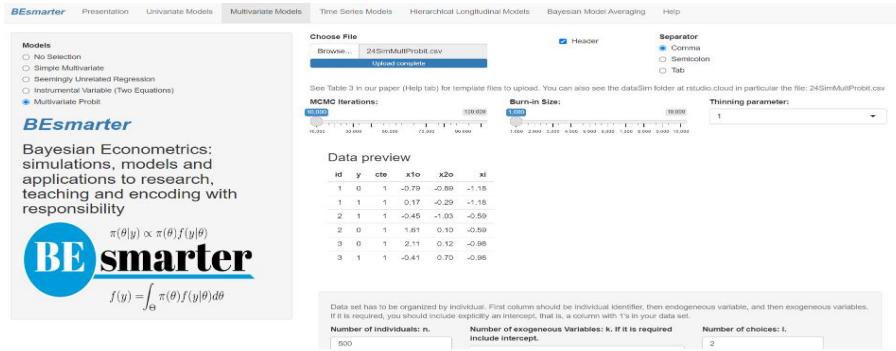
Figure 5.6 illustrates the setup for multivariate regression models. The input file should first include the dependent variables, followed by the regressors. If each equation includes an intercept, a column of ones must be added after the dependent variables in the input file. Once the file is uploaded, users can preview the data.

The user must then specify the number of dependent variables and regressors, indicate whether an intercept should be included, and define the hyperparameter values (see Figure 5.6).

In seemingly unrelated regressions, the input file should first include the dependent variables, followed by the regressors for each equation, including the intercept (a column of ones) if applicable. Users must specify the number of dependent variables (equations), the total number of regressors (i.e., the sum of all regressors across equations), and the number of regressors per equation (including the intercept if relevant). Users may also define the values of the hyperparameters if prior information is available.

For both simple multivariate and seemingly unrelated regression models, the results first display the posterior estimates of the location parameters by equation, followed by the posterior covariance matrix.

In the instrumental variable setting, users must specify both the main equation and the instrumental equation. Intercepts are included by default. The first variable on the right-hand side of the main equation must correspond to the endogenous regressor. In the instrumental equation, this endogenous variable serves as the dependent variable and is modeled as a function of

**FIGURE 5.7**

Multivariate models: Multivariate probit.

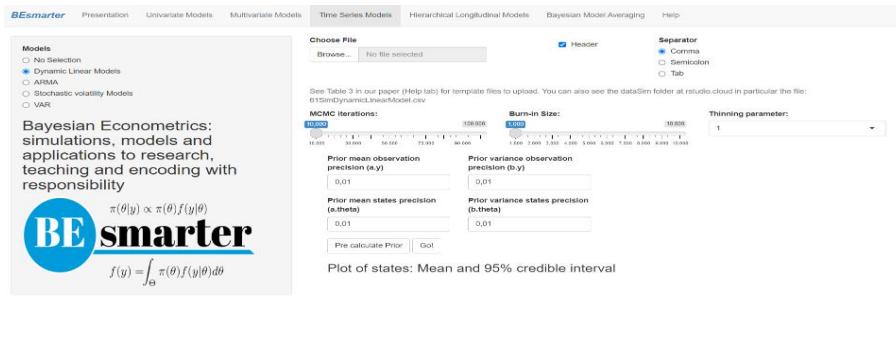
the instruments. Users may also specify hyperparameter values if prior information is available. The input file should contain the dependent variable, the endogenous regressor, the instruments, and the exogenous regressors. The results first report the posterior estimates of the endogenous regressor, followed by the location parameters of the auxiliary (instrumental) regression, the location parameters of the exogenous regressors, and finally, the posterior covariance matrix.

The multivariate probit model requires the input dataset to be ordered by unit. For example, if there are three choices, each unit should appear three times in the dataset. The first column must contain a unique identifier for each unit, using ordered integers. The next column should contain the dependent variable, represented as a single vector of 0s and 1s, followed by the regressors, which must include a column of ones for the intercepts. Users must specify the number of units, the number of regressors, and the number of choices (see Figure 5.7). The results first display the posterior estimates of the location parameters by equation, followed by the posterior covariance matrix.

## 5.4 Time series models

After our GUI is deployed (see Figure 5.1), the user should select *Time Series Models* from the top panel. Then, Figure 5.8 will be displayed, and the user will see the radio button on the left-hand side, which shows the specific models within this general class.

Users can perform inference using dynamic linear models (DLM), autoregressive moving average (ARMA) models, stochastic volatility models (SVM), and vector autoregressive (VAR) models. Users should upload a dataset, which



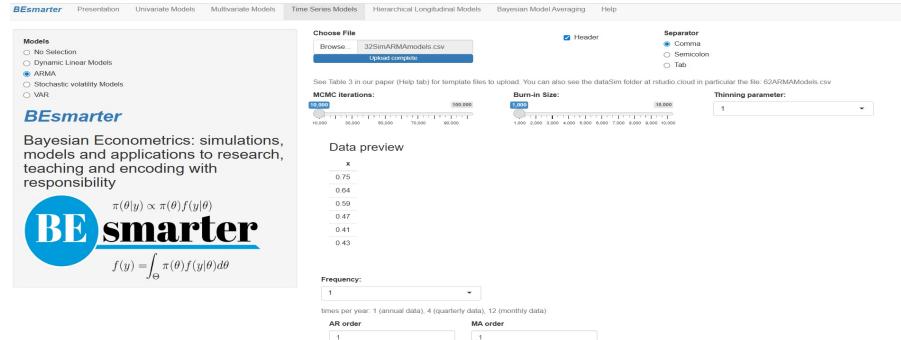
**FIGURE 5.8**  
Time series models.

must be a *csv* file with headers in the first row. The files for DLMs and SVMs have the same structure: the first column contains the dependent variable, followed by the independent variables. For ARMA models, there is only one column with the modeled variable, while VAR models have each modeled variable in a separate column. Note that this version of the GUI does not allow for exogenous variables in VAR models. Users should specify the separator used in the input file: comma, semicolon, or tab. A dataset preview is displayed once the file is uploaded. Dataset templates can be found in the folders **DataSim** (see Table A.2 for details) and **DataApp** (see Table A.3 for details) in our *GitHub* repository.

Next, users should set the MCMC and burn-in iterations using the range sliders and the thinning parameter using the input box.

To estimate DLMs, users should set the hyperparameters for the precision of the observation equation and the state equations (means and variances) if prior information is available. Otherwise, users can click the *Pre Calculate Prior* button, where these hyperparameters are estimated based on a recursive model estimation using ordinary least squares (OLS). The sample size is progressively increased, and the location parameters are saved. The GUI then computes the covariance matrix of this sequence and uses it to set the prior mean for the precision of the state vector, which is equal to the inverse of the maximum element on the main diagonal of the covariance matrix (*a.theta*). The prior variance is set to ten times this value (*b.theta*). For the observation equation, the prior mean of the precision is set to the inverse of the OLS variance estimate (*a.y*), and the prior variance is set to ten times this value (*b.y*). This is a rudimentary approach to setting these hyperparameters, and users are encouraged to use a more thoughtful process.

Next, users should click the *Go!* button to start estimating the model. This may take a few minutes, as DLMs are complex to estimate. Users should be patient. Once the estimation is complete, the GUI will display graphs of the states (mean and 95% credible intervals), summary statistics of the posterior

**FIGURE 5.9**

Time series models: ARMA specification

chains for the observation and state variances, and convergence diagnostics. Users can download the mean and the lower and upper limits of the 95% credible intervals of the states, as well as the posterior chains for the variances.

For ARMA models, users need to set the frequency (annual -1-, quarterly -4-, and monthly -12-), as well as the AR and MA orders (see Figure 5.9). Then, users should set the location and scale hyperparameters for the intercept, autoregressive (AR), moving average (MA), and standard deviation terms. Note that there is only one set of hyperparameters for the AR and MA coefficients. This step is optional, as the GUI uses non-informative priors by default.

Then, users should click the *Go!* button, and the GUI will start estimating the model. The GUI will display the summary statistics of the posterior draws and the convergence diagnostics. The order is AR coefficients (if any), MA coefficients (if any), intercept, and standard deviation. Users can download the posterior chains and figures (density, autocorrelation, and trace plots).

Estimation of the SVMs requires setting the coefficients of the mean and standard deviation of the Gaussian prior for the regression parameters, the mean and standard deviation for the Gaussian prior distribution of the level of the log-volatility, shape parameters for the Beta prior distribution of the transformed persistence parameter, and a positive real number representing the scaling of the transformed volatility of log-volatility. However, this step is not necessary, as by default our GUI uses the default values in the *stochvol* package.

Then, click the *Go!* button, wait for the estimation to be completed, and the GUI will display the stochastic volatility plot (mean and 95% credible interval). Users can also view the summary and diagnostics of the posterior chains (see Figure 5.10). In addition, users can download the mean and the lower and upper limits of the 95% credible intervals of the stochastic volatility, as well as the posterior chains of the variances.

**FIGURE 5.10**

Time series models: Stochastic volatility results.

To estimate VAR models, users should set the number of lags, the impulse response and forecast periods, the three coefficients of the Minnesota prior, and the type of impulse response (*forecast error impulse response* -feir- or *orthogonalized impulse response* -other-, both cumulative or non-cumulative). See Chapter 8 for details, Section 8.4.

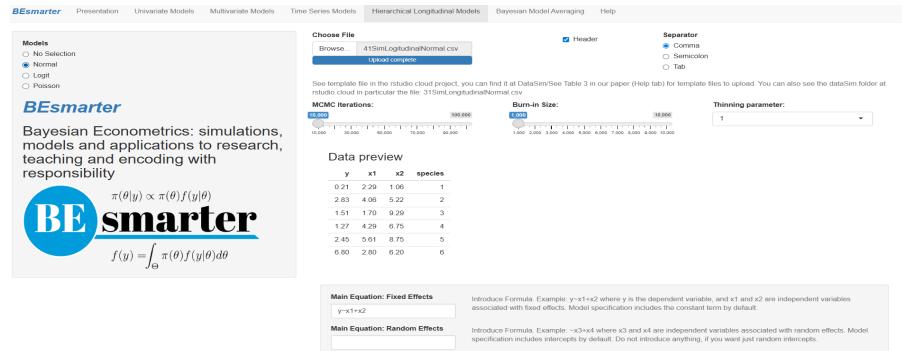
Click the *Go!* button, and after a few minutes, users will be able to see the plots of the impulse responses and forecasts (means and 95% credible intervals). Click the *Download Results* button, and a zip file with .csv files containing the impulse responses and forecasts, along with their plots, will be downloaded.

## 5.5 Longitudinal/panel models

After our GUI is deployed (see Figure 5.1), the user should select *Hierarchical Longitudinal Models* in the top panel. Then, Figure 5.11 will be displayed, and the user can see the radio button on the left-hand side that shows the specific models inside this generic class.

The hierarchical longitudinal models tab allows for estimating models that account for within-subject correlation when the dependent variable is continuous (Normal), binary (Logit), or a count (Poisson).

The input files for hierarchical longitudinal models should first include the dependent variable, followed by the regressors and a cross-sectional identifier ( $i = 1, 2, \dots, N$ ). It is not a requirement to have a balanced dataset:  $T_i$  can be different for each  $i$  (see Chapter 9 for technical details). Users can see templates of datasets in the folders **DataSim** (see Table A.2 for details) and



**FIGURE 5.11**  
Hierarchical longitudinal models: Specification.

**DataApp** (see Table A.3 for details) in our *GitHub* repository. When the dataset is uploaded, users will have a preview of it.

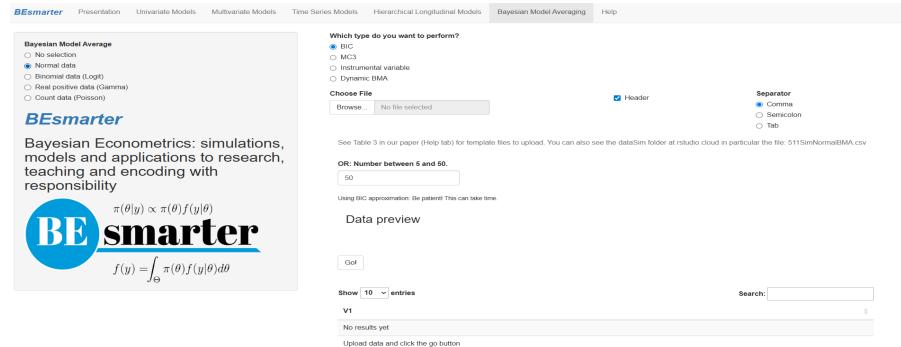
Users should also specify the fixed part equation and the random part equation, both in **R** format. If only random intercepts are required, do not enter anything in the latter part (see Figure 5.11). Users should also type the name of the cross-sectional identifier variable (*species* in Figure 5.11). The results displayed and the posterior graphs are associated with the fixed effects and covariance matrix. However, users can download the posterior chains of all posterior estimates: fixed and random effects, and the covariance matrix.

## 5.6 Bayesian model average

After our GUI is deployed (see Figure 5.1), the user should select *Bayesian Model Averaging* in the top panel. Then, Figure 5.12 will be displayed, and the user can see the radio button on the left-hand side that shows the specific models inside this generic class.

Bayesian model averaging (BMA) based on a Gaussian distribution can be carried out using the Bayesian information criterion (BIC) approximation, Markov chain Monte Carlo model composition (MC3), instrumental variables, and dynamic BMA (see Figure 5.12). The first two approaches require an input dataset where the first column is the dependent variable, followed by the potentially important regressors.

Users should set the bandwidth model selection parameter ( $O_R$ ) and the number of iterations for BIC and MC3, respectively (see Chapter 10 for technical details). The results include the posterior inclusion probability ( $p \neq 0$ ), expected value (EV), and standard deviation (SD) of the coefficients associ-

**FIGURE 5.12**

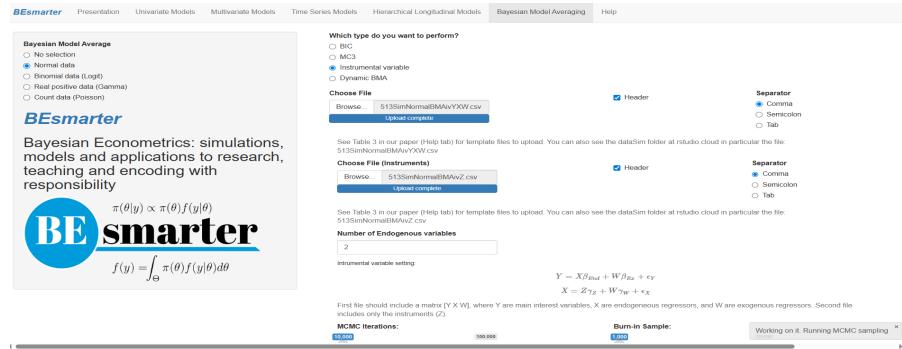
Bayesian model averaging: Specification and results.

ated with each regressor. The BIC framework also displays the most relevant models with their posterior model probabilities (PMP). Users can download two *csv* files: *Best models* and *Descriptive statistics coefficients*. The former is a 0-1 matrix such that the columns are the regressors and the rows are the models; a 1 indicates the presence of a specific regressor in a specific model, and 0 indicates its absence. Note that the last column of this file is the posterior model probability for each model (row). The latter file shows the posterior inclusion probabilities, expected values, and standard deviations associated with each regressor, taking into account the BMA procedure based on the best models.

Bayesian model averaging with endogeneity issues requires two input files. The first file should have the dependent variable in the first column, followed by the regressors with endogeneity issues, and then the exogenous regressors. The user should include a column of 1's if an intercept is required. The second input file contains all the instruments. Users should also specify the number of regressors with endogeneity issues (see Figure 5.13).

The results include the posterior inclusion probabilities and expected values for each regressor. The user can find the results of the main equation, and then of the auxiliary equations. Users can download *csv* files of BMA results for both the second stage (main equation) and the first stage (auxiliary equations). In addition, users can download the posterior chains of the location parameters of the main equation,  $\beta_l$ ,  $l = 1, 2, \dots, \dim\{\beta\}$ , the location parameters of the auxiliary equations,  $\gamma_{j,i}$ ,  $j = 1, 2, \dots, \dim\{\beta_s\}$  where  $\dim\{\beta_s\}$  is the number of regressors with endogeneity issues,  $i = 1, 2, \dots, \dim\{\gamma\}$ , where  $\dim\{\gamma\}$  is the number of regressors in the auxiliary regressors (exogeneous regressors + instruments), and the elements of the covariance matrix  $\sigma_{j,k}$  (see Chapter 10 for technical details).

Dynamic BMA also requires two files. The first is the dataset with the dependent variable and potential regressors, and the second file describes the

**FIGURE 5.13**

Bayesian model averaging: Instrumental variable specification.

competing models. There is one column for each regressor and one row for each competing model; 0 indicates that the regressor is not in the model, and 1 indicates that it is in the model. Users can see templates of this file in the folders **DataSim** (see Table A.2 for details) and **DataApp** (see Table A.3 for details) of our *GitHub* repository.

Then, the users should set the *forgetting parameters* of the covariance and transition matrices and click the *Go!* button. A plot of the PMPs of the competing models is displayed, and users can click the *Download the results for DBMA*. Two files are downloaded: the first contains the dynamic Bayesian average filtering recursions for each state, and the second contains the PMP of each model and the dynamic Bayesian model averaging prediction.

Bayesian model averaging based on BIC approximation for non-linear models (Logit, Gamma, and Poisson) requires an input dataset where the first column is the dependent variable, and the other columns are the potentially relevant regressors. Users should specify the bandwidth model selection parameters, also referred to as Occam's window parameters ( $O_R$  and  $O_L$ ). Our GUI displays the posterior inclusion probabilities ( $p! = 0$ ), the expected value of the posterior coefficients (EV), and the standard deviation (SD). In addition, users can view the results associated with the models with the highest posterior model probabilities and download *csv* files with the results of specifications of the best models and descriptive statistics of the posterior coefficients from the BMA procedure. These files are similar to the results of the BIC approximation for the Gaussian model.

## 5.7 Help

The last tab in our GUI is *Help*. There, users can find links to both the online *HTML* version of the book,<sup>4</sup> and the *PDF* version.<sup>5</sup> Users are also welcome to contact me at *aramir21@gmail.com* with any questions, comments, or suggestions.

## 5.8 Warning

Users should also note that sometimes our GUI shuts down. In our experience, this is due to computational issues arising from the implicit commands we call when estimating certain models. These issues may include computationally singular systems, missing values where TRUE/FALSE are needed, L-BFGS-B requiring finite values for “fn”, NA/Nan/Inf values, or errors in `backsolve`. These issues can sometimes be resolved by adjusting the dataset, such as avoiding high levels of multicollinearity.

It should also be noted that when warning messages are displayed in our GUI, there is a high likelihood of convergence issues with the posterior chains. Therefore, the results may not be trustworthy. Users can identify these problems by checking the console in their *RStudio* sections, where the specific folder/file where the issue occurred will be specified. In any case, we would appreciate your feedback to improve and enhance our GUI.

We should also mention that there are many ways to improve the codes presented in this book. For instance, the *MCMCpack* and *bayesm* packages perform most of the matrix operations in C++ using the *Rcpp* package. This substantially speeds up the algorithms compared to the codes presented in the next chapters when we program from scratch the samplers. We could further improve the computational times of our codes using parallel computing and the *Rcpp* package, but this requires more advanced skills that are not covered in this book.

---

<sup>4</sup><https://bookdown.org/aramir21/IntroductionBayesianEconometricsGuidedTour/>

<sup>5</sup>[https://drive.google.com/file/d/1JBKe7vS5a2XLnvg74T\\_UNESRoKNDUUt/edit](https://drive.google.com/file/d/1JBKe7vS5a2XLnvg74T_UNESRoKNDUUt/edit)



# 6

## Univariate models

In this chapter, we describe how to perform Bayesian inference in several of the most common univariate models, including the normal, logit, probit, multinomial probit and logit, ordered probit, negative binomial, tobit, quantile regression, and Bayesian bootstrap models for linear regression. Our point of departure is a random sample of cross-sectional units. We then present the posterior distributions of the parameters and illustrate their use through selected applications. In addition, we demonstrate how to perform inference in these models using three levels of programming engagement: through our graphical user interface (GUI), via **R** packages, and by directly programming the posterior distributions. The first requires no programming skills, the second requires an intermediate level, and the third demands advanced skills. We also provide mathematical and computational exercises.

We can run our GUI typing

### *R code. How to display our graphical user interface*

```
1 shiny::runGitHub("besmarter/BSTApp", launch.browser = T)
```

in the **R** package console or any **R** code editor. However, users should see Chapter 5 for details.

### 6.1 The Gaussian linear model

The Gaussian linear model is specified as  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu}$ , where  $\boldsymbol{\mu} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_N)$  is a stochastic error term,  $\mathbf{X}$  is an  $N \times K$  matrix of regressors,  $\boldsymbol{\beta}$  is a  $K$ -dimensional vector of location parameters,  $\sigma^2$  is the variance of the model (scale parameter),  $\mathbf{y}$  is an  $N$ -dimensional vector of the dependent variable, and  $N$  is the sample size. In Section 3.3, we describe this model under the

conjugate prior family, that is,  $\pi(\boldsymbol{\beta}, \sigma^2) = \pi(\boldsymbol{\beta} | \sigma^2) \times \pi(\sigma^2)$ , which allows us to derive the marginal posterior distributions of  $\boldsymbol{\beta}$  and  $\sigma^2$ .

In this section, we assume independent priors, that is,  $\pi(\boldsymbol{\beta}, \sigma^2) = \pi(\boldsymbol{\beta}) \times \pi(\sigma^2)$ , where  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$  and  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ , with  $\alpha_0/2$  and  $\delta_0/2$  denoting the shape and rate parameters, respectively. This specification allows us to obtain the conditional posterior distributions,  $\pi(\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X})$  and  $\pi(\sigma^2 | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X})$ , which can be used to implement the Gibbs sampler for posterior inference on  $\boldsymbol{\beta}$  and  $\sigma^2$ .

The likelihood function of this model is given by

$$p(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbf{X}) = (2\pi\sigma^2)^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}.$$

The conditional posterior distributions are

$$\begin{aligned} \boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} &\sim N(\boldsymbol{\beta}_n, \mathbf{B}_n), \\ \sigma^2 | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X} &\sim IG(\alpha_n/2, \delta_n/2), \end{aligned}$$

where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \sigma^{-2} \mathbf{X}^\top \mathbf{X})^{-1}$ ,  $\boldsymbol{\beta}_n = \mathbf{B}_n (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sigma^{-2} \mathbf{X}^\top \mathbf{y})$ ,  $\alpha_n = \alpha_0 + N$ , and  $\delta_n = \delta_0 + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  (see Exercise 1 in this chapter).<sup>1</sup>

### Example: The market value of soccer players in Europe

Let's analyze the determinants of the market value of soccer players in Europe. In particular, we use the dataset `1ValueFootballPlayers.csv`, located in the **DataApp** folder of our GitHub repository <https://github.com/besmarter/BSTApp>. This dataset was originally used by [406] to identify the determinants of high-performance soccer players in the five major European national leagues.

The model specification is as follows:

$$\begin{aligned} \log(\text{Value}_i) &= \beta_1 + \beta_2 \text{Perf}_i + \beta_3 \text{Age}_i + \beta_4 \text{Age}_i^2 + \beta_5 \text{NatTeam}_i \\ &\quad + \beta_6 \text{Goals}_i + \beta_7 \text{Exp}_i + \beta_8 \text{Exp}_i^2 + \mu_i, \end{aligned}$$

where *Value* is the market value in euros (2017), *Perf* is a measure of player performance, *Age* is the player's age in years, *NatTeam* is an indicator variable that equals 1 if the player has been a member of the national team, *Goals* is the total number of goals scored during the player's career, and *Exp* represents playing experience in years.

We assume that the dependent variable follows a normal distribution. Accordingly, we estimate a normal-inverse-gamma model using vague conjugate priors with  $\mathbf{B}_0 = 1000\mathbf{I}_8$ ,  $\boldsymbol{\beta}_0 = \mathbf{0}_8$ ,  $\alpha_0 = 0.001$ , and  $\delta_0 = 0.001$ . Posterior inference is performed using a Gibbs sampler with 5,000 MCMC iterations, a burn-in period of 5,000, and a thinning parameter of 1.

<sup>1</sup>This model can be extended to account for heteroskedasticity such that  $y_i \sim N(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2/\tau_i)$ , where  $\tau_i \sim G(v/2, v/2)$ . See Exercise 2 for details.

Once the GUI is launched (see the beginning of this chapter), users can follow Algorithm A6 to run linear Gaussian models in the interface (see Chapter 5 for further details).

---

**Algorithm A6** Gaussian linear model

---

- 1: Select *Univariate Models* on the top panel
  - 2: Select *Normal* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select dependent and independent variables using the *Formula builder* table
  - 6: Click the *Build formula* button to generate the formula in **R** syntax. You can modify the formula in the **Main equation** box using valid arguments of the *formula* command structure in **R**
  - 7: Set the hyperparameters: mean vector, covariance matrix, shape and scale parameters. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Click the *Go!* button
  - 9: Analyze results
  - 10: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

The following **R** code examples illustrate how to estimate the linear Gaussian model using the *MCMCregress* command from the *MCMCpack* package, as well as how to implement the Gibbs sampler manually. We should obtain similar results across all three approaches —GUI, package, and custom function— since our GUI internally relies on the *MCMCregress* command. For instance, the market value of a top soccer player in Europe increases by approximately 134% ( $\exp(0.85) - 1$ ) on average when he has played for the national team, with a 95% credible interval of (86%, 197%).

**R code. The value of soccer players, using R packages**

```

1 rm(list = ls()); set.seed(10101)
2 # Linear regression: Value of soccer players #
3 Data <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/1
  ValueFootballPlayers.csv", sep = ",", header = TRUE,
  quote = "")
4 attach(Data)
5 y <- log(Value)
6 # Value: Market value in Euros (2017) of soccer players
7 # Regressors quantity including intercept
8 X <- cbind(1, Perf, Age, Age2, NatTeam, Goals, Exp, Exp2)
9 # Perf: Performance. Perf2: Performance squared. Age: Age;
  Age: Age squared.
10 # NatTeam: Indicator of national team. Goals: Scored goals.
  Goals2: Scored goals squared
11 # Exp: Years of experience. Exp2: Years of experience
  squared. Assists: Number of assists
12 k <- dim(X)[2]
13 N <- dim(X)[1]
14 # Hyperparameters
15 d0 <- 0.001/2
16 a0 <- 0.001/2
17 b0 <- rep(0, k)
18 c0 <- 1000
19 B0 <- c0*diag(k)
20 B0i <- solve(B0)
21 # MCMC parameters
22 mcmc <- 5000
23 burnin <- 5000
24 tot <- mcmc + burnin
25 thin <- 1
26 # Posterior distributions using packages: MCMCpack sets the
  model in terms of the precision matrix
27 posterior <- MCMCpack::MCMCregress(y~X-1, b0=b0, B0 = B0i,
  c0 = a0, d0 = d0, burnin = burnin, mcmc = mcmc, thin =
  thin)
28 summary(coda::mcmc(posterior))
29 # Posterior distributions programming the Gibbs sampling
30 # Auxiliary parameters
31 XtX <- t(X)%*%X
32 bhat <- solve(XtX)%*%t(X)%*%y
33 an <- a0 + N

```

**R. code. The value of soccer players, programming our Gibbs sampler**

```

1 # Gibbs sampling functions
2 PostSig2 <- function(Beta){
3   dn <- d0 + t(y - X%*%Beta)%*%(y - X%*%Beta)
4   sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
5   return(sig2)
6 }
7 PostBeta <- function(sig2){
8   Bn <- solve(B0i + sig2^(-1)*XtX)
9   bn <- Bn%*%(B0i%*%b0 + sig2^(-1)*XtX%*%bhat)
10  Beta <- MASS::mvrnorm(1, bn, Bn)
11  return(Beta)
12 }
13 PostBetas <- matrix(0, mcmc+burnin, k)
14 PostSigma2 <- rep(0, mcmc+burnin)
15 Beta <- rep(0, k)
16 for(s in 1:tot){
17   sig2 <- PostSig2(Beta = Beta)
18   PostSigma2[s] <- sig2
19   Beta <- PostBeta(sig2 = sig2)
20   PostBetas[s,] <- Beta
21 }
22 keep <- seq((burnin+1), tot, thin)
23 PosteriorBetas <- PostBetas[keep,]
24 colnames(PosteriorBetas) <- c("Intercept", "Perf", "Age", "
  Age2", "NatTeam", "Goals", "Exp", "Exp2")
25 summary(coda::mcmc(PosteriorBetas))
26 PosteriorSigma2 <- PostSigma2[keep]
27 summary(coda::mcmc(PosteriorSigma2))

```

## 6.2 The logit model

In the logit model, the dependent variable is binary,  $y_i \in \{0, 1\}$ , and follows a Bernoulli distribution,  $y_i \stackrel{\text{ind}}{\sim} B(\pi_i)$ , where  $p(y_i = 1) = \pi_i$  and

$$\pi_i = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}.$$

Here,  $\mathbf{x}_i$  is a  $K$ -dimensional vector of regressors.

The likelihood function of the logit model is given by

$$\begin{aligned} p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{X}) &= \prod_{i=1}^N \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \\ &= \prod_{i=1}^N \left( \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})} \right)^{y_i} \left( \frac{1}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})} \right)^{1-y_i}. \end{aligned}$$

We specify a normal prior distribution for the coefficients,  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$ . The posterior distribution is therefore

$$\begin{aligned} \pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) &\propto \prod_{i=1}^N \left( \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})} \right)^{y_i} \left( \frac{1}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})} \right)^{1-y_i} \\ &\quad \times \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \right\}. \end{aligned}$$

Since the logit model does not yield a standard posterior distribution, a random-walk Metropolis–Hastings algorithm can be employed to obtain draws from the posterior. A suitable proposal distribution is a multivariate normal centered at the current draw, with covariance matrix  $\tau^2 (\mathbf{B}_0^{-1} + \widehat{\boldsymbol{\Sigma}}^{-1})^{-1}$ , where  $\tau > 0$  is a tuning parameter and  $\widehat{\boldsymbol{\Sigma}}$  is the sample covariance matrix obtained from the maximum likelihood estimates [294].<sup>2</sup>

The log-likelihood function can be expressed as

$$\log(p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{X})) = \sum_{i=1}^N [y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \log(1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta}))].$$

This expression is useful for computing the acceptance probability in the Metropolis–Hastings algorithm, given by

$$\begin{aligned} \alpha &= \min \left\{ 1, \exp \left( \log(p(\mathbf{y} | \boldsymbol{\beta}^c, \mathbf{X})) + \log(\pi(\boldsymbol{\beta}^c)) \right. \right. \\ &\quad \left. \left. - \left[ \log(p(\mathbf{y} | \boldsymbol{\beta}^{(s-1)}, \mathbf{X})) + \log(\pi(\boldsymbol{\beta}^{(s-1)})) \right] \right) \right\}, \end{aligned}$$

where  $\boldsymbol{\beta}^c$  and  $\boldsymbol{\beta}^{(s-1)}$  denote the candidate draw from the proposal distribution and the previous iteration of the Markov chain, respectively.<sup>3</sup>

### Example: Simulation exercise

Let's conduct a simulation exercise to evaluate the performance of the algorithm. Set  $\boldsymbol{\beta} = [0.5 \ 0.8 \ -1.2]^\top$ , with  $x_{ik} \sim N(0, 1)$  for  $k = 2, 3$  and  $i = 1, 2, \dots, 10,000$ .

---

<sup>2</sup>Tuning parameters should be chosen to ensure satisfactory convergence diagnostics and acceptance rates.

<sup>3</sup>Expressing the acceptance probability in logarithmic form helps mitigate numerical underflow and overflow issues during computation.

We specify the hyperparameters as  $\beta_0 = [0 \ 0 \ 0]^\top$  and  $B_0 = 1000\mathbf{I}_3$ . The tuning parameter for the Metropolis–Hastings algorithm is set to 1.

Once the GUI is launched (see the beginning of this chapter), users should follow Algorithm A7 to estimate logit models in the interface (see Chapter 5 for further details).

---

**Algorithm A7** Logit model
 

---

- 1: Select *Univariate Models* on the top panel
  - 2: Select *Logit* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select dependent and independent variables using the *Formula builder* table
  - 6: Click the *Build formula* button to generate the formula in **R** syntax. You can modify the formula in the **Main equation** box using valid arguments of the *formula* command structure in **R**
  - 7: Set the hyperparameters: mean vector and covariance matrix. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Select the tuning parameter for the Metropolis-Hastings algorithm. This step is not necessary as by default our GUI sets the tuning parameter at 1.1
  - 9: Click the *Go!* button
  - 10: Analyze results
  - 11: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

The following **R** code illustrates how to estimate the logit model using the *MCMClogit* command from the *MCMCpack* package, as well as by implementing the Metropolis–Hastings algorithm manually.

We should obtain similar results across the three approaches —GUI, package, and custom function— since our GUI internally relies on the *MCMClogit* command. In this example, the acceptance rate is 0.46, and the diagnostics indicate that the posterior chains exhibit satisfactory convergence behavior. Overall, the 95% credible intervals contain the true population values, and both the posterior means and medians are very close to those values.

### R. code. Simulation of the logit model estimation using R packages

```

1 ##### Logit: Simulation #####
2 # Simulate data
3 rm(list = ls()); set.seed(10101)
4 N <- 10000 # Sample size
5 B <- c(0.5, 0.8, -1.2) # Population location parameters
6 x2 <- rnorm(N) # Regressor
7 x3 <- rnorm(N) # Regressor
8 X <- cbind(1, x2, x3) # Regressors
9 XB <- X%*%B
10 PY <- exp(XB)/(1 + exp(XB)) # Probability of Y = 1
11 Y <- rbinom(N, 1, PY) # Draw Y's
12 table(Y) # Frequency
13 # write.csv(cbind(Y, x2, x3), file = "DataSimulations/
    LogitSim.csv") # Export data
14 # MCMC parameters
15 iter <- 5000; burnin <- 1000; thin <- 5; tune <- 1
16 # Hyperparameters
17 K <- dim(X)[2]
18 b0 <- rep(0, K)
19 c0 <- 1000
20 B0 <- c0*diag(K)
21 B0i <- solve(B0)
22 # Posterior distributions using packages: MCMCpack sets the
    model in terms of the precision matrix
23 RegLog <- MCMCpack::MCMClogit(Y~X-1, mcmc = iter, burnin =
    burnin, thin = thin, b0 = b0, B0 = B0i, tune = tune)
24 summary(RegLog)
25 # Metropolis-Hastings algorithm
26 MHfunc <- function(y, X, b0 = rep(0, dim(X)[2] + 1), B0 =
    1000*diag(dim(X)[2] + 1), tau = 1, iter = 6000, burnin =
    1000, thin = 5){
27   Xm <- cbind(1, X) # Regressors
28   K <- dim(Xm)[2] # Number of location parameters
29   BETAS <- matrix(0, iter + burnin, K) # Space for posterior
    chains
30   Reg <- glm(y ~ Xm - 1, family = binomial(link = "logit"))
    # Maximum likelihood estimation
31   BETA <- Reg$coefficients # Maximum likelihood parameter
    estimates
32   tot <- iter + burnin # Total iterations M-H algorithm
33   COV <- vcov(Reg) # Maximum likelihood covariance matrix
34   COVt <- tau^2*solve(solve(B0) + solve(COV)) # Covariance
    matrix for the proposal distribution
35   Accep <- rep(0, tot) # Space for calculating the
    acceptance rate
36   pb <- txtProgressBar(min = 0, max = tot, style = 3)

```

**R. code. Simulation of the logit model estimation  
programming our M-H algorithm**

```

1   for(it in 1:tot){
2     BETAc <- BETA + MASS::mvrnorm(n = 1, mu = rep(0, K),
3       Sigma = COVt) # Candidate location parameter
4     likecand <- sum((Xm%*%BETAc) * Y - apply(Xm%*%BETAc, 1,
5       function(x) log(1 + exp(x)))) # Log likelihood candidate
6     likepast <- sum((Xm%*%BETA) * Y - apply((Xm%*%BETA), 1,
7       function(x) log(1 + exp(x)))) # Log likelihood actual
8     draw
9     priorcand <- (-1/2)*crossprod((BETAc - b0), solve(B0))%*
10    %(BETAc - b0) # Log prior candidate
11    priorpast <- (-1/2)*crossprod((BETA - b0), solve(B0))%*%
12    (BETA - b0) # Log prior actual draw
13    alpha <- min(1, exp((likecand + priorcand) - (likepast +
14      priorpast))) #Probability of selecting candidate
15    u <- runif(1) # Decision rule for selecting candidate
16    if(u < alpha){
17      BETA <- BETAc # Changing reference for candidate
18      Accep[it] <- 1 # Indicator candidate is accepted
19    }
20    close(pb)
21    keep <- seq(burnin, tot, thin)
22    return(list(Bs = BETAS[keep[-1], ], AceptRate = mean(Accep
23      [keep[-1]])))
24  }
25  Posterior <- MHfunc(y = Y, X = cbind(x2, x3), iter = iter,
26    burnin = burnin, thin = thin) # Running our M-H function
27    changing some default parameters.
28  paste("Acceptance rate equal to", round(Posterior$AceptRate,
29    2), sep = " ")
30  "Acceptance rate equal to 0.46"
31  PostPar <- coda::mcmc(Posterior$Bs)
32  # Names
33  colnames(PostPar) <- c("Cte", "x1", "x2")
34  # Summary posterior draws
35  summary(PostPar)
36  # Trace and density plots
37  plot(PostPar)
38  # Autocorrelation plots
39  coda::autocorr.plot(PostPar)
40  # Convergence diagnostics
41  coda::geweke.diag(PostPar)
42  coda::raftery.diag(PostPar, q=0.5, r=0.05, s = 0.95)
43  coda::heidel.diag(PostPar)

```

### 6.3 The probit model

The probit model also features a binary dependent variable. In this case, an unobserved latent variable ( $y_i^*$ ) defines the underlying structure of the estimation problem.

Specifically,

$$y_i = \begin{cases} 0, & y_i^* \leq 0 \\ 1, & y_i^* > 0 \end{cases},$$

such that  $y_i^* = \mathbf{x}_i^\top \boldsymbol{\beta} + \mu_i$ , where  $\mu_i \stackrel{i.i.d.}{\sim} N(0, 1)$ .<sup>4</sup> This implies that  $P(y_i = 1) = \pi_i = \Phi(\mathbf{x}_i^\top \boldsymbol{\beta})$ , where  $\mathbf{x}_i$  is a  $K$ -dimensional vector of regressors.

[5] introduced a data augmentation approach [428] to implement a Gibbs sampling algorithm for this model. By augmenting the data with the latent variable  $y_i^*$ , the likelihood contribution from observation  $i$  can be written as

$$p(y_i | y_i^*) = \mathbb{1}(y_i = 0)\mathbb{1}(y_i^* \leq 0) + \mathbb{1}(y_i = 1)\mathbb{1}(y_i^* > 0),$$

where  $\mathbb{1}(A)$  is an indicator function that equals 1 when the condition  $A$  is satisfied.

The joint posterior distribution is given by

$$\begin{aligned} \pi(\boldsymbol{\beta}, \mathbf{y}^* | \mathbf{y}, \mathbf{X}) &\propto \prod_{i=1}^N [\mathbb{1}(y_i = 0)\mathbb{1}(y_i^* \leq 0) + \mathbb{1}(y_i = 1)\mathbb{1}(y_i^* > 0)] \\ &\quad \times N_N(\mathbf{y}^* | \mathbf{X}\boldsymbol{\beta}, \mathbf{I}_N) \times N_K(\boldsymbol{\beta} | \boldsymbol{\beta}_0, \mathbf{B}_0), \end{aligned}$$

where we assume a Gaussian prior for  $\boldsymbol{\beta}$ :  $\boldsymbol{\beta} \sim N_K(\boldsymbol{\beta}_0, \mathbf{B}_0)$ .

The corresponding conditional posterior distributions are:

$$y_i^* | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X} \sim \begin{cases} TN_{(-\infty, 0]}(\mathbf{x}_i^\top \boldsymbol{\beta}, 1), & y_i = 0, \\ TN_{(0, \infty)}(\mathbf{x}_i^\top \boldsymbol{\beta}, 1), & y_i = 1, \end{cases}$$

$$\boldsymbol{\beta} | \mathbf{y}^*, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \mathbf{B}_n),$$

where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}$  and  $\boldsymbol{\beta}_n = \mathbf{B}_n(\mathbf{B}_0^{-1}\boldsymbol{\beta}_0 + \mathbf{X}^\top \mathbf{y}^*)$ .

#### Example: Determinants of hospitalization

We use the dataset **2HealthMed.csv**, located in the **DataApp** folder of our GitHub repository <https://github.com/besmarter/BSTApp>, which was previously analyzed by [354]. The dependent variable is a binary indicator that equals 1 if an individual was hospitalized in 2007 and 0 otherwise.

---

<sup>4</sup>The variance in this model is normalized to one due to identification restrictions. Note that  $P(y_i = 1 | \mathbf{x}_i) = P(y_i^* > 0 | \mathbf{x}_i) = P(\mathbf{x}_i^\top \boldsymbol{\beta} + \mu_i > 0 | \mathbf{x}_i) = P(\mu_i > -\mathbf{x}_i^\top \boldsymbol{\beta} | \mathbf{x}_i) = P(c \times \mu_i > -c \times \mathbf{x}_i^\top \boldsymbol{\beta} | \mathbf{x}_i)$  for all  $c > 0$ . Multiplying by a positive constant does not affect the probability of  $y_i = 1$ .

The model specification is as follows:

$$\text{Hosp}_i = \beta_1 + \beta_2 \text{SHI}_i + \beta_3 \text{Female}_i + \beta_4 \text{Age}_i + \beta_5 \text{Age}_i^2 + \beta_6 \text{Est2}_i + \beta_7 \text{Est3}_i + \beta_8 \text{Fair}_i + \beta_9 \text{Good}_i + \beta_{10} \text{Excellent}_i,$$

where  $\text{SHI}$  is a binary variable equal to 1 if the individual is enrolled in a subsidized health care program and 0 otherwise,  $\text{Female}$  is a gender indicator,  $\text{Age}$  is measured in years,  $\text{Est2}$  and  $\text{Est3}$  are indicators of socioeconomic status (with  $\text{Est1}$  as the reference category representing the lowest status), and  $\text{HealthStatus}$  is the individual's self-perceived health status, where  $\text{Bad}$  serves as the reference category.

We set  $\beta_0 = \mathbf{0}_{10}$  and  $\mathbf{B}_0 = \mathbf{I}_{10}$ , with 10,000 iterations, a burn-in period of 1,000, and a thinning parameter of 1. Algorithm A8 can be used to estimate the probit model within our GUI. The GUI internally relies on the *rbprobitGibbs* command from the *bayesm* package to perform Bayesian inference for the probit model. The following **R** code illustrates how to run this example using the *rbprobitGibbs* command. In the exercises, you are also asked to implement a Gibbs sampling algorithm to perform inference in the probit model manually.

---

#### Algorithm A8 Probit model

---

- 1: Select *Univariate Models* on the top panel
  - 2: Select *Probit* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select dependent and independent variables using the *Formula builder* table
  - 6: Click the *Build formula* button to generate the formula in **R** syntax. You can modify the formula in the **Main equation** box using valid arguments of the *formula* command structure in **R**
  - 7: Set the hyperparameters: mean vector and covariance matrix. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Click the *Go!* button
  - 9: Analyze results
  - 10: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

Our analysis finds evidence that gender and self-perceived health status significantly affect the probability of hospitalization. Women have a higher probability of being hospitalized than men, and individuals with a better perception of their health status have a lower probability of hospitalization.

### R. code. Determinants of hospitalization

```

1 mydata <- read.csv("https://raw.githubusercontent.com/
  besmaster/BSTApp/refs/heads/master/DataApp/2HealthMed.
  csv", sep = ",", header = TRUE, quote = "")
2 attach(mydata); str(mydata)
3 K <- 10 ; b0 <- rep(0, K) # Regressors and prior mean
4 B0i <- diag(K) # Prior precision (inverse of covariance)
5 Prior <- list(betabar = b0, A = B0i) # Prior list
6 y <- Hosp # Dependent variables
7 X <- cbind(1, SHI, Female, Age, Age2, Est2, Est3, Fair, Good
  , Excellent) # Regressors
8 Data <- list(y = y, X = X) # Data list
9 Mcmc <- list(R = 10000, keep = 1, nprint = 0) # MCMC
10 RegProb <- bayesm::rbprobitGibbs(Data = Data, Prior = Prior,
  Mcmc = Mcmc) # Inference using bayesm package
11 PostPar <- coda::mcmc(RegProb$betadraw) # Posterior draws
12 colnames(PostPar) <- c("Cte", "SHI", "Female", "Age", "Age2"
  , "Est2", "Est3", "Fair", "Good", "Excellent") # Names
13 summary(PostPar) # Posterior summary

```

## 6.4 The multinomial probit model

The multinomial probit model is used to explain the choice of the  $l$ -th alternative from a set of  $L$  mutually exclusive options. We observe

$$y_{il} = \begin{cases} 1, & \text{if } y_{il}^* \geq \max\{\mathbf{y}_i^*\}, \\ 0, & \text{otherwise,} \end{cases}$$

where  $\mathbf{y}_i^* = \mathbf{X}_i \boldsymbol{\delta} + \boldsymbol{\mu}_i$ , and  $\boldsymbol{\mu}_i \stackrel{i.i.d.}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$ . The vector  $\mathbf{y}_i^*$  is an unobserved latent vector of dimension  $L$ . The matrix  $\mathbf{X}_i = [(1 \mathbf{c}_i^\top) \otimes \mathbf{I}_L \mathbf{A}_i]$  is an  $L \times j$  matrix of regressors for each alternative, where  $l = 1, 2, \dots, L$ , and  $j = L \times (1 + \dim(\mathbf{c}_i)) + a$ . Here,  $\mathbf{c}_i$  is a vector of individual-specific characteristics,  $\mathbf{A}_i$  is an  $L \times a$  matrix of alternative-varying regressors,  $a$  denotes the number of alternative-varying regressors, and  $\boldsymbol{\delta}$  is a  $j$ -dimensional vector of parameters.

The model jointly accounts for alternative-varying regressors (alternative attributes) and alternative-invariant regressors (individual characteristics).<sup>5</sup>

---

<sup>5</sup>This model is not identified if  $\boldsymbol{\Sigma}$  is unrestricted. The likelihood function remains unchanged if a scalar random variable is added to each of the  $L$  latent regressions.

The vector  $\mathbf{y}_i^*$  can be stacked into a multivariate regression model with correlated stochastic errors, i.e.,  $\mathbf{y}^* = \mathbf{X}\boldsymbol{\delta} + \boldsymbol{\mu}$ , where  $\mathbf{y}^* = [\mathbf{y}_1^{*\top}, \mathbf{y}_2^{*\top}, \dots, \mathbf{y}_N^{*\top}]^\top$ ,  $\mathbf{X} = [\mathbf{X}_1^\top, \mathbf{X}_2^\top, \dots, \mathbf{X}_N^\top]^\top$ , and  $\boldsymbol{\mu} = [\boldsymbol{\mu}_1^\top, \boldsymbol{\mu}_2^\top, \dots, \boldsymbol{\mu}_N^\top]^\top$ .

Following standard practice, we express  $y_{il}^*$  relative to  $y_{iL}^*$  by defining  $\mathbf{w}_i = [w_{i1}, w_{i2}, \dots, w_{iL-1}]^\top$ , where  $w_{il} = y_{il}^* - y_{iL}^*$ . Then,  $\mathbf{w}_i = \mathbf{R}_i\boldsymbol{\beta} + \boldsymbol{\epsilon}_i$ , with  $\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \boldsymbol{\Omega})$ , where  $\mathbf{R}_i = [(1 \mathbf{c}_i^\top) \otimes \mathbf{I}_{L-1} \Delta \mathbf{A}_i]$  is an  $(L-1) \times K$  matrix,  $\Delta \mathbf{A}_i = \mathbf{A}_{li} - \mathbf{A}_{Li}$  for  $l = 1, 2, \dots, L-1$ , and  $\boldsymbol{\beta}$  is a  $K$ -dimensional parameter vector, with  $K = (L-1) \times (1 + \dim(\mathbf{c}_i)) + a$ . That is, the last row of  $\mathbf{A}_i$  is subtracted from each row to normalize the model.

The vector  $\boldsymbol{\beta}$  contains the same last  $a$  elements as  $\boldsymbol{\delta}$ , corresponding to the coefficients on the alternative-specific attributes. However, the first  $(L-1) \times (1 + \dim(\mathbf{c}_i))$  elements of  $\boldsymbol{\beta}$  are differences of the form  $\delta_{jl} - \delta_{jL}$ , for  $j = 1, \dots, \dim(\mathbf{c}_i)$  and  $l = 1, \dots, L-1$ . These represent the differences between the coefficients of each qualitative response and those of the  $L$ -th (base) alternative for the individual characteristics, making the interpretation of the multinomial probit coefficients less straightforward.

In multinomial models, only one coefficient per alternative-specific attribute needs to be estimated across alternatives. However, for individual characteristics (non-alternative-specific regressors),  $L-1$  coefficients must be estimated, since the coefficient for the base alternative is set to zero for identification.

The likelihood function of this model is

$$p(\boldsymbol{\beta}, \boldsymbol{\Omega} | \mathbf{y}, \mathbf{R}) = \prod_{i=1}^N \prod_{l=1}^L p_{il}^{y_{il}},$$

where  $p_{il} = P(y_{il}^* \geq \max(\mathbf{y}_i^*))$ .

We assume independent priors for the parameters:

$$\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0) \quad \text{and} \quad \boldsymbol{\Omega}^{-1} \sim W(\alpha_0, \boldsymbol{\Sigma}_0),^6$$

and employ Gibbs sampling to estimate the model. With data augmentation for  $\mathbf{w}$ , this model becomes a standard Bayesian linear regression. The conditional posterior distributions are

$$\begin{aligned} \boldsymbol{\beta} | \boldsymbol{\Omega}, \mathbf{w} &\sim N(\boldsymbol{\beta}_n, \mathbf{B}_n), \\ \boldsymbol{\Omega}^{-1} | \boldsymbol{\beta}, \mathbf{w} &\sim W(\alpha_n, \boldsymbol{\Sigma}_n), \end{aligned}$$

where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \mathbf{X}^{*\top} \mathbf{X}^*)^{-1}$ ,  $\boldsymbol{\beta}_n = \mathbf{B}_n(\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{X}^{*\top} \mathbf{w}^*)$ ,  $\boldsymbol{\Omega}^{-1} = \mathbf{C}^\top \mathbf{C}$ ,  $\mathbf{X}_i^{*\top} = \mathbf{C}^\top \mathbf{R}_i$ ,  $\mathbf{w}_i^* = \mathbf{C}^\top \mathbf{w}_i$ ,  $\mathbf{X}^* = [\mathbf{X}_1^*, \mathbf{X}_2^*, \dots, \mathbf{X}_N^*]^\top$ ,  $\alpha_n = \alpha_0 + N$ , and  $\boldsymbol{\Sigma}_n = (\boldsymbol{\Sigma}_0 + \sum_{i=1}^N (\mathbf{w}_i - \mathbf{R}_i \boldsymbol{\beta})^\top (\mathbf{w}_i - \mathbf{R}_i \boldsymbol{\beta}))^{-1}$ .

We can collapse the multinomial outcome vector  $\mathbf{y}_i$  into an indicator variable  $d_i = \sum_{l=1}^{L-1} l \times \mathbb{1}(\max(\mathbf{w}_l) = w_{il})$ .<sup>7</sup> Then, the conditional distribution of

<sup>6</sup>  $W$  denotes the Wishart distribution.

<sup>7</sup> The identification issue arises because scaling  $w_{il}$  by a positive constant does not alter the value of  $d_i$ .

$\mathbf{w}_i | \boldsymbol{\beta}, \boldsymbol{\Omega}^{-1}, d_i$  is an  $(L - 1)$ -dimensional Gaussian truncated over the appropriate cone in  $\mathbb{R}^{L-1}$ .

[304] propose sampling from the univariate conditional distributions  $w_{il} | \mathbf{w}_{i,-l}, \boldsymbol{\beta}, \boldsymbol{\Omega}^{-1}, d_i \sim TN_{I_{il}}(m_{il}, \tau_{ll}^2)$ , where

$$I_{il} = \begin{cases} w_{il} > \max(\mathbf{w}_{i,-l}, 0), & d_i = l, \\ w_{il} < \max(\mathbf{w}_{i,-l}, 0), & d_i \neq l, \end{cases}$$

and permuting the columns and rows of  $\boldsymbol{\Omega}^{-1}$  so that the  $l$ -th column and row appear last:

$$\boldsymbol{\Omega}^{-1} = \begin{bmatrix} \boldsymbol{\Omega}_{-l,-l} & \boldsymbol{\omega}_{-l,l} \\ \boldsymbol{\omega}_{l,-l} & \omega_{l,l} \end{bmatrix}^{-1} = \begin{bmatrix} \boldsymbol{\Omega}_{-l,-l}^{-1} + \tau_{ll}^{-2} \mathbf{f}_l \mathbf{f}_l^\top & -\mathbf{f}_l \tau_{ll}^{-2} \\ -\tau_{ll}^{-2} \mathbf{f}_l^\top & \tau_{ll}^{-2} \end{bmatrix},$$

where  $\mathbf{f}_l = \boldsymbol{\Omega}_{-l,-l}^{-1} \boldsymbol{\omega}_{-l,l}$ ,  $\tau_{ll}^2 = \omega_{ll} - \boldsymbol{\omega}_{l,-l} \boldsymbol{\Omega}_{-l,-l}^{-1} \boldsymbol{\omega}_{-l,l}$ , and  $m_{il} = \mathbf{r}_{il}^\top \boldsymbol{\beta} + \mathbf{f}_l^\top (\mathbf{w}_{i,-l} - \mathbf{R}_{i,-l} \boldsymbol{\beta})$ . Here,  $\mathbf{w}_{i,-l}$  is an  $(L - 2)$ -dimensional vector excluding  $w_{il}$ , and  $\mathbf{r}_{il}$  is the  $l$ -th row of  $\mathbf{R}_i$ , for  $l = 1, 2, \dots, L - 1$ .

The identified parameters are obtained by normalizing with respect to one of the diagonal elements, such that  $\boldsymbol{\beta}/\omega_{11}^{0.5}$  and  $\boldsymbol{\Omega}/\omega_{11}$ .<sup>8</sup>

**Warning:** This model illustrates the importance of deciding whether to specify the model in an identified or unidentified parameter space. The posterior draws often exhibit better mixing properties in the unidentified space [305], reducing computational burden. However, it is crucial to recover the identified space in the final stage. Moreover, defining priors in the unidentified space may have unintended effects on the posterior distributions in the identified space [325]. The multinomial probit model presented in this section is specified in the unidentified space following [304], whereas a version formulated in the identified space is proposed by [305].

#### Example: Choice of fishing mode

In this application, we use the dataset *3Fishing.csv* from [65, p. 491]. The dependent variable represents mutually exclusive fishing mode choices (*mode*), where beach equals 1, pier equals 2, private boat equals 3, and chartered boat (the baseline alternative) equals 4. The model is specified as

$$\mathbf{X}_i = \begin{bmatrix} 1 & 0 & 0 & 0 & \text{Income}_i & 0 & 0 & 0 & \text{Price}_{i,1} & \text{Catch rate}_{i,1} \\ 0 & 1 & 0 & 0 & 0 & \text{Income}_i & 0 & 0 & \text{Price}_{i,2} & \text{Catch rate}_{i,2} \\ 0 & 0 & 1 & 0 & 0 & 0 & \text{Income}_i & 0 & \text{Price}_{i,3} & \text{Catch rate}_{i,3} \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & \text{Income}_i & \text{Price}_{i,4} & \text{Catch rate}_{i,4} \end{bmatrix}.$$

In this example, the chartered boat is the base category, and the number of choice alternatives is four. There are two alternative-specific regressors (price and catch rate) and one non-alternative-specific regressor (income). This specification involves the estimation of eight location parameters ( $\boldsymbol{\beta}$ ):

<sup>8</sup>Our GUI is based on the *bayesm* package, which applies this identification restriction when displaying the posterior outcomes.

three intercepts, three for income, one for price, and one for catch rate. This is the order in which the posterior chains are displayed in our GUI. Note that the location coefficients are set to zero for the baseline category. For multinomial models, we strongly recommend using the last category as the baseline.

Posterior inference also provides estimates of a  $3 \times 3$  covariance matrix (number of alternatives minus one). The element (1, 1) is normalized to 1 due to identification restrictions, and elements (2,4), (3,7), and (6,8) are equal due to symmetry.<sup>9</sup> Note that this identification restriction leads to *NaN* values for the element (1, 1) of the covariance matrix in the [170] and [202] diagnostics, and only eight dependence factors are reported for the remaining covariance elements.

Once the GUI is launched (see the beginning of this chapter), users should follow Algorithm A9 to estimate multinomial probit models in the interface (see Chapter 5 for additional details). The GUI internally relies on the *rmpGibbs* command from the *bayesm* package.

---

**Algorithm A9** Multinomial probit models

---

- 1: Select *Univariate Models* on the top panel
  - 2: Select *Multinomial Probit* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select dependent and independent variables using the *Formula builder* table
  - 6: Select the number of the **Base Alternative**
  - 7: Select the **Number of choice categorical alternatives**
  - 8: Select the **Number of alternative specific variables**
  - 9: Select the **Number of Non-alternative specific variables**
  - 10: Click the *Build formula* button to generate the formula in **R** syntax.
  - 11: Set the hyperparameters: mean vector, covariance matrix, scale matrix and degrees of freedom. This step is not necessary as by default our GUI uses non-informative priors
  - 12: Click the *Go!* button
  - 13: Analyze results
  - 14: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

We ran 100,000 MCMC iterations with an additional 10,000 as burn-in and a thinning parameter of 5, using the default hyperparameter values for

---

<sup>9</sup>This order corresponds to the layout in the *pdf*, *eps*, and *csv* files that can be downloaded from our GUI.

all priors in our GUI. The 95% credible intervals for the coefficients associated with income in the beach and private boat alternatives are (8.58e-06, 8.88e-05) and (3.36e-05, 1.45e-04), respectively. These results suggest that the probability of choosing these alternatives increases relative to the chartered boat option as income rises. Furthermore, an increase in price or a decrease in the catch rate for specific fishing alternatives leads to lower probabilities of selecting those options, as indicated by the corresponding 95% credible intervals of (-9.91e-03, -3.83e-03) and (1.40e-01, 4.62e-01), respectively. However, diagnostic checks of the posterior chains indicate some convergence issues with the posterior draws (see Exercise 5).

## 6.5 The multinomial logit model

The multinomial logit model is used to analyze mutually exclusive discrete outcomes or qualitative response variables. However, this model assumes the independence of irrelevant alternatives (IIA), meaning that the relative odds of choosing between two alternatives are unaffected by the presence or characteristics of a third option; this assumption is not required in the multinomial probit. We consider the multinomial mixed logit model (not to be confused with the random-parameters logit model), which simultaneously accommodates both alternative-varying regressors (conditional) and alternative-invariant regressors (multinomial).<sup>10</sup>

In this setting, there are  $L$  mutually exclusive alternatives, and the dependent variable  $y_{il}$  equals 1 if individual  $i$  chooses the  $l$ -th alternative and 0 otherwise, where  $l = 1, 2, \dots, L$ . The likelihood function is

$$p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \prod_{i=1}^N \prod_{l=1}^L p_{il}^{y_{il}},$$

where the probability that individual  $i$  chooses alternative  $l$  is given by

$$p_{il} = p(y_i = l | \boldsymbol{\beta}, \mathbf{X}) = \frac{\exp \{ \mathbf{x}_{il}^\top \boldsymbol{\beta}_l \}}{\sum_{j=1}^L \exp \{ \mathbf{x}_{ij}^\top \boldsymbol{\beta}_j \}}.$$

Here,  $\mathbf{y}$  and  $\mathbf{X}$  denote the vector and matrix of dependent and explanatory variables, respectively, and  $\boldsymbol{\beta}$  is the vector containing all model coefficients.

Note that coefficients associated with alternative-invariant regressors are normalized to zero for the baseline category, while the coefficients on alternative-specific regressors are constrained to be identical across all categories. We assume a prior distribution  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$ . The posterior distribution is then

$$\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) \propto p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) \times \pi(\boldsymbol{\beta}).$$

<sup>10</sup>The multinomial mixed logit model can be implemented as a conditional logit model.

Because the multinomial logit model does not yield a standard posterior distribution, [379] propose a tailored independent Metropolis–Hastings algorithm in which the proposal distribution follows a multivariate Student- $t$  distribution with  $v$  degrees of freedom (a tuning parameter), mean equal to the maximum likelihood estimate, and scale matrix equal to the inverse of the observed Hessian.

#### **Example: Simulation exercise**

Let's conduct a simulation exercise to evaluate the performance of the Metropolis–Hastings algorithm for inference in the multinomial logit model. We consider a setting with three alternatives, one alternative-invariant regressor (in addition to the intercept), and three alternative-varying regressors. The population parameters are specified as  $\beta_1 = [1 \ -2.5 \ 0.5 \ 0.8 \ -3]^\top$ ,  $\beta_2 = [1 \ -3.5 \ 0.5 \ 0.8 \ -3]^\top$ , and  $\beta_3 = [0 \ 0 \ 0.5 \ 0.8 \ -3]^\top$ , where the first two elements of each vector correspond to the intercept and the alternative-invariant regressor, while the last three elements correspond to the alternative-varying regressors. The sample size is 1,000, and all regressors are simulated from standard normal distributions.

We can launch our GUI using the command line introduced at the beginning of this chapter. To estimate the multinomial logit model, we follow Algorithm A10 in our GUI (see Chapter 5 for further details).

The following **R** code demonstrates how to implement the Metropolis–Hastings (M–H) algorithm from scratch. The first part simulates the dataset, the second defines the log-likelihood function, and the third implements the M–H sampling procedure. We use vague priors centered at zero, with a covariance matrix of  $1000\mathbf{I}_7$ . The posterior estimates closely approximate the true population parameters, and all 95% credible intervals include the corresponding population values.

**Algorithm A10** Multinomial logit models

- 1: Select *Univariate Models* on the top panel
- 2: Select *Multinomial Logit* model using the left radio button
- 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
- 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
- 5: Select dependent and independent variables using the *Formula builder* table
- 6: Select the **Base Alternative**
- 7: Select the **Number of choice categorical alternatives**
- 8: Select the **Number of alternative specific variables**
- 9: Select the **Number of Non-alternative specific variables**
- 10: Click the *Build formula* button to generate the formula in **R** syntax.
- 11: Set the hyperparameters: mean vector and covariance matrix. This step is not necessary as by default our GUI uses non-informative priors
- 12: Select the tuning parameter for the Metropolis-Hastings algorithm, that is, the **Degrees of freedom: Multivariate Student's t distribution**
- 13: Click the *Go!* button
- 14: Analyze results
- 15: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons

***R. code. Simulation of the multinomial logit model***

```

1 remove(list = ls()); set.seed(12345)
2 # Simulation of data
3 N<-1000 # Sample Size
4 B<-c(0.5,0.8,-3); B1<-c(-2.5,-3.5,0); B2<-c(1,1,0)
5 # Alternative specific attributes of choice 1
6 X1<-matrix(cbind(rnorm(N,0,1),rnorm(N,0,1),rnorm(N,0,1)),N,
  length(B))
7 # Alternative specific attributes of choice 2
8 X2<-matrix(cbind(rnorm(N,0,1),rnorm(N,0,1),rnorm(N,0,1)),N,
  length(B))
9 # Alternative specific attributes of choice 3
10 X3<-matrix(cbind(rnorm(N,0,1),rnorm(N,0,1),rnorm(N,0,1)),N,
   length(B))
11 X4<-matrix(rnorm(N,1,1),N,1)
12 V1<-B2[1]+X1%*%B+B1[1]*X4; V2<-B2[2]+X2%*%B+B1[2]*X4; V3<-B2
  [3]+X3%*%B+B1[3]*X4

```

### R. code. Simulation of the multinomial logit model

```

1 suma<-exp(V1)+exp(V2)+exp(V3)
2 p1<-exp(V1)/suma; p2<-exp(V2)/suma; p3<-exp(V3)/suma
3 p<-cbind(p1,p2,p3)
4 y<- apply(p,1, function(x)sample(1:3, 1, prob = x, replace =
TRUE))
5 y1<-y==1; y2<-y==2; y3<-y==3
6 # Log likelihood
7 log.L<- function(Beta){
8   V1<-Beta[1]+Beta[3]*X4+X1%*%Beta[5:7]
9   V2<-Beta[2]+Beta[4]*X4+X2%*%Beta[5:7]
10  V3<- X3%*%Beta[5:7]
11  suma<-exp(V1)+exp(V2)+exp(V3)
12  p11<-exp(V1)/suma; p22<-exp(V2)/suma; p33<-exp(V3)/suma
13  suma2<-NULL
14  for(i in 1:N){
15    suma1<-y1[i]*log(p11[i])+y2[i]*log(p22[i])+y3[i]*log(p33
     [i])
16    suma2<-c(suma2,suma1)}
17  logL<-sum(suma2)
18  return(-logL)
19 }
20 # Parameters: Proposal
21 k <- 7
22 res.optim<-optim(rep(0, k), log.L, method="BFGS", hessian=
TRUE)
23 MeanT <- res.optim$par
24 ScaleT <- as.matrix(Matrix::forceSymmetric(solve(res.optim$ 
  hessian))) # Force this matrix to be symmetric
25 # Hyperparameters: Priors
26 B0 <- 1000*diag(k); b0 <- rep(0, k)
27 MHfunction <- function(iter, tuning){
28   Beta <- rep(0, k); Acept <- NULL
29   BetasPost <- matrix(NA, iter, k)
30   pb <- txtProgressBar(min = 0, max = iter, style = 3)
31   for(s in 1:iter){
32     LogPostBeta <- -log.L(Beta) + mvtnorm::dmvnorm(Beta,
      mean = b0, sigma = B0, log = TRUE)
33     BetaC <- c(LaplaceDemon::rmvt(n=1, mu = MeanT, S =
      ScaleT, df = tuning))
34     LogPostBetaC <- -log.L(BetaC) + mvtnorm::dmvnorm(BetaC,
      mean = b0, sigma = B0, log = TRUE)

```

### R. code. Simulation of the multinomial logit model

```

1   alpha <- min(exp((LogPostBetaC-mvtnorm::dmvt(BetaC,
2   delta = MeanT, sigma = ScaleT, df = tuning, log = TRUE))
3   -(LogPostBeta-mvtnorm::dmvt(Beta, delta = MeanT, sigma =
4   ScaleT, df = tuning, log = TRUE))), 1)
5   u <- runif(1)
6   if(u <= alpha){
7     Acepti <- 1; Beta <- BetaC
8   }else{
9     Acepti <- 0; Beta <- Beta
10  }
11  BetasPost[s, ] <- Beta; Acept <- c(Acept, Acepti)
12  setTxtProgressBar(pb, s)
13  }
14  close(pb); AcepRate <- mean(Acept)
15  Results <- list(AcepRate = AcepRate, BetasPost = BetasPost
16  )
17  return(Results)
18 }

# MCMC parameters
mcmc <- 10000; burnin <- 1000; thin <- 5; iter <- mcmc +
burnin; keep <- seq(burnin, iter, thin); tuning <- 6 # Degrees of freedom
ResultsPost <- MHfunction(iter = iter, tuning = tuning)
summary(coda::mcmc(ResultsPost$BetasPost[keep[-1], ]))
```

## 6.6 Ordered probit model

The ordered probit model is used when the categorical response variable has a natural ordering. In this case, there exists an unobserved latent variable  $y_i^* = \mathbf{x}_i^\top \boldsymbol{\beta} + \mu_i$ , where  $\mathbf{x}_i$  is a  $K$ -dimensional vector of regressors and  $\mu_i \stackrel{i.i.d.}{\sim} N(0, 1)$ . The observed outcome is defined as  $y_i = l$  if and only if  $\alpha_{l-1} < y_i^* \leq \alpha_l$ , for  $l \in \{1, 2, \dots, L\}$ , where  $\alpha_0 = -\infty$ ,  $\alpha_1 = 0$ , and  $\alpha_L = \infty$ .<sup>11</sup>

The probability of observing category  $l$  is given by

$$p(y_i = l) = \Phi(\alpha_l - \mathbf{x}_i^\top \boldsymbol{\beta}) - \Phi(\alpha_{l-1} - \mathbf{x}_i^\top \boldsymbol{\beta}),$$

---

<sup>11</sup>Identification requires fixing the variance of the latent error term to 1 and setting  $\alpha_1 = 0$ . Note that multiplying  $y_i^*$  by a positive constant, or adding a constant to all the cut-offs and subtracting the same constant from the intercept, does not affect the observed outcome  $y_i$ .

and the likelihood function is

$$p(\boldsymbol{\beta}, \boldsymbol{\alpha} | \mathbf{y}, \mathbf{X}) = \prod_{i=1}^N p(y_i = l | \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{X}).$$

The priors are assumed to be independent, i.e.,  $\pi(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \pi(\boldsymbol{\beta}) \times \pi(\boldsymbol{\gamma})$ , where  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$  and  $\boldsymbol{\gamma} \sim N(\boldsymbol{\gamma}_0, \boldsymbol{\Gamma}_0)$ , with  $\boldsymbol{\gamma} = [\gamma_2 \ \gamma_3 \ \dots \ \gamma_{L-1}]^\top$  such that

$$\boldsymbol{\alpha} = \begin{bmatrix} \exp\{\gamma_2\} \\ \sum_{l=2}^3 \exp\{\gamma_l\} \\ \vdots \\ \sum_{l=2}^{L-1} \exp\{\gamma_l\} \end{bmatrix}.$$

This parameterization enforces the ordering constraint on the cut-offs.

Although the model does not yield a standard conditional posterior distribution for  $\boldsymbol{\gamma}$  (and therefore for  $\boldsymbol{\alpha}$ ), it admits a standard conditional posterior for  $\boldsymbol{\beta}$  once data augmentation is introduced. Thus, inference can be performed using a Metropolis-within-Gibbs sampling algorithm. In particular, we use Gibbs sampling to draw  $\boldsymbol{\beta}$  and  $\mathbf{y}^*$ , where

$$\boldsymbol{\beta} | \mathbf{y}^*, \boldsymbol{\alpha}, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \mathbf{B}_n),$$

with  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}$  and  $\boldsymbol{\beta}_n = \mathbf{B}_n(\mathbf{B}_0^{-1}\boldsymbol{\beta}_0 + \mathbf{X}^\top \mathbf{y}^*)$ , and

$$y_i^* | \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{y}, \mathbf{X} \sim TN_{(\alpha_{y_i-1}, \alpha_{y_i})}(\mathbf{x}_i^\top \boldsymbol{\beta}, 1).$$

For  $\boldsymbol{\gamma}$ , we employ a random-walk Metropolis–Hastings algorithm with a Gaussian proposal distribution centered at the current value and covariance matrix  $s^2(\boldsymbol{\Gamma}_0^{-1} + \hat{\boldsymbol{\Sigma}}_\gamma^{-1})^{-1}$ , where  $s > 0$  is a tuning parameter and  $\hat{\boldsymbol{\Sigma}}_\gamma$  denotes the sample covariance matrix of  $\boldsymbol{\gamma}$  obtained from the maximum likelihood estimation.

### Example: Determinants of preventive health care visits

We use the file *2HealthMed.csv* for this application. The dependent variable is *MedVisPrevOr*, an ordered variable equal to 1 if the individual did not visit a physician for preventive reasons, 2 if the individual visited once during the year, and so on, up to 6 if the individual visited five or more times. The latter category represents 1.6% of the sample. Thus, the dependent variable has six ordered categories.

The set of regressors includes *SHI*, an indicator variable equal to 1 if the individual is enrolled in the subsidized health care system; *Female*, a gender indicator; *Age* and its square; socioeconomic status indicators (*Est2* and *Est3*), with the lowest level as the reference category; self-perceived health status (*Fair*, *Good*, and *Excellent*), where *Bad* serves as the baseline; and education level (*PriEd*, *HighEd*, *VocEd*, *UnivEd*), with *No education* as the reference category.

We ran this application with 50,000 MCMC iterations, an additional 10,000 as burn-in, and a thinning parameter of 5, resulting in 10,000 effective posterior draws. The prior parameters were set as  $\beta_0 = \mathbf{0}_{11}$ ,  $B_0 = 1000\mathbf{I}_{11}$ ,  $\gamma_0 = \mathbf{0}_4$ ,  $\Gamma_0 = \mathbf{I}_4$ , and the tuning parameter was set to 1.

The ordered probit model can be estimated in our GUI by following the steps outlined in Algorithm A11.

---

**Algorithm A11** Ordered probit models

---

- 1: Select *Univariate Models* on the top panel
  - 2: Select *Ordered Probit* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select dependent and independent variables using the *Formula builder* table
  - 6: Click the *Build formula* button to generate the formula in **R** syntax. Remember that this formula must have -1 to omit the intercept in the specification.
  - 7: Set the hyperparameters: mean vector and covariance matrix. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Select the number of ordered alternatives
  - 9: Set the hyperparameters: mean and covariance matrix of the cutoffs. This step is not necessary as by default our GUI uses non-informative prior
  - 10: Select the tuning parameter for the Metropolis-Hastings algorithm
  - 11: Click the *Go!* button
  - 12: Analyze results
  - 13: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

The following **R** code shows how to perform inference in this model using the command *rordprobitGibbs* from the *bayesm* library, which is the command that our GUI uses.

### R. code. Determinants of preventive health care visits

```

1 rm(list = ls()); set.seed(10101)
2 Data <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/2HealthMed.
  csv", sep = ",", header = TRUE, quote = "")
3 attach(Data)
4 y <- MedVisPrevOr
5 # MedVisPrevOr: Ordered variable for preventive visits to
  # doctors in one year: 1 (none), 2 (once), ... 6 (five or
  # more)
6 X <- cbind(SHI, Female, Age, Age2, Est2, Est3, Fair, Good,
  Excellent, PriEd, HighEd, VocEd, UnivEd)
7 k <- dim(X)[2]
8 L <- length(table(y))
9 # Hyperparameters
10 b0 <- rep(0, k); c0 <- 1000; B0 <- c0*diag(k)
11 gamma0 <- rep(0, L-2); Gamma0 <- diag(L-2)
12 # MCMC parameters
13 mcmc <- 60000+1; thin <- 5; tuningPar <- 1/(L-2)^0.5
14 DataApp <- list(y = y, X = X, k = L)
15 Prior <- list(betabar = b0, A = solve(B0), dstarbar = gamma0
  , Ad = Gamma0)
16 mcmcpar <- list(R = mcmc, keep = 5, s = tuningPar)
17 PostBeta <- bayesm::rordprobitGibbs(Data = DataApp, Prior =
  Prior, Mcmc = mcmcpar)
18 BetasPost <- coda::mcmc(PostBeta[["betadraw"]])
19 colnames(BetasPost) <- c("SHI", "Female", "Age", "Age2", " "
  , "Est2", "Est3", "Fair", "Good", "Excellent", "PriEd", " "
  , "HighEd", "VocEd", "UnivEd")
20 summary(BetasPost)
21 # Convergence diagnostics
22 coda::geweke.diag(BetasPost)
23 coda::raftery.diag(BetasPost, q=0.5, r=0.05, s = 0.95)
24 coda::heidel.diag(BetasPost)
25 # Cut offs
26 Cutoffs <- PostBeta[["cutdraw"]]
27 summary(Cutoffs)
28 coda::geweke.diag(Cutoffs)
29 coda::heidel.diag(Cutoffs)
30 coda::raftery.diag(Cutoffs[,-1], q=0.5, r=0.05, s = 0.95)

```

The results suggest that older individuals (at a decreasing rate), those enrolled in the subsidized health care program, belonging to the second socioeconomic stratum, with a better self-perception of health, and without a high school degree as their highest level of education, have a higher probability of visiting a physician for preventive health purposes. Convergence diagnostics

appear satisfactory for most parameters, except for the draws associated with self-perceived health status.

We also obtained the posterior estimates of the cut-offs in the ordered probit model. These estimates are necessary to compute the probability that an individual falls into a specific category of physician visits. Due to identification restrictions, the first cut-off is fixed at zero. Consequently, *Nan* values appear in the [170] and [202] diagnostics, and only four values are reported in the [352] test, corresponding to the remaining unrestricted cut-offs. However, the [352] diagnostics indicate potential convergence issues for these parameters, as their dependence factors are relatively high.

## 6.7 Negative binomial model

The dependent variable in the negative binomial model is a nonnegative integer (count). Unlike the Poisson model, the negative binomial model accounts for over-dispersion, whereas the Poisson model assumes equi-dispersion—that is, the mean and variance are equal.

We assume that  $y_i \stackrel{i.i.d.}{\sim} \text{NB}(\gamma, \theta_i)$ , where the probability mass function for individual  $i$  is

$$\frac{\Gamma(y_i + \gamma)}{\Gamma(\gamma)y_i!}(1 - \theta_i)^{y_i}\theta_i^\gamma,$$

with success probability  $\theta_i = \frac{\gamma}{\lambda_i + \gamma}$ , where  $\lambda_i = \exp\{\boldsymbol{x}_i^\top \boldsymbol{\beta}\}$  is the conditional mean,  $\gamma = \exp\{\alpha\}$  is the target number of successful trials (or dispersion parameter), and  $\boldsymbol{x}_i$  is a  $K$ -dimensional vector of regressors.

We assume independent priors for the model parameters:  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \boldsymbol{B}_0)$  and  $\alpha \sim G(\alpha_0, \delta_0)$ .

Since this model does not yield standard conditional posterior distributions, [379] propose using a random-walk Metropolis–Hastings algorithm. The proposal distribution for  $\boldsymbol{\beta}$  is Gaussian, centered at the current draw, with covariance matrix  $s_{\boldsymbol{\beta}}^2 \hat{\Sigma}_{\boldsymbol{\beta}}$ , where  $s_{\boldsymbol{\beta}}$  is a tuning parameter and  $\hat{\Sigma}_{\boldsymbol{\beta}}$  is the covariance matrix from the maximum likelihood estimation. Similarly, the proposal distribution for  $\alpha$  is normal, centered at the current value, with variance  $s_\alpha^2 \hat{\sigma}_\alpha^2$ , where  $s_\alpha$  is a tuning parameter and  $\hat{\sigma}_\alpha^2$  is the variance from the maximum likelihood estimation.

### Example: Simulation exercise

Let's conduct a simulation exercise to evaluate the performance of the Metropolis–Hastings algorithm in the negative binomial model. There are two regressors,  $x_{i1} \sim U(0, 1)$  and  $x_{i2} \sim N(0, 1)$ , along with an intercept. The dispersion parameter is set to  $\gamma = \exp\{1.2\}$ , and the coefficient vector is  $\boldsymbol{\beta} = [1 \ 1 \ 1]^\top$ . The sample size is 1,000.

We run this simulation using 10,000 MCMC iterations, a burn-in period of 1,000, and a thinning parameter of 5. Vague priors are assigned to the location parameters:  $\beta_0 = \mathbf{0}_3$ ,  $B_0 = 1000\mathbf{I}_3$ ,  $\alpha_0 = 0.5$ , and  $\delta_0 = 0.1$ , which are the default values in the *rnegbinRw* command from the *bayesm* package in **R**. In addition, the tuning parameters for the Metropolis–Hastings algorithm are  $s_\beta = 2.93/k^{1/2}$  and  $s_\alpha = 2.93$ , which are also the default values in *rnegbinRw*, where  $k$  denotes the number of location parameters.

The negative binomial model can be estimated in our GUI by following the steps outlined in Algorithm A12.

---

**Algorithm A12** Negative binomial models

---

- 1: Select *Univariate Models* on the top panel
  - 2: Select *Negative Binomial (Poisson)* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select dependent and independent variables using the *Formula builder* table
  - 6: Click the *Build formula* button to generate the formula in **R** syntax. You can modify the formula in the **Main equation** box using valid arguments of the *formula* command structure in **R**
  - 7: Set the hyperparameters: mean vector, covariance matrix, shape and scale parameters. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Select the tuning parameters for the Metropolis-Hastings algorithms
  - 9: Click the *Go!* button
  - 10: Analyze results
  - 11: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

The following **R** code demonstrates how to perform inference in the negative binomial model by programming the Metropolis–Hastings algorithm from scratch. In Exercise 8, we ask you to estimate this example using the *rnegbinRw* command.

The results show that all 95% credible intervals contain the true population parameters, and the posterior means are very close to those values.

*R. code. Simulation of the negative binomial model*

```

1 rm(list = ls()); set.seed(10101)
2 N <- 2000 # Sample size
3 x1 <- runif(N); x2 <- rnorm(N)
4 X <- cbind(1, x1, x2); k <- dim(X)[2]; B <- rep(1, k)
5 alpha <- 1.2; gamma <- exp(alpha); lambda <- exp(X%*%B)
6 y <- rnbnom(N, mu = lambda, size = gamma)
7 # log likelihood
8 logLik <- function(par){
9   alpha <- par[1]; beta <- par[2:(k+1)]
10  gamma <- exp(alpha)
11  lambda <- exp(X%*%beta)
12  logLikNB <- sum(sapply(1:N, function(i){dnbinom(y[i], size
13    = gamma, mu = lambda[i], log = TRUE)}))
14  return(-logLikNB)
15 }
16 # Parameters: Proposal
17 par0 <- rep(0.5, k+1)
18 res.optim <- optim(par0, logLik, method="BFGS", hessian=TRUE
19 )
20 res.optim$par
21 res.optim$convergence
22 Covar <- solve(res.optim$hessian)
23 CovarBetas <- Covar[2:(k+1),2:(k+1)]
24 VarAlpha <- Covar[1:1]
25 # Hyperparameters: Priors
26 B0 <- 1000*diag(k); b0 <- rep(0, k)
27 alpha0 <- 0.5; delta0 <- 0.1
28 # Metropolis-Hastings function
29 MHfunction <- function(iter, sbeta, salpha){
30   Beta <- rep(0, k); Acept1 <- NULL; Acept2 <- NULL
31   BetasPost <- matrix(NA, iter, k); alpha <- 1
32   alphaPost <- rep(NA, iter); par <- c(alpha, Beta)
33   pb <- txtProgressBar(min = 0, max = iter, style = 3)
34   for(s in 1:iter){
35     LogPostBeta <- -logLik(par) + dgamma(alpha, shape =
36       alpha0, scale = delta0, log = TRUE) + mvtnorm::dmvnorm(
37       Beta, mean = b0, sigma = B0, log = TRUE)
38     BetaC <- c(MASS::mvrnorm(1, mu = Beta, Sigma = sbeta^2*
39       CovarBetas))
40     parC <- c(alpha, BetaC)
41     LogPostBetaC <- -logLik(parC) + dgamma(alpha, shape =
42       alpha0, scale = delta0, log = TRUE) + mvtnorm::dmvnorm(
43       BetaC, mean = b0, sigma = B0, log = TRUE)
44     alpha1 <- min(exp((LogPostBetaC - mvtnorm::dmvnorm(BetaC
45       , mean = Beta, sigma = sbeta^2*CovarBetas, log = TRUE))
46       -(LogPostBeta - mvtnorm::dmvnorm(Beta, mean = Beta,
47       sigma = sbeta^2*CovarBetas, log = TRUE))),1)

```

**R. code. Simulation of the negative binomial model, M-H algorithm**

```

1   u1 <- runif(1)
2   if(u1 <= alpha1){Acept1i <- 1; Beta <- BetaC}else{
3     Acept1i <- 0; Beta <- Beta
4   }
5   par <- c(alpha, Beta)
6   LogPostBeta <- -logLik(par) + dgamma(alpha, shape =
alpha0, scale = delta0, log = TRUE) + mvtnorm::dmvnorm(
Beta, mean = b0, sigma = B0, log = TRUE)
7   alphaC <- rnorm(1, mean = alpha, sd = salpha*VarAlpha
^0.5)
8   parC <- c(alphaC, Beta)
9   LogPostBetaC <- -logLik(parC) + dgamma(alphaC, shape =
alpha0, scale = delta0, log = TRUE) + mvtnorm::dmvnorm(
Beta, mean = b0, sigma = B0, log = TRUE)
10  alpha2 <- min(exp((LogPostBetaC - dnorm(alphaC, mean =
alpha, sd = salpha*VarAlpha^0.5, log = TRUE))-
LogPostBeta - dnorm(alpha, mean = alpha, sd = salpha*
VarAlpha^0.5, log = TRUE))),1)
11  u2 <- runif(1)
12  if(u2 <= alpha2){Acept2i <- 1; alpha <- alphaC}else{
13    Acept2i <- 0; alpha <- alpha
14  }
15
16  BetasPost[s, ] <- Beta; alphaPost[s] <- alpha
17  Acept1 <- c(Acept1, Acept1i); Acept2 <- c(Acept2,
Acept2i)
18  setTxtProgressBar(pb, s)
19 }
20 close(pb)
21 AcepRateBeta <- mean(Acept1); AcepRateAlpha <- mean(Acept2
)
22 Results <- list(AcepRateBeta = AcepRateBeta, AcepRateAlpha
= AcepRateAlpha, BetasPost = BetasPost, alphaPost =
alphaPost)
23 return(Results)
24 }
25 # MCMC parameters
26 mcmc <- 10000; burnin <- 1000; thin <- 5
27 iter <- mcmc + burnin; keep <- seq(burnin, iter, thin)
28 sbeta <- 2.93/sqrt(k); salpha <- 2.93
29 # Run M-H
30 ResultsPost <- MHfunction(iter = iter, sbeta = sbeta, salpha
= salpha)
31 ResultsPost$AcepRateBeta; ResultsPost$AcepRateAlpha
32 summary(coda::mcmc(ResultsPost$BetasPost[keep[-1], ]))
33 summary(coda::mcmc(ResultsPost$alphaPost[keep[-1]]))

```

## 6.8 Tobit model

The dependent variable in Tobit models is partially observed due to sampling or censoring, while the regressors are fully observed. In particular,

$$y_i = \begin{cases} L, & \text{if } y_i^* < L, \\ y_i^*, & \text{if } L \leq y_i^* < U, \\ U, & \text{if } y_i^* \geq U, \end{cases}$$

where  $y_i^* \stackrel{i.i.d.}{\sim} N(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2)$ , and  $\mathbf{x}_i$  is a  $K$ -dimensional vector of regressors.<sup>12</sup>

We assume independent conjugate priors:  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$  and  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ . Using data augmentation, we introduce  $\mathbf{y}_C^*$  such that  $y_{C_i}^* \stackrel{i.i.d.}{\sim} N(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2)$ , where  $y_{C_i} = \{y_{C_i}^L, y_{C_i}^U\}$  corresponds to the lower- and upper-censored observations. This formulation allows us to implement a Gibbs sampling algorithm [86].

The joint posterior distribution is

$$\begin{aligned} \pi(\boldsymbol{\beta}, \sigma^2, \mathbf{y}^* | \mathbf{y}, \mathbf{X}) \propto & \prod_{i=1}^N \left[ \mathbb{1}(y_i = L) \mathbb{1}(y_{C_i}^L < L) + \mathbb{1}(L \leq y_i < U) \right. \\ & \left. + \mathbb{1}(y_i = U) \mathbb{1}(y_{C_i}^U \geq U) \right] \times N(y_i^* | \mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2) \\ & \times N(\boldsymbol{\beta} | \boldsymbol{\beta}_0, \mathbf{B}_0) \times IG(\sigma^2 | \alpha_0/2, \delta_0/2). \end{aligned}$$

The conditional posterior distributions are:

$$\begin{aligned} y_{C_i}^* | \boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X} &\sim \begin{cases} TN_{(-\infty, L)}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2), & y_i = L, \\ TN_{[U, \infty)}(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2), & y_i = U, \end{cases} \\ \boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} &\sim N(\boldsymbol{\beta}_n, \sigma^2 \mathbf{B}_n), \\ \sigma^2 | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X} &\sim IG(\alpha_n/2, \delta_n/2), \end{aligned}$$

where

$$\mathbf{B}_n = (\mathbf{B}_0^{-1} + \sigma^{-2} \mathbf{X}^\top \mathbf{X})^{-1}, \quad \boldsymbol{\beta}_n = \mathbf{B}_n (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sigma^{-2} \mathbf{X}^\top \mathbf{y}^*), \quad \alpha_n = \alpha_0 + N,$$

and  $\delta_n = \delta_0 + (\mathbf{y}^* - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y}^* - \mathbf{X}\boldsymbol{\beta})$ .

### Example: The market value of soccer players in Europe continues

We continue with the example of the market value of soccer players from Section 6.1. We specify the same model but now assume that the sample is left-censored, meaning we only observe soccer players whose market value

<sup>12</sup>We can set  $L$  or  $U$  equal to  $-\infty$  or  $\infty$  to model data censored on only one side.

exceeds one million euros. The dependent variable is  $\log(ValueCens)$ , and the left-censoring point is 13.82.

Algorithm A13 illustrates how to estimate Tobit models in our GUI. Our GUI implements the *MCMCTobit* command from the *MCMCpack* package.

---

**Algorithm A13** Tobit models

---

- 1: Select *Univariate Models* on the top panel
  - 2: Select *Tobit* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select dependent and independent variables using the *Formula builder* table
  - 6: Click the *Build formula* button to generate the formula in **R** syntax. You can modify the formula in the **Main equation** box using valid arguments of the *formula* command structure in **R**
  - 7: Set the left and right censoring points. To censor above only, specify *-Inf* in the left censoring box, and to censor below only, specify *Inf* in the right censoring box
  - 8: Set the hyperparameters: mean vector, covariance matrix, shape and scale parameters. This step is not necessary as by default our GUI uses non-informative priors
  - 9: Click the *Go!* button
  - 10: Analyze results
  - 11: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

We run this application using the same hyperparameters that we set in the example of Section 6.1. All results seem similar to those in the example of linear models. In addition, the posterior chains seem to achieve good diagnostics.

*R. code. The value of soccer player with left censoring*

```

1 rm(list = ls()); set.seed(10101)
2 Data <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/1
  ValueFootballPlayers.csv", sep = ",", header = TRUE,
  quote = "")
3 attach(Data)
4 y <- log(ValueCens)
5 X <- cbind(1, Perf, Age, Age2, NatTeam, Goals, Exp, Exp2)
6 k <- dim(X)[2]
7 N <- dim(X)[1]
8 # Hyperparameters
9 d0 <- 0.001; a0 <- 0.001
10 b0 <- rep(0, k); c0 <- 1000; B0 <- c0*diag(k)
11 B0i <- solve(B0)
12 # MCMC parameters
13 mcmc <- 50000; burnin <- 10000
14 tot <- mcmc + burnin; thin <- 1
15 # Posterior distributions using packages: MCMCpack sets the
  # model in terms of the precision matrix
16 posterior <- MCMCpack::MCMCTobit(y~X-1, b0=b0, B0 = B0i, c0
  = a0, d0 = d0, burnin = burnin, mcmc = mcmc, thin =
  thin, below = 13.82, above = Inf)
17 summary(coda::mcmc(posterior))
18 # Gibbs sampling functions
19 XtX <- t(X) %*% X
20 PostBeta <- function(Yl, sig2){
21   Bn <- solve(B0i + sig2^(-1)*XtX)
22   bn <- Bn %*% (B0i %*% b0 + sig2^(-1)*t(X) %*% Yl)
23   Beta <- MASS::mvrnorm(1, bn, Bn)
24   return(Beta)
25 }
26 PostYl <- function(Beta, sig2, L, U, i){
27   Ylmean <- X[i,] %*% Beta
28   if(y[i] == L){
29     Yli <- truncnorm::rtruncnorm(1, a = -Inf, b = L, mean =
      Ylmean, sd = sig2^0.5)
30   }else{
31     if(y[i] == U){
32       Yli <- truncnorm::rtruncnorm(1, a = U, b = Inf, mean =
      Ylmean, sd = sig2^0.5)
33     }else{
34       Yli <- y[i]
35     }
36   }
37   return(Yli)
38 }
```

**R. code. The value of soccer player with left censoring, Gibbs sampler**

```

1 PostSig2 <- function(Beta, Yl){
2   an <- a0 + length(y)
3   dn <- d0 + t(Yl - X%*%Beta)%*%(Yl - X%*%Beta)
4   sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
5   return(sig2)
6 }
7 PostBetas <- matrix(0, mcmc+burnin, k); Beta <- rep(0, k)
8 PostSigma2 <- rep(0, mcmc+burnin); sig2 <- 1
9 L <- log(1000000); U <- Inf
10 # create progress bar in case that you want to see
    iterations progress
11 pb <- txtProgressBar(min = 0, max = tot, style = 3)
12 for(s in 1:tot){
13   Yl <- sapply(1:N, function(i){PostYl(Beta = Beta, sig2 =
      sig2, L = L, U = U, i)})
14   Beta <- PostBeta(Yl = Yl, sig2 = sig2)
15   sig2 <- PostSig2(Beta = Beta, Yl = Yl)
16   PostBetas[s,] <- Beta; PostSigma2[s] <- sig2
17   setTxtProgressBar(pb, s)
18 }
19 close(pb)
20 keep <- seq((burnin+1), tot, thin)
21 PosteriorBetas <- PostBetas[keep,]
22 colnames(PosteriorBetas) <- c("Intercept", "Perf", "Age", "
  Age2", "NatTeam", "Goals", "Exp", "Exp2")
23 summary(coda::mcmc(PosteriorBetas))
24 summary(coda::mcmc(PostSigma2[keep]))

```

## 6.9 Quantile regression

In quantile regression, the location parameters vary across quantiles of the dependent variable. Let  $q_\tau(\mathbf{x}_i) = \mathbf{x}_i^\top \boldsymbol{\beta}_\tau$  denote the  $\tau$ -th quantile regression function of  $y_i$  given  $\mathbf{x}_i$ , where  $\mathbf{x}_i$  is a  $K$ -dimensional vector of regressors and  $0 < \tau < 1$ . Specifically, we consider the model  $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_\tau + \mu_i$ , subject to the condition  $\int_{-\infty}^0 f_\tau(\mu_i) d\mu_i = \tau$ , meaning that the  $\tau$ -th quantile of  $\mu_i$  is 0.

[251] propose modeling  $f_\tau(\mu_i)$  using the asymmetric Laplace distribution, given by

$$f_\tau(\mu_i) = \tau(1 - \tau) \exp \{ -\mu_i (\tau - \mathbb{1}(\mu_i < 0)) \},$$

where  $\mu_i(\tau - \mathbb{1}(\mu_i < 0))$  is the check (loss) function. These authors also show that this distribution admits a location-scale mixture of normals representation:

$$\mu_i = \theta e_i + \psi \sqrt{e_i} z_i,$$

where  $\theta = \frac{1-2\tau}{\tau(1-\tau)}$ ,  $\psi^2 = \frac{2}{\tau(1-\tau)}$ ,  $e_i \sim E(1)$ , and  $z_i \sim N(0, 1)$ , with  $e_i \perp z_i$ .<sup>13</sup>

Given this representation and assuming that the sample is i.i.d., the likelihood function becomes

$$p(\mathbf{y} | \boldsymbol{\beta}_\tau, \mathbf{e}, \mathbf{X}) \propto \left( \prod_{i=1}^N e_i^{-1/2} \right) \exp \left\{ - \sum_{i=1}^N \frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_\tau - \theta e_i)^2}{2\psi^2 e_i} \right\}.$$

Assuming a normal prior for  $\boldsymbol{\beta}_\tau$ , i.e.,  $\boldsymbol{\beta}_\tau \sim N(\boldsymbol{\beta}_{\tau 0}, \mathbf{B}_{\tau 0})$ , and applying data augmentation for  $\mathbf{e}$ , a Gibbs sampling algorithm can be implemented for this model. The posterior distributions are:

$$\begin{aligned} \boldsymbol{\beta}_\tau | \mathbf{e}, \mathbf{y}, \mathbf{X} &\sim N(\boldsymbol{\beta}_{n\tau}, \mathbf{B}_{n\tau}), \\ e_i | \boldsymbol{\beta}_\tau, \mathbf{y}, \mathbf{X} &\sim \text{GIG} \left( \frac{1}{2}, \alpha_{ni}, \delta_{ni} \right), \end{aligned} \quad ^{14}$$

where

$$\begin{aligned} \mathbf{B}_{n\tau} &= \left( \mathbf{B}_{\tau 0}^{-1} + \sum_{i=1}^N \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\psi^2 e_i} \right)^{-1}, \quad \boldsymbol{\beta}_{n\tau} = \mathbf{B}_{n\tau} \left( \mathbf{B}_{\tau 0}^{-1} \boldsymbol{\beta}_{\tau 0} + \sum_{i=1}^N \frac{\mathbf{x}_i (y_i - \theta e_i)}{\psi^2 e_i} \right), \\ \alpha_{ni} &= \frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_\tau)^2}{\psi^2}, \quad \delta_{ni} = 2 + \frac{\theta^2}{\psi^2}. \end{aligned}$$

#### Example: The market value of soccer players in Europe continues

We continue with the example of the market value of soccer players from Section 6.1. We now examine whether the marginal effect of having played on the national team varies across the quantiles of the market value distribution of top European soccer players. We use the same set of regressors as in the previous example but focus on the effects at the 0.5 and 0.9 quantiles of *NatTeam*.

Algorithm A14 illustrates how to estimate quantile regression models in our GUI. Our GUI implements the *MCMCquantreg* command from the *MCMCpack* package. The following code demonstrates how to conduct this analysis using that package.

The results show that, at the median of the market value distribution (0.5 quantile), the 95% credible interval for the coefficient associated with *NatTeam* is (0.34, 1.02), with a posterior mean of 0.69. At the 0.9 quantile, these values are (0.44, 1.59) and 1.03, respectively. These results suggest that being on the national team increases the market value of higher-valued players more strongly, although there is some overlap in the credible intervals.

<sup>13</sup>*E* denotes an exponential density.

<sup>14</sup>GIG denotes a generalized inverse Gaussian density.

**Algorithm A14** Quantile regression

- 1: Select *Univariate Models* on the top panel
- 2: Select *Quantile* model using the left radio button
- 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
- 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
- 5: Select dependent and independent variables using the *Formula builder* table
- 6: Click the *Build formula* button to generate the formula in **R** syntax. You can modify the formula in the **Main equation** box using valid arguments of the *formula* command structure in **R**
- 7: Set the quantile to be analyzed, by default is 0.5
- 8: Set the hyperparameters: mean vector and covariance matrix. This step is not necessary as by default our GUI uses non-informative priors
- 9: Click the *Go!* button
- 10: Analyze results
- 11: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons

**R. code. The value of soccer player, quantile regression**

```

1 rm(list = ls()); set.seed(10101)
2 Data <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/1
  ValueFootballPlayers.csv", sep = ",", header = TRUE,
  quote = "")
3 attach(Data)
4 y <- log(ValueCens)
5 X <- cbind(1, Perf, Age, Age2, NatTeam, Goals, Exp, Exp2)
6 k <- dim(X)[2]; N <- dim(X)[1]
7 # Hyperparameters
8 b0 <- rep(0, k); c0 <- 1000; B0<-c0*diag(k); B0i<-solve(B0)
9 # MCMC parameters
10 mcmc <- 50000; burnin <- 10000
11 tot <- mcmc + burnin; thin <- 1
12 q <- 0.5 # quantile
13 posterior05 <- MCMCpack::MCMCquantreg(y~X-1, tau = q, b0=b0
  , B0 = B0i, burnin = burnin, mcmc = mcmc, thin = thin)
14 summary(coda::mcmc(posterior05))
15 q <- 0.9 # quantile
16 posterior09 <- MCMCpack::MCMCquantreg(y~X-1, tau = q, b0=b0
  , B0 = B0i, burnin = burnin, mcmc = mcmc, thin = thin)
17 summary(coda::mcmc(posterior09))

```

## 6.10 Bayesian bootstrap regression

We implement the Bayesian bootstrap [391] for linear regression models. In particular, the Bayesian bootstrap simulates posterior distributions by assuming that the sample cumulative distribution function (CDF) represents the population CDF, an assumption also implicit in the frequentist bootstrap [132].

Let  $y_i \stackrel{i.i.d.}{\sim} \mathcal{F}$ , where  $\mathcal{F}$  does not correspond to any specific parametric family of distributions but satisfies  $\mathbb{E}(y_i | \mathbf{x}_i) = \mathbf{x}_i^\top \boldsymbol{\beta}$ . Here,  $\mathbf{x}_i$  is a  $K$ -dimensional vector of regressors, and  $\boldsymbol{\beta}$  is a  $K$ -dimensional vector of parameters. The Bayesian bootstrap generates posterior probabilities associated with each  $y_i$ , assigning zero probability to unobserved values of  $\mathbf{y}$ .

The algorithm to implement the Bayesian bootstrap is the following:

---

**Algorithm A15** Bayesian bootstrap from scratch in linear regression

---

- 1: Draw  $\mathbf{g} \sim Dir(\alpha_1, \alpha_2, \dots, \alpha_N)$  such that  $\alpha_i = 1 \forall i$ .
- 2:  $\mathbf{g} = (g_1, g_2, \dots, g_N)$  is the vector of probabilities to attach to  $(y_1, \mathbf{x}_1^\top), (y_2, \mathbf{x}_2^\top), \dots, (y_n, \mathbf{x}_N^\top)$  for each Bayesian bootstrap replication.
- 3: **for**  $s = 1, \dots, S$  **do**
- 4:     Sample  $(y_i, \mathbf{x}_i^\top)$   $N$  times with replacement and probabilities  $g_i$ ,  $i = 1, 2, \dots, N$ .
- 5:     Estimate  $\boldsymbol{\beta}^{(s)}$  using weighted least squares in the model  $\mathbb{E}(\mathbf{y} | \mathbf{X}) = \mathbf{X}\boldsymbol{\beta}$ , where the weights are based on  $g_i$ .\*
- 6: **end for**
- 7: The distribution of  $\boldsymbol{\beta}^{(s)}$  is the posterior distribution of  $\boldsymbol{\beta}$ .

---

\*Ordinary least squares is the posterior mean of  $\boldsymbol{\beta}$  using Jeffrey's prior in a linear regression.

---

### Example: Simulation exercise

Let's perform a simulation exercise to evaluate the performance of Algorithm A15 for inference using the Bayesian bootstrap. The data-generating process includes two regressors, each drawn from a standard normal distribution. The location vector is  $\boldsymbol{\beta} = [1 \ 1]^\top$ , with variance  $\sigma^2 = 1$ , and the sample size is 1,000.

Algorithm A16 illustrates how to use our GUI to implement the Bayesian bootstrap. Our GUI relies on the *bayesboot* command from the *bayesboot* package in **R**. Exercise 11 asks you to use this package to perform inference in this simulation and compare the results with those obtained from our GUI using  $S = 10,000$ .

The following **R** code demonstrates how to implement the Bayesian bootstrap from scratch. The results show that all 95% credible intervals contain the true population parameters, and the posterior means are close to these values.

**Algorithm A16** Bayesian bootstrap in linear regression

- 1: Select *Univariate Models* on the top panel
- 2: Select *Bootstrap* model using the left radio button
- 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend. You should see a preview of the dataset
- 4: Select number of bootstrap replications using the *Range sliders*
- 5: Select dependent and independent variables using the *Formula builder* table
- 6: Click the *Build formula* button to generate the formula in **R** syntax. You can modify the formula in the **Main equation** box using valid arguments of the *formula* command structure in **R**
- 7: Click the *Go!* button
- 8: Analyze results
- 9: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons

***R. code. Bayesian bootstrap***

```

1 rm(list = ls()); set.seed(10101)
2 #--- Data
3 N <- 1000; x1 <- runif(N); x2 <- rnorm(N)
4 X <- cbind(1, x1, x2); B <- c(1, 1, 1); sig2 <- 1
5 y <- as.numeric(X %*% B + rnorm(N, 0, sqrt(sig2)))
6 data <- data.frame(y, x1, x2)
7
8 #--- Bayesian bootstrap (Rubin) for OLS coefficients
9 BB <- function(S, df, alpha = 1) {
10   N <- nrow(df)
11   Betas <- matrix(NA_real_, nrow = S, ncol = 3)
12   colnames(Betas) <- c("(Intercept)", "x1", "x2")
13   for (s in 1:S) {
14     # One Dirichlet weight vector over the N observations
15     w <- as.numeric(LaplacesDemon::rdirichlet(1, rep(alpha,
16       N)))
17     # Weighted least squares = Bayesian bootstrap draw of
18     # the OLS functional
19     fit <- lm(y ~ x1 + x2, data = df, weights = w)
20     Betas[s, ] <- coef(fit)
21   }
22   Betas
23 }
24 S <- 10000
25 BBs <- BB(S = S, df = data, alpha = 1)
26 summary(coda::mcmc(BBs))

```

---

## 6.11 Summary

In this chapter, we present the core univariate regression models and demonstrate how to perform Bayesian inference using Markov Chain Monte Carlo (MCMC) methods. Specifically, we cover several key algorithms: Gibbs sampling, Metropolis–Hastings, nested Metropolis–Hastings, and Metropolis–Hastings-within-Gibbs. These algorithms provide the foundation for conducting Bayesian inference in more complex frameworks based on cross-sectional data.

---

## 6.12 Exercises

1. Derive the posterior conditional distributions of the Gaussian linear model assuming independent priors  $\pi(\boldsymbol{\beta}, \sigma^2) = \pi(\boldsymbol{\beta}) \times \pi(\sigma^2)$ , where  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$  and  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ .
2. Consider the model  $y_i \sim N(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2/\tau_i)$  (Gaussian linear model with heteroskedasticity) and independent priors  $\pi(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\tau}) = \pi(\boldsymbol{\beta}) \times \pi(\sigma^2) \times \prod_{i=1}^N \pi(\tau_i)$ , where  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$ ,  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ , and  $\tau_i \sim G(v/2, v/2)$ .  
Show that  $\boldsymbol{\beta} \mid \sigma^2, \boldsymbol{\tau}, \mathbf{y}, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \mathbf{B}_n)$ ,  $\sigma^2 \mid \boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{y}, \mathbf{X} \sim IG(\alpha_n/2, \delta_n/2)$ , and  $\tau_i \mid \boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X} \sim G(v_{1n}/2, v_{2n}/2)$ , where  $\boldsymbol{\tau} = [\tau_1 \dots \tau_N]^\top$ ,  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \sigma^{-2} \mathbf{X}^\top \Psi \mathbf{X})^{-1}$ ,  $\boldsymbol{\beta}_n = \mathbf{B}_n (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sigma^{-2} \mathbf{X}^\top \Psi \mathbf{y})$ ,  $\alpha_n = \alpha_0 + N$ ,  $\delta_n = \delta_0 + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \Psi (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ ,  $v_{1n} = v + 1$ ,  $v_{2n} = v + \sigma^{-2}(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2$ , and  $\Psi = \text{diag}(\tau_i)$ .
3. **The market value of soccer players in Europe continues.**  
Use the setting of the previous exercise to perform inference using a Gibbs sampling algorithm for the market value of soccer players in Europe, setting  $v = 5$  and using the same hyperparameters as in the homoscedastic case. Is there any meaningful difference in the coefficient associated with the national team compared to the homoscedastic case?
4. **Example: Determinants of hospitalization continues.**  
Program a Gibbs sampling algorithm for the application on the determinants of hospitalization.
5. **Choice of fishing mode continues.**
  - Run Algorithm A9 to report the results of the Geweke [170], Raftery [352], and Heidelberger [202] tests using our GUI.

- Use the command `rmnpGibbs` to replicate the fishing mode choice example.
6. **Simulation exercise: Multinomial logit model continues.**  
Perform inference in the simulated multinomial logit model using the command `rmnlIndepMetrop` from the *bayesm* package in **R**, and compare the results with those obtained from our GUI.
  7. **Simulation of the ordered probit model.**  
Simulate an ordered probit model where the first regressor follows  $N(6, 5)$  and the second follows  $G(1, 1)$ . The location vector is  $\beta = [0.5 \ -0.25 \ 0.5]^\top$ , and the cutoffs are  $\alpha = [0 \ 1 \ 2.5]^\top$ . Program a Metropolis-within-Gibbs sampling algorithm from scratch to perform inference in this simulation.
  8. **Simulation of the negative binomial model continues.**  
Perform inference in the simulated negative binomial model using the *bayesm* package in **R**.
  9. **The market value of soccer players in Europe continues.**  
Perform the application on soccer players' market value with left censoring at one million euros in our GUI using Algorithm A13 and the hyperparameters of the example.
  10. **The market value of soccer players in Europe continues.**  
Program from scratch the Gibbs sampling algorithm for the example of soccer players' market value at the 0.75 quantile.
  11. Use the *bayesboot* package to perform inference in the simulation exercise of Section 6.10 and compare the results with those obtained from our GUI, setting  $S = 10,000$ .



# 7

## Multivariate models

We describe how to perform Bayesian inference in multivariate response models, including multivariate regression, seemingly unrelated regression, instrumental variable, and multivariate probit models. In particular, we present the posterior distributions of the parameters and illustrate their use through selected applications and simulations. We also demonstrate how to perform inference in these models using three levels of programming engagement: through the GUI, via **R** packages, and by directly programming the algorithms from scratch. Finally, we provide a set of mathematical and computational exercises.

Remember that we can run our GUI typing

### ***R code. How to display our graphical user interface***

```
1 shiny::runGitHub("besmarter/BSTApp", launch.browser = T)
```

in the **R** package console or any **R** code editor. However, users should see Chapter 5 for seeing other options.

### 7.1 Multivariate regression

A complete presentation of this model is provided in Section 3.4. Here, we summarize the setting and posterior distributions for ease of exposition. There are  $M$  jointly dependent variables that share the same set of regressors, and their stochastic errors are contemporaneously correlated. Specifically, let  $\mathbf{Y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_M]$  denote an  $N \times M$  matrix generated by

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{U},$$

where  $\mathbf{X}$  is an  $N \times K$  matrix of regressors,  $\mathbf{B} = [\boldsymbol{\beta}_1 \ \boldsymbol{\beta}_2 \ \dots \ \boldsymbol{\beta}_M]$  is a  $K \times M$  matrix of coefficients, and  $\mathbf{U} = [\boldsymbol{\mu}_1 \ \boldsymbol{\mu}_2 \ \dots \ \boldsymbol{\mu}_M]$  is an  $N \times M$  matrix of stochastic errors such that each row  $\boldsymbol{\mu}_i \stackrel{i.i.d.}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$  for  $i = 1, 2, \dots, N$ .

The prior distributions are specified as  $\mathbf{B} | \boldsymbol{\Sigma} \sim N(\mathbf{B}_0, \mathbf{V}_0, \boldsymbol{\Sigma})$  and  $\boldsymbol{\Sigma} \sim IW(\boldsymbol{\Psi}_0, \alpha_0)$ . The conditional posterior distributions are then given by

$$\mathbf{B} | \boldsymbol{\Sigma}, \mathbf{Y}, \mathbf{X} \sim N(\mathbf{B}_n, \mathbf{V}_n, \boldsymbol{\Sigma}),$$

$$\boldsymbol{\Sigma} | \mathbf{Y}, \mathbf{X} \sim IW(\boldsymbol{\Psi}_n, \alpha_n),$$

where  $\mathbf{V}_n = (\mathbf{X}^\top \mathbf{X} + \mathbf{V}_0^{-1})^{-1}$ ,  $\mathbf{B}_n = \mathbf{V}_n(\mathbf{V}_0^{-1}\mathbf{B}_0 + \mathbf{X}^\top \mathbf{X} \hat{\mathbf{B}})$ ,  $\hat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ ,  $\mathbf{S} = (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}})^\top (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}})$ ,  $\boldsymbol{\Psi}_n = \boldsymbol{\Psi}_0 + \mathbf{S} + \mathbf{B}_0^\top \mathbf{V}_0^{-1} \mathbf{B}_0 + \hat{\mathbf{B}}^\top \mathbf{X}^\top \mathbf{X} \hat{\mathbf{B}} - \mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n$ , and  $\alpha_n = \alpha_0 + N$ .

Since the conditional posterior distributions are standard, inference in this model can be carried out using a Gibbs sampling algorithm.

### Example: The effect of institutions on per capita gross domestic product

To illustrate multivariate regression models, we use the dataset from [4], who analyzed the effect of property rights on economic growth.

We begin with the following *simultaneous structural* economic model:<sup>1</sup>

$$\log(\text{pcGDP95}_i) = \beta_1 + \beta_2 \text{PAER}_i + \beta_3 \text{Africa}_i + \beta_4 \text{Asia}_i + \beta_5 \text{Other}_i + u_{1i}, \quad (7.1)$$

$$\text{PAER}_i = \alpha_1 + \alpha_2 \log(\text{pcGDP95}_i) + \alpha_3 \log(\text{Mort}_i) + u_{2i}. \quad (7.2)$$

Here, *pcGDP95*, *PAER*, and *Mort* denote, respectively, the per capita gross domestic product in 1995, the average index of protection against expropriation between 1985 and 1995, and the settler mortality rate during the period of colonization. *Africa*, *Asia*, and *Other* are continent indicators, with *America* serving as the baseline category.

In this model, there is *simultaneous causality* because *GDP* affects *PAER*, and *PAER* simultaneously affects *GDP*. Estimating Equations 7.1 and 7.2 without accounting for this simultaneity leads to posterior mean estimates that are *biased* and *inconsistent* from a frequentist perspective.<sup>2</sup>

A standard approach to address this issue is to estimate the *reduced-form* model—one without simultaneous causality—in which all *endogenous variables* are expressed as functions of the *exogenous variables*. The former are determined within the model (e.g.,  $\log(\text{pcGDP95}_i)$  and *PAER* in this example), while the latter are determined outside the model (e.g.,  $\log(\text{Mort}_i)$ , *Africa*, *Asia*, and *Other*).

<sup>1</sup>This model captures the potential underlying economic relationships among the variables.

<sup>2</sup>Specifically,  $\mathbb{E}[u_{1i}\text{PAER}_i] \neq 0$ , violating a key condition for obtaining unbiased and consistent estimators of  $\boldsymbol{\beta}$ . See Exercise 1.

Substituting Equation 7.2 into Equation 7.1 and solving for  $\log(\text{pcGDP95}_i)$  yields:

$$\log(\text{pcGDP95}_i) = \pi_1 + \pi_2 \log(\text{Mort}_i) + \pi_3 \text{Africa}_i + \pi_4 \text{Asia}_i + \pi_5 \text{Other}_i + e_{1i}. \quad (7.3)$$

Then, substituting Equation 7.3 into Equation 7.2 and solving for  $PAER$  gives:

$$PAER_i = \gamma_1 + \gamma_2 \log(\text{Mort}_i) + \gamma_3 \text{Africa}_i + \gamma_4 \text{Asia}_i + \gamma_5 \text{Other}_i + e_{2i}, \quad (7.4)$$

where  $\pi_2 = \frac{\beta_2 \alpha_3}{1 - \beta_2 \alpha_2}$  and  $\gamma_2 = \frac{\alpha_3}{1 - \beta_2 \alpha_2}$ , assuming  $\beta_2 \alpha_2 \neq 1$ , which ensures model identification (see Exercise 2).

Observe that Equations 7.3 and 7.4 take the form of a multivariate regression model, where the common set of regressors is  $\mathbf{X} = [\log(\text{Mort}) \text{ Africa} \text{ Asia} \text{ Other}]$ , and the common set of dependent variables is  $\mathbf{Y} = [\log(\text{pcGDP95}) \text{ PAER}]$ . Therefore, we can estimate this model using the framework described in this section.

In the first stage, we estimate the parameters of the *reduced-form* model (Equations 7.3 and 7.4); however, our primary interest lies in the parameters of the *structural* model (Equations 7.1 and 7.2). A natural question arises as to whether the *structural* parameters can be recovered (identified) from the *reduced-form* parameters. There are two key criteria to address this question: the *order condition*, which is necessary, and the *rank condition*, which is both necessary and sufficient.<sup>3</sup>

### The order condition

Given a system of equations with  $M$  endogenous variables and  $K$  exogenous variables (including the intercept), there are two ways to assess the order condition:

- The parameters of an equation in the system are identified if there are at least  $M - 1$  variables excluded from that equation (*exclusion restrictions*). The equation is *exactly identified* if the number of excluded variables equals  $M - 1$ , and *over-identified* if the number of excluded variables is greater than  $M - 1$ .
- The parameters of equation  $m$  in the system are identified if  $K - K_m \geq M_m - 1$ , where  $K_m$  and  $M_m$  denote the number of exogenous and endogenous variables in equation  $m$ , respectively. The  $m$ th equation is *exactly identified* if  $K - K_m = M_m - 1$ , and *over-identified* if  $K - K_m > M_m - 1$ .

<sup>3</sup>The posterior distribution of the structural parameters,  $\pi(\boldsymbol{\beta}, \boldsymbol{\alpha} | \boldsymbol{\gamma}, \boldsymbol{\pi}, \mathbf{Y}, \mathbf{X})$ , is proportional to the prior distribution of the structural parameters conditional on the reduced-form parameters,  $\pi(\boldsymbol{\beta}, \boldsymbol{\alpha} | \boldsymbol{\gamma}, \boldsymbol{\pi})$ . In other words, it is the prior distribution of the reduced-form parameters that is updated by the sample information, while the updating of the structural parameters occurs solely through the reduced-form parameters. See Section 9.3 in [465] for further details.

Otherwise, the structural parameters are said to be *under-identified*.

From Equations 7.1 and 7.2 in this example, we have  $K = 5$ ,  $M = 2$ ,  $K_1 = 4$ ,  $K_2 = 2$ ,  $M_1 = 2$ , and  $M_2 = 2$ . Therefore,  $K - K_1 = 1 = M - 1$  and  $K - K_2 = 3 > M - 1 = 1$ . According to the order condition, both equations satisfy the necessary condition for identification: the first equation is *exactly identified*, and the second equation is *over-identified*. Note that one variable is excluded from the first equation, while three variables are excluded from the second equation.

### The rank condition

The rank condition (necessary and sufficient) states that, given a *structural* model with  $M$  equations (and  $M$  endogenous variables), an equation is identified if and only if there exists at least one nonzero determinant of an  $(M-1) \times (M-1)$  matrix constructed from the variables that are excluded from the equation under analysis but included in at least one other equation of the system. Otherwise, the structural parameters are said to be *under-identified*.

To verify this condition, it is useful to construct the *identification matrix*. Table 7.1 presents this matrix for the current example.

**TABLE 7.1**

Identification matrix.

log(pcGDP95)	PAER	Constant	log(Mort)	Africa	Asia	Other
1	$-\beta_2$	$-\beta_1$	0	$-\beta_3$	$-\beta_4$	$-\beta_5$
$-\alpha_2$	1	$-\alpha_1$	$-\alpha_3$	0	0	0

The only excluded variable in the log(pcGDP95) equation is log(Mort). Therefore, only one matrix can be constructed using the excluded variables from this equation, namely  $[-\alpha_3]$  (see column 4 in Table 7.1). The determinant of this matrix is  $-\alpha_3$ , and as long as this coefficient is nonzero (i.e.,  $\alpha_3 \neq 0$ )—meaning that the mortality rate is relevant in the PAER equation—the coefficients in the log(pcGDP95) equation are *exactly identified*. For example,  $\beta_2 = \pi_2/\gamma_2$ , which represents the effect of property rights on GDP, is exactly identified.

It is important to note the role of excluding log(Mort) from the log(pcGDP95) equation while including it in the PAER equation. This exclusion is known as an *exclusion restriction*, which requires an exogenous source of variation in the PAER equation to identify the log(pcGDP95) equation. The presence of relevant exogenous variation is essential for the identification, estimation, and inference of *structural* parameters.

Regarding the identification of the *structural* parameters in the PAER equation, three possible matrices can be constructed:  $[-\beta_3]$ ,  $[-\beta_4]$ , and  $[-\beta_5]$  (see columns 5–7 in Table 7.1). As long as any of these parameters are relevant in the log(pcGDP95) equation, the PAER equation is identified. In this case, the PAER equation is *over-identified*, meaning that multiple ex-

clusion restrictions are available to estimate its parameters. For instance,  $\alpha_2 = \gamma_3/\pi_3 = \gamma_4/\pi_4 = \gamma_5/\pi_5$  (see Exercise 2).

In general, recovering the *structural* parameters from the *reduced-form* parameters can be challenging, as it requires relevant and credible identification restrictions, which may be difficult to justify or find in some empirical applications.<sup>4</sup>

For this example, we set noninformative priors:  $B_0 = [\mathbf{0}_5 \ \mathbf{0}_5]$ ,  $V_0 = 100I_K$ ,  $\Psi_0 = 5I_2$ , and  $\alpha_0 = 5$ .<sup>5</sup> Once the GUI is displayed (see the beginning of this chapter), follow Algorithm A17 to estimate multivariate linear models in the GUI (see Chapter 5 for details, particularly on data setup).

---

**Algorithm A17** Multivariate linear model

---

- 1: Select *Multivariate Models* on the top panel
  - 2: Select *Simple Multivariate* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select the number of dependent variables in the box **Number of endogenous variables: m**
  - 6: Select the number of independent variables (including the intercept) in the box **Number of exogenous variables: k**
  - 7: Set the hyperparameters: mean vectors, covariance matrix, degrees of freedom, and the scale matrix. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Click the *Go!* button
  - 9: Analyze results
  - 10: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

The following **R** code shows how to perform the Gibss sampling algorithm in this example using the dataset *4Institutions.csv*. We ask to run this example using the *rmultireg* command from the *bayesm* package as an exercise. We find that the posterior mean *structural* effect of property rights on GDP is 0.98, and the 95% credible interval is (0.56, 2.87). This means that there is evidence supporting a positive effect of property rights on gross domestic product.

<sup>4</sup>For accessible introductions to identification in linear systems, see [185, Chap. 19] and [457, Chap. 16].

<sup>5</sup>Note that these priors are specified for the *reduced-form* model. This choice may have unintended consequences for the posterior distributions of the *structural* parameters, which are ultimately of interest to researchers. For further discussion, see [250, p. 302].

**R code. The effect of institutions on per capita GDP**

---

```

1 rm(list = ls())
2 set.seed(12345)
3 DataInst <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/4Institutions
  .csv", sep = ",", header = TRUE, quote = "")
4 attach(DataInst)
5 Y <- cbind(logpcGDP95, PAER)
6 X <- cbind(1, logMort, Africa, Asia, Other)
7 M <- dim(Y)[2]
8 K <- dim(X)[2]
9 N <- dim(Y)[1]
10 # Hyperparameters
11 B0 <- matrix(0, K, M)
12 c0 <- 100
13 V0 <- c0*diag(K)
14 Psi0 <- 5*diag(M)
15 a0 <- 5
16 # Posterior parameters
17 Bhat <- solve(t(X)%*%X)%*%t(X)%*%Y
18 S <- t(Y - X%*%Bhat)%*%(Y - X%*%Bhat)
19 Vn <- solve(solve(V0) + t(X)%*%X)
20 Bn <- Vn%*%(solve(V0)%*%B0 + t(X)%*%X%*%Bhat)
21 Psin <- Psi0 + S + t(B0)%*%solve(V0)%*%B0 + t(Bhat)%*%t(X)%*
  %X%*%Bhat - t(Bn)%*%solve(Vn)%*%Bn
22 an <- a0 + N
23 #Posterior draws
24 s <- 10000 #Number of posterior draws
25 SIGs <- replicate(s, LaplacesDemon::rinvwishart(an, Psin))
26 BsCond <- sapply(1:s, function(s) {MixMatrix::rmatrixnorm(n
  = 1, mean=Bn, U = Vn, V = SIGs[, , s])})
27 summary(coda::mcmc(t(BsCond)))
28 SIGMs <- t(sapply(1:s, function(l) {gdata::lowerTriangle(
  SIGs[,,l], diag=TRUE, byrow=FALSE)}))
29 summary(coda::mcmc(SIGMs))
30 hdiBs <- HDInterval::hdi(t(BsCond), credMass = 0.95) #
  Highest posterior density credible interval
31 hdiBs
32 hdiSIG <- HDInterval::hdi(SIGMs, credMass = 0.95) # Highest
  posterior density credible interval
33 hdiSIG
34 beta2 <- BsCond[2,] / BsCond[7,]
35 summary(coda::mcmc(beta2)) # Effect of property rights on
  GDP

```

---

## 7.2 Seemingly unrelated regression

In seemingly unrelated regression (SUR) models, there are  $M$  dependent variables, each with potentially different sets of regressors, and their stochastic errors are contemporaneously correlated. The model is given by

$$\mathbf{y}_m = \mathbf{X}_m \boldsymbol{\beta}_m + \boldsymbol{\mu}_m,$$

where  $\mathbf{y}_m$  is an  $N$ -dimensional vector of observations,  $\mathbf{X}_m$  is an  $N \times K_m$  matrix of regressors,  $\boldsymbol{\beta}_m$  is a  $K_m$ -dimensional vector of coefficients, and  $\boldsymbol{\mu}_m$  is an  $N$ -dimensional vector of stochastic errors, for  $m = 1, 2, \dots, M$ .

Let  $\boldsymbol{\mu}_i = [\mu_{i1} \ \mu_{i2} \ \dots \ \mu_{iM}]^\top$ , where  $\boldsymbol{\mu}_i \sim N(\mathbf{0}, \boldsymbol{\Sigma})$ . Stacking the  $M$  equations, the system can be written as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu},$$

where  $\mathbf{y} = [\mathbf{y}_1^\top \ \mathbf{y}_2^\top \ \dots \ \mathbf{y}_M^\top]^\top$  is an  $MN$ -dimensional vector,  $\boldsymbol{\beta} = [\boldsymbol{\beta}_1^\top \ \boldsymbol{\beta}_2^\top \ \dots \ \boldsymbol{\beta}_M^\top]^\top$  is a  $K$ -dimensional vector with  $K = \sum_{m=1}^M K_m$ , and  $\mathbf{X}$  is an  $MN \times K$  block-diagonal matrix composed of the individual  $\mathbf{X}_m$  matrices:

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{X}_M \end{bmatrix}.$$

Similarly, the error vector is  $\boldsymbol{\mu} = [\boldsymbol{\mu}_1^\top \ \boldsymbol{\mu}_2^\top \ \dots \ \boldsymbol{\mu}_M^\top]^\top$ , an  $MN$ -dimensional vector with  $\boldsymbol{\mu} \sim N(\mathbf{0}, \boldsymbol{\Sigma} \otimes \mathbf{I}_N)$ .

The likelihood function for the parameters is

$$p(\boldsymbol{\beta}, \boldsymbol{\Sigma} | \mathbf{y}, \mathbf{X}) \propto |\boldsymbol{\Sigma}|^{-N/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I}_N) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}.$$

Assuming independent priors,  $\pi(\boldsymbol{\beta}) \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$  and  $\pi(\boldsymbol{\Sigma}^{-1}) \sim W(\alpha_0, \boldsymbol{\Psi}_0)$ , the posterior distributions are

$$\boldsymbol{\beta} | \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \mathbf{B}_n),$$

$$\boldsymbol{\Sigma}^{-1} | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X} \sim W(\alpha_n, \boldsymbol{\Psi}_n),$$

where

$$\mathbf{B}_n = (\mathbf{X}^\top (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I}_N) \mathbf{X} + \mathbf{B}_0^{-1})^{-1}, \quad \boldsymbol{\beta}_n = \mathbf{B}_n (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{X}^\top (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I}_N) \mathbf{y}),$$

$$\alpha_n = \alpha_0 + N, \quad \text{and} \quad \boldsymbol{\Psi}_n = (\boldsymbol{\Psi}_0^{-1} + \mathbf{U}^\top \mathbf{U})^{-1},$$

where  $\mathbf{U}$  is an  $N \times M$  matrix whose columns are given by  $\mathbf{y}_m - \mathbf{X}_m \boldsymbol{\beta}_m$ .

We can show, through straightforward but tedious algebra, that by defining  $\mathbf{y}_i = [y_{i1} \ y_{i2} \ \dots \ y_{iM}]^\top$  and

$$\mathbf{X}_i = \begin{bmatrix} \mathbf{x}_{1i}^\top & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{x}_{2i}^\top & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{x}_{Mi}^\top \end{bmatrix},$$

we can alternatively express

$$\mathbf{B}_n = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^N \mathbf{X}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}_i \right)^{-1}, \quad \boldsymbol{\beta}_n = \mathbf{B}_n \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^N \mathbf{X}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{y}_i \right),$$

and

$$\boldsymbol{\Psi}_n = \left( \boldsymbol{\Psi}_0^{-1} + \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})^\top \right)^{-1}.$$

Since the conditional posteriors are standard, we can use a Gibbs sampling algorithm to obtain the posterior draws.

### Example: Utility demand

Let's use the dataset *Utilities.csv* to estimate a seemingly unrelated regression (SUR) model for utility demand. We adopt the same setting as in Exercise 14 of Chapter 3, where a multivariate regression model is estimated after omitting households with zero consumption in any utility. In this exercise, we observe that not all regressors are relevant for the demand for electricity, water, and gas. Therefore, we estimate the following model:

$$\begin{aligned} \log(\text{electricity}_i) &= \beta_1 + \beta_2 \log(\text{electricity price}_i) + \beta_3 \log(\text{water price}_i) \\ &\quad + \beta_4 \log(\text{gas price}_i) + \beta_5 \text{IndSocio1}_i + \beta_6 \text{IndSocio2}_i + \beta_7 \text{Altitude}_i \\ &\quad + \beta_8 \text{Nrooms}_i + \beta_9 \text{HouseholdMem}_i + \beta_{10} \log(\text{Income}_i) + \mu_{i1} \\ \log(\text{water}_i) &= \alpha_1 + \alpha_2 \log(\text{electricity price}_i) + \alpha_3 \log(\text{water price}_i) \\ &\quad + \alpha_4 \log(\text{gas price}_i) + \alpha_5 \text{IndSocio1}_i + \alpha_6 \text{IndSocio2}_i \\ &\quad + \alpha_7 \text{Nrooms}_i + \alpha_8 \text{HouseholdMem}_i + \mu_{i2} \\ \log(\text{gas}_i) &= \gamma_1 + \gamma_2 \log(\text{electricity price}_i) + \gamma_3 \log(\text{water price}_i) \\ &\quad + \gamma_4 \log(\text{gas price}_i) + \gamma_5 \text{IndSocio1}_i + \gamma_6 \text{IndSocio2}_i + \gamma_7 \text{Altitude}_i \\ &\quad + \gamma_8 \text{Nrooms}_i + \gamma_9 \text{HouseholdMem}_i + \mu_{i3}, \end{aligned}$$

where electricity, water, and gas represent the monthly consumption of electricity (kWh), water (m<sup>3</sup>), and gas (m<sup>3</sup>) by Colombian households. The dataset includes information for 2,103 households, with details on the average prices of electricity (USD/kWh), water (USD/m<sup>3</sup>), and gas (USD/m<sup>3</sup>),

as well as indicators of the socioeconomic conditions of the neighborhoods where the households are located (*IndSocio1* being the lowest and *IndSocio3* the highest). Additionally, the dataset contains information on whether the household is located in a municipality situated above 1,000 meters above sea level, the number of rooms in the dwelling, the number of household members, and monthly income (USD).

Because each equation includes a different set of regressors and we suspect correlation among the stochastic errors of the three equations, a SUR model is appropriate. Unobserved correlation across equations is expected since we are modeling utilities, and in many cases, a single provider manages all three services and issues a single bill.

Algorithm A18 illustrates how to estimate SUR models using our GUI. The GUI implements the command *rsurGibbs* from the *bayesm* package in **R**. See Chapter 5 for additional details, including instructions on how to set up the dataset, and consult the templates available in our GitHub repository (<https://github.com/besmarter/BSTAApp>) under the **DataApp** and **DataSim** folders.

---

**Algorithm A18** Seemingly unrelated regression

---

- 1: Select *Multivariate Models* on the top panel
  - 2: Select *Seemingly Unrelated Regression* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Select the number of dependent variables in the box **Number of endogenous variables: m**
  - 6: Select the number of independent variables in the box **TOTAL number Exogenous Variables: k**. This is the sum of all exogenous variables over all equations including intercepts. In the example of **Utility demand**, it is equal to 27
  - 7: Set the hyperparameters: mean vectors, covariance matrix, degrees of freedom, and the scale matrix. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Click the *Go!* button
  - 9: Analyze results
  - 10: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

The following code shows how to program this application using the *bayesm* package. We use 10,000 MCMC iterations, with  $\beta_0 = \mathbf{0}_{27}$ ,  $B_0 = 100I_{27}$ ,  $\alpha_0 = 5$ , and  $\Psi_0 = 5I_3$ .

The posterior median estimates of the own-price elasticities of demand for electricity, water, and gas are  $-1.88$ ,  $-0.36$ , and  $-0.62$ , respectively, and none

of the 95% credible intervals include zero. This implies that a 1% increase in the prices of electricity, water, and gas results in a 1.88%, 0.36%, and 0.62% decrease in monthly consumption of these utilities, respectively.<sup>6</sup> Overall, the results provide evidence supporting the relevance of most regressors in these equations and indicate the presence of unobserved correlation in the demand for these services, further justifying the use of the SUR model in this application.

### R code. Utility demand in Colombia

```

1 rm(list = ls()); set.seed(10101); library(dplyr)
2 DataUt <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/references/master/DataApp/Utilities.csv"
  , sep = ",", header = TRUE, quote = "")
3 DataUtEst <- DataUt %>% filter(Electricity != 0 & Water !=0
  & Gas != 0)
4 attach(DataUtEst)
5 y1 <- log(Electricity); y2 <- log(Water); y3 <- log(Gas)
6 X1 <- cbind(1, LnPriceElect, LnPriceWater, LnPriceGas,
  IndSocio1, IndSocio2, Altitude, Nrooms, HouseholdMem,
  Lnincome)
7 X2 <- cbind(1, LnPriceElect, LnPriceWater, LnPriceGas,
  IndSocio1, IndSocio2, Nrooms, HouseholdMem)
8 X3 <- cbind(1, LnPriceElect, LnPriceWater, LnPriceGas,
  IndSocio1, IndSocio2, Altitude, Nrooms, HouseholdMem)
9 regdata <- NULL
10 regdata[[1]] <- list(y = y1, X = X1); regdata[[2]] <- list(y
  = y2, X = X2); regdata[[3]] <- list(y = y3, X = X3)
11 M <- length(regdata); K1 <- dim(X1)[2]; K2 <- dim(X2)[2]; K3
  <- dim(X3)[2]
12 K <- K1 + K2 + K3
13 # Hyperparameters
14 b0 <- rep(0, K); c0 <- 100; B0 <- c0*diag(K); V <- 5*diag(M)
  ; a0 <- M
15 Prior <- list(betabar = b0, A = solve(B0), nu = a0, V = V)
16 #Posterior draws
17 S <- 10000; keep <- 1; Mcmc <- list(R = S, keep = keep)
18 PosteriorDraws <- bayesm::rsurGibbs(Data = list(regdata =
  regdata), Mcmc = Mcmc, Prior = Prior)

```

<sup>6</sup>This is an example where concerns about *biased* and *inconsistent* posterior mean estimates may arise, for instance, due to *simultaneous causality* between quantity and price. These concerns are valid; however, we use micro-level data, which implies no price-quantity simultaneity. Moreover, the utility providers operate in regulated natural monopoly markets, mitigating potential endogeneity arising from searching provider strategies. Finally, prices were obtained directly from provider records, reducing possible measurement errors [360].

### R code. Utility demand in Colombia, results

```

1 Bs <- PosteriorDraws[["betadraw"]]
2 Names <- c("Const", "LnPriceElect", "LnPriceWater", "
   LnPriceGas", "IndSocio1", "IndSocio2", "Altitude", "
   Nrooms", "HouseholdMem", "Lnincome", "Const", "
   LnPriceElect", "LnPriceWater", "LnPriceGas", "IndSocio1"
   , "IndSocio2",
3 "Nrooms", "HouseholdMem", "Const", "LnPriceElect", "
   LnPriceWater", "LnPriceGas", "IndSocio1", "IndSocio2", "
   Altitude", "Nrooms", "HouseholdMem")
4 colnames(Bs) <- Names
5 summary(coda::mcmc(Bs))
6 summary(coda::mcmc(PosteriorDraws[["Sigmadraw"]]))

```

In Exercise 5, we ask you to run this application using our GUI and the dataset *Utilities.csv*. Note that this file must be modified to match the structure required by our GUI (see the dataset *5Institutions.csv* in the *DataApp* folder of our GitHub repository — <https://github.com/besmarter/BSTApp> — for a template). In addition, we ask you to program the Gibbs sampling algorithm for this application from scratch.

### 7.3 Instrumental variable

This inferential approach is used when there are *endogeneity* issues, that is, when the stochastic error term is not independent of the regressors. This dependence generates *bias* in posterior mean estimates when using an inferential approach that does not account for it. *Endogeneity* can arise from *simultaneous causality*, *omitted relevant correlated variables*, or *measurement error* in the regressors.<sup>7</sup>

Let the dependent variable be modeled as a linear function of one endogenous regressor and a set of exogenous regressors:

$$y_i = \mathbf{x}_{ei}^\top \boldsymbol{\beta}_1 + \beta_s x_{si} + \mu_i,$$

where

$$x_{si} = \mathbf{x}_{ei}^\top \boldsymbol{\gamma}_1 + \mathbf{z}_i^\top \boldsymbol{\gamma}_2 + v_i.$$

---

<sup>7</sup>See [457, Chap. 15] for an introductory treatment of instrumental variables under the frequentist inferential approach.

Here,  $x_{si}$  is the variable that generates endogeneity ( $\mathbb{E}[\mu | x_s] \neq 0$ ),  $\mathbf{x}_e$  is a vector of  $K_1$  exogenous regressors ( $\mathbb{E}[\mu | \mathbf{x}_e] = \mathbf{0}$ ), and  $\mathbf{z}$  is a vector of  $K_2$  instruments. The instruments influence  $x_s$  ( $\mathbb{E}[x_s \mathbf{z}] \neq \mathbf{0}$ ) but do not directly affect  $y$  ( $\mathbb{E}[yz | x_s] = \mathbf{0}$ ). The equation for  $y$  is known as the *structural equation*, which is the main object of inference.

Assuming  $(\mu_i, v_i)^\top \stackrel{i.i.d.}{\sim} N(\mathbf{0}, \Sigma)$ , with  $\Sigma = [\sigma_{lm}]$ ,  $l, m = 1, 2$ , the likelihood function is

$$\begin{aligned} p(\beta, \gamma, \Sigma | \mathbf{y}, \mathbf{X}, \mathbf{Z}) &= \frac{1}{(2\pi)^{N/2} |\Sigma|^{N/2}} \\ &\times \exp \left\{ -\frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{x}_i^\top \beta, x_{si} - \mathbf{w}_i^\top \gamma) \Sigma^{-1} \begin{pmatrix} y_i - \mathbf{x}_i^\top \beta \\ x_{si} - \mathbf{w}_i^\top \gamma \end{pmatrix} \right\}, \end{aligned}$$

where  $\beta = [\beta_1^\top \ \beta_s]^\top$ ,  $\gamma = [\gamma_1^\top \ \gamma_2^\top]^\top$ ,  $\mathbf{x}_i = [\mathbf{x}_{ei}^\top \ x_{si}]^\top$ , and  $\mathbf{w}_i = [\mathbf{x}_{ei}^\top \ \mathbf{z}_i^\top]^\top$ .

We obtain standard conditional posterior densities by specifying independent priors:

$$\gamma \sim N(\gamma_0, \mathbf{G}_0), \quad \beta \sim N(\beta_0, \mathbf{B}_0), \quad \Sigma^{-1} \sim W(\alpha_0, \Psi_0).$$

The conditional posterior distributions are:

$$\beta | \gamma, \Sigma, \mathbf{y}, \mathbf{X}, \mathbf{Z} \sim N(\beta_n, \mathbf{B}_n),$$

$$\gamma | \beta, \Sigma, \mathbf{y}, \mathbf{X}, \mathbf{Z} \sim N(\gamma_n, \mathbf{G}_n),$$

$$\Sigma^{-1} | \beta, \gamma, \mathbf{y}, \mathbf{X}, \mathbf{Z} \sim W(\alpha_n, \Psi_n),$$

where  $\beta_n = \mathbf{B}_n \left( \mathbf{B}_0^{-1} \beta_0 + \omega_1^{-1} \sum_{i=1}^N \mathbf{x}_i \left( y_i - \frac{\sigma_{12}(x_{si} - \mathbf{w}_i^\top \gamma)}{\sigma_{22}} \right) \right)$ ,  $\mathbf{B}_n = \left( \omega_1^{-1} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top + \mathbf{B}_0^{-1} \right)^{-1}$ ,  $\omega_1 = \sigma_{11} - \frac{\sigma_{12}^2}{\sigma_{22}}$ ,  $\mathbf{G}_n = \left( \omega_2^{-1} \sum_{i=1}^N \mathbf{w}_i \mathbf{w}_i^\top + \mathbf{G}_0^{-1} \right)^{-1}$ ,  $\gamma_n = \mathbf{G}_n \left( \mathbf{G}_0^{-1} \gamma_0 + \omega_2^{-1} \sum_{i=1}^N \mathbf{w}_i \left( x_{si} - \frac{\sigma_{12}(y_i - \mathbf{x}_i^\top \beta)}{\sigma_{11}} \right) \right)$ ,  $\omega_2 = \sigma_{22} - \frac{\sigma_{12}^2}{\sigma_{11}}$ ,  $\Psi_n = \left[ \Psi_0^{-1} + \sum_{i=1}^N \begin{pmatrix} y_i - \mathbf{x}_i^\top \beta \\ x_{si} - \mathbf{w}_i^\top \gamma \end{pmatrix} (y_i - \mathbf{x}_i^\top \beta, x_{si} - \mathbf{w}_i^\top \gamma) \right]^{-1}$ ,  $\alpha_n = \alpha_0 + N$ , and  $\sigma_{lm}$  are the elements of  $\Sigma$ .

We also use a Gibbs sampling algorithm in this model since we have standard conditional posterior distributions.

### Example: Simulation exercise

Let's simulate the simple process  $y_i = \beta_1 + \beta_2 x_{si} + \mu_i$  and  $x_{si} = \gamma_1 + \gamma_2 z_i + v_i$  where  $[\mu_i \ v_i]^\top \sim N(\mathbf{0}, \Sigma)$ ,  $\Sigma = [\sigma_{lj}]$  such that  $\sigma_{12} \neq 0$ ,  $i = 1, 2, \dots, 100$ .

Observe that  $\mu | v \sim N\left(\frac{\sigma_{12}}{\sigma_{22}}v, \ \sigma_{11} - \frac{\sigma_{21}^2}{\sigma_{22}}\right)$ . This implies that  $\mathbb{E}[\mu | x_s] = \mathbb{E}[\mu | v] = \frac{\sigma_{12}}{\sigma_{22}}v \neq 0$  when  $\sigma_{12} \neq 0$ , while  $\mathbb{E}[\mu | z] = 0$ .

Let all location parameters be equal to 1, and set  $\sigma_{11} = \sigma_{22} = 1$ ,  $\sigma_{12} = 0.8$ , and  $z \sim N(0, 1)$ . From the large-sample properties of the posterior mean, we

know that it converges to the maximum likelihood estimator (see Section 1.1 and [260, 439]). In this setting, the estimator is

$$\hat{\beta}_2 = \frac{\widehat{\text{Cov}}(x_s, y)}{\widehat{\text{Var}}(x_s)},$$

which converges in probability to

$$\beta_2 + \frac{\sigma_{12}}{\sigma_{22}\text{Var}(x_s)} = \beta_2 + \frac{\sigma_{12}}{\sigma_{22}(\gamma_2^2\text{Var}(z) + \sigma_{22})} = 1.4.$$

Thus, the asymptotic bias when using the posterior mean of a linear regression without accounting for endogeneity is 0.4 in this example.

We assess the sampling performance of the Bayesian estimators by simulating this setting 100 times. The following code demonstrates how to perform this exercise using a linear model that does not account for the *endogeneity* problem (see Section 6.1), as well as how to implement the instrumental variable model using the function *rivGibbs* from the *bayesm* package.<sup>8</sup>

In this setup, we use  $B_0 = 1000I_2$ ,  $\beta_0 = 0_2$ , and set the parameters of the inverse-gamma distribution to 0.0005. For the instrumental variable model, we additionally specify  $\gamma_0 = 0_2$ ,  $G_0 = 1000I_2$ ,  $\alpha_0 = 3$ , and  $\Psi_0 = 3I_2$ .

### R code. Simulation exercise, sampling properties ordinary and instrumental models

```

1 rm(list = ls()); set.seed(10101)
2 N <- 100; k <- 2
3 B <- rep(1, k); G <- rep(1, 2); s12 <- 0.8
4 SIGMA <- matrix(c(1, s12, s12, 1), 2, 2)
5 z <- rnorm(N); Z <- cbind(1, z); w <- matrix(1, N, 1); S <-
6 100
7 U <- replicate(S, MASS::mvrnorm(n = N, mu = rep(0, 2), SIGMA
8 ))
9 x <- G[1] + G[2]*z + U[,2,]; y <- B[1] + B[2]*x + U[,1,]
# Hyperparameters
10 d0 <- 0.001/2; a0 <- 0.001/2
11 b0 <- rep(0, k); c0 <- 1000; B0 <- c0*diag(k)
12 Boi <- solve(B0); g0 <- rep(0, 2)
13 Go <- 1000*diag(2); Goi <- solve(Go)
14 nu <- 3; Psi0 <- nu*diag(2)
# MCMC parameters
15 mcmc <- 5000; burnin <- 1000
16 tot <- mcmc + burnin; thin <- 1

```

<sup>8</sup>It appears that this function does not explicitly account for the effects of exogenous regressors in the equation for the endogenous regressor.

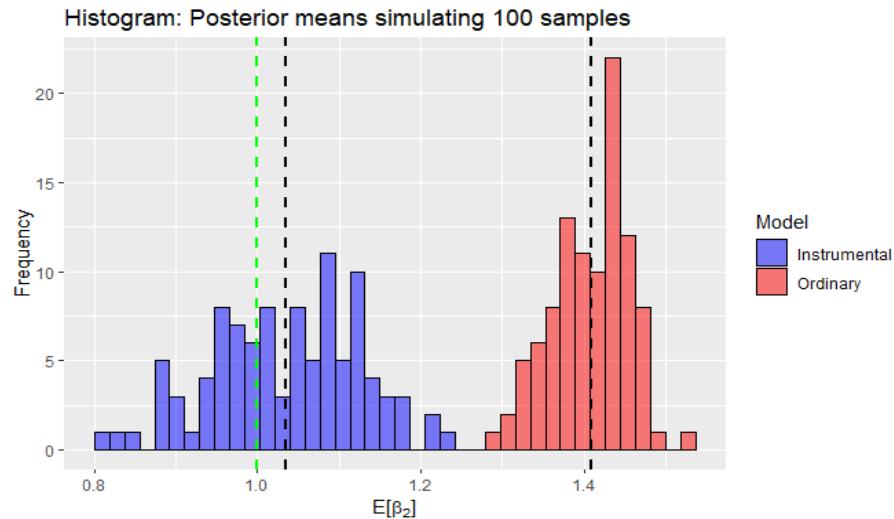
**R code. Simulation exercise, sampling properties ordinary and instrumental models**

```

1 # Gibbs sampling
2 Gibbs <- function(x, y){
3   Data <- list(y = y, x = x, w = w, z = Z)
4   Mcmc <- list(R = mcmc, keep = thin, nprint = 0)
5   Prior <- list(md = g0, Ad = G0i, mbg = b0, Abg = B0i, nu =
       nu, V = Psi0)
6   RestIV <- bayesm::rivGibbs(Data = Data, Mcmc = Mcmc, Prior
       = Prior)
7   PostBIV <- mean(RestIV[["betadraw"]])
8   ResLM <- MCMCpack::MCMCregress(y ~ x + w - 1, b0 = b0, B0
       = B0i, c0 = a0, d0 = d0)
9   PostB <- mean(ResLM[,1]); Res <- c(PostB,PostBIV)
10  return(Res)
11 }
12 PosteriorMeans <- sapply(1:S, function(s) {Gibbs(x = x[,s],
       y = y[,s])})
13 rowMeans(PosteriorMeans)
14 Model <- c(replicate(S, "Ordinary"), replicate(S, "
       Instrumental"))
15 postmeans <- c(t(PosteriorMeans))
16 df <- data.frame(postmeans, Model, stringsAsFactors = FALSE)
17 library(ggplot2); library(latex2exp)
18 histExo <- ggplot(df, aes(x = postmeans, fill = Model)) +
       geom_histogram(bins = 40, position = "identity", color =
       "black", alpha = 0.5) + labs(title = "Overlaid
       Histograms", x = "Value", y = "Count") + scale_fill_
       manual(values = c("blue", "red")) + geom_vline(aes(
       xintercept = mean(postmeans[1:S])), color = "black",
       linewidth = 1, linetype = "dashed") + geom_vline(aes(
       xintercept = mean(postmeans[101:200])), color = "black",
       linewidth = 1, linetype = "dashed") + geom_vline(aes(
       xintercept = B[2]), color = "green", linewidth = 1,
       linetype = "dashed") + xlab(TeX("$E[\beta_2]$")) + ylab
       ("Frequency") + ggtitle("Histogram: Posterior means
       simulating 100 samples")
19 histExo
20

```

Figure 7.1 displays the histograms of the posterior means of  $\beta_2$  obtained from the ordinary model, which does not account for endogeneity, and from the instrumental variable model. On one hand, the mean of the posterior means for the ordinary model is 1.41 (black dashed line in the red histogram), implying a bias of 0.41, very close to the theoretical bias of 0.40. On the other hand, the mean of the posterior means for the instrumental variable model



**FIGURE 7.1**  
Histogram of posterior means: Ordinary and instrumental models.

is 1.04 (black dashed line in the blue histogram), which is close to the true population value of  $\beta_2 = 1$  (green dashed line).

We also observe that the histogram of the posterior means for the ordinary model is less dispersed. In other words, this estimator is more efficient, a well-known result in the frequentist inferential framework when comparing ordinary least squares (OLS) and two-stage least squares (2SLS) estimators (see [455, Chap. 5]).

Two key aspects in the instrumental variables literature are the *weakness* and *exogeneity* of the instruments. The former refers to the strength of the relationship between the instruments and the endogenous regressors, while the latter concerns the independence of the instruments from the stochastic error in the *structural equation*. In Exercise 6, you are asked to use the previous code as a baseline to study these two aspects.

Observe the link between the *weakness* and *exogeneity* of the instruments and the *exclusion restrictions* ( $\mathbb{E}[x_s z] \neq \mathbf{0}$  and  $\mathbb{E}[yz | x_s] = \mathbf{0}$ ). This is the point of departure for [104], who propose assessing the plausibility of the *exclusion restrictions* by defining *plausible exogeneity* as having prior information suggesting that the effect of the instrument in the *structural equation* is close to zero, though perhaps not exactly zero. See Chapter 13, particularly Section 13.4, for further discussion of instrumental variables.

Algorithm A19 can be used to estimate the instrumental variable model using our GUI. We ask in Exercise 8 to replicate the example of the effect of institutions on per capita GDP using our GUI.

**Algorithm A19** Instrumental variable model

- 
- 1: Select *Multivariate Models* on the top panel
  - 2: Select *Variable instrumental (two equations)* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Write down the formula of the structural equation in the **Main Equation** box. This formula must be written using the syntax of the *formula* command of **R** software. This equation includes intercept by default, do not include it in the equation
  - 6: Write down the formula of the endogenous regressor in the **Instrumental Equation** box. This formula must be written using the syntax of the *formula* command of **R** software. This equation includes intercept by default, do not include it in the equation
  - 7: Set the hyperparameters: mean vectors, covariance matrices, degrees of freedom, and the scale matrix. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Click the *Go!* button
  - 9: Analyze results
  - 10: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

---

## 7.4 Multivariate probit model

In the multivariate probit model [131], the response variable  $y_{il} \in \{0, 1\}$  indicates whether individual  $i$  selects option  $l$  among  $L$  mutually exclusive binary choices, where  $l = 1, 2, \dots, L$  and  $i = 1, 2, \dots, N$ . Specifically,

$$y_{il} = \begin{cases} 0, & \text{if } y_{il}^* \leq 0, \\ 1, & \text{if } y_{il}^* > 0, \end{cases}$$

where  $\mathbf{y}_i^* = \mathbf{X}_i \boldsymbol{\beta} + \boldsymbol{\mu}_i \stackrel{iid}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$ . Here,  $\mathbf{y}_i^*$  is an unobserved latent  $L$ -dimensional vector,  $\mathbf{X}_i = \mathbf{x}_i^\top \otimes \mathbf{I}_L$  is an  $L \times K$  design matrix of regressors with  $K = L \times k$ , and  $k$  is the number of regressors (i.e., the dimension of  $\mathbf{x}_i$ ). In addition,  $\boldsymbol{\beta} = [\boldsymbol{\beta}_1^\top \ \boldsymbol{\beta}_2^\top \ \dots \ \boldsymbol{\beta}_k^\top]^\top$ , where each  $\boldsymbol{\beta}_j$  is an  $L$ -dimensional vector of coefficients for  $j = 1, 2, \dots, k$ .

The likelihood function for this model is given by

$$p(\boldsymbol{\beta}, \boldsymbol{\Sigma} | \mathbf{y}, \mathbf{X}) = \prod_{i=1}^N \prod_{l=1}^L p_{il}^{y_{il}},$$

where  $p_{il} = p(y_{il}^* \geq 0)$ .

Note that  $p(y_{il}^* \geq 0) = p(\lambda_{ll}y_{il}^* \geq 0)$  for any  $\lambda_{ll} > 0$ . This leads to identification issues because only the correlation matrix can be identified, similar to the univariate probit model, where the variance is fixed to 1 for identification. We follow the post-processing strategy proposed by [131] to obtain identified parameters, namely,

$$\tilde{\beta} = \text{vec}\{\Lambda B\} \quad \text{and} \quad R = \Lambda \Sigma \Lambda,$$

where  $\Lambda = \text{diag}\{\sigma_{ll}\}^{-1/2}$  and  $B = [\beta_1 \ \beta_2 \dots \beta_k]$ .<sup>9</sup>

We assume independent priors  $\beta \sim N(\beta_0, B_0)$  and  $\Sigma^{-1} \sim W(\alpha_0, \Psi_0)$ . Using data augmentation for  $\mathbf{y}^*$ , we can apply Gibbs sampling, since this becomes a standard Bayesian linear regression model.

The posterior conditional distributions are:

$$\begin{aligned} \beta \mid \Sigma, \mathbf{w} &\sim N(\beta_n, B_n), \\ \Sigma^{-1} \mid \beta, \mathbf{w} &\sim W(\alpha_n, \Psi_n), \\ y_{il}^* \mid \mathbf{y}_{i,-l}^*, \beta, \Sigma^{-1}, \mathbf{y}_i &\sim TN_{I_{il}}(m_{il}, \tau_{ll}^2), \end{aligned}$$

where

$$\begin{aligned} B_n &= (B_0^{-1} + \mathbf{X}^{*\top} \mathbf{X}^*)^{-1}, \quad \beta_n = B_n(B_0^{-1}\beta_0 + \mathbf{X}^{*\top} \mathbf{y}^{**}), \\ \Sigma^{-1} &= \mathbf{C}^\top \mathbf{C}, \quad \mathbf{X}_i^* = \mathbf{C}^\top \mathbf{X}_i, \quad \mathbf{y}_i^{**} = \mathbf{C}^\top \mathbf{y}_i^*, \\ \alpha_n &= \alpha_0 + N, \quad \Psi_n = \left( \Psi_0 + \sum_{i=1}^N (\mathbf{y}_i^* - \mathbf{X}_i \beta)(\mathbf{y}_i^* - \mathbf{X}_i \beta)^\top \right)^{-1}, \\ m_{il} &= \mathbf{x}_{il}^\top \beta + \mathbf{f}_l^\top (\mathbf{y}_{i,-l}^* - \mathbf{X}_{i,-l} \beta), \quad \tau_{ll}^2 = \sigma_{ll} - \omega_{l,-l}^\top \Sigma_{-l,-l}^{-1} \omega_{-l,l}, \end{aligned}$$

and

$\mathbf{f}_l^\top = \omega_{l,-l}^\top \Sigma_{-l,-l}^{-1}$ ,  $\mathbf{y}_{i,-l}^*$  is the  $(L-1)$ -dimensional subvector of  $\mathbf{y}_i^*$  excluding  $y_{il}^*$ ,

$\mathbf{x}_{il}^\top$  is the  $l$ -th row of  $\mathbf{X}_i$ ,  $\mathbf{X}_{i,-l}$  is  $\mathbf{X}_i$  with the  $l$ -th row removed,

$$\mathbf{X}^* = \begin{bmatrix} \mathbf{X}_1^* \\ \mathbf{X}_2^* \\ \vdots \\ \mathbf{X}_N^* \end{bmatrix}, \quad I_{il} = \begin{cases} y_{il}^* > 0, & \text{if } y_{il} = 1, \\ y_{il}^* \leq 0, & \text{if } y_{il} = 0, \end{cases} \quad \text{and} \quad \Sigma = \begin{bmatrix} \omega_1^\top \\ \omega_2^\top \\ \vdots \\ \omega_L^\top \end{bmatrix}.$$

The setting in our GUI has same regressors in each binary decision. However, we can see that the multivariate probit model is similar to a SUR model

---

<sup>9</sup>In a Bayesian framework, a model can be non-identified; however, the posterior distribution of the parameters exists as long as proper prior distributions are specified [131].

in latent variables. We ask in Exercise 9 to implement a Gibbs sampling algorithm for a multivariate probit model with different regressors in each equation.

**Example: Self selection in hospitalization due to a subsidized health care program**

We use the dataset *7HealthMed.csv*, where the dependent variable is  $y = [\text{Hosp } \text{SHI}]^\top$ , with Hosp = 1 if an individual was hospitalized in the year prior to the survey (0 otherwise), and SHI = 1 if the individual was enrolled in a subsidized health insurance program (0 otherwise).

Recall that in our previous application of binary response models, we examined the determinants of hospitalization in Medellín (Colombia), where one of the regressors was a binary indicator of participation in a subsidized health care program (see Section 6.3). If we suspect dependence between the decisions related to these two variables, a bivariate probit model is appropriate.

A priori, we would expect that being enrolled in a subsidized health care program increases the probability of hospitalization *ceteris paribus*, since it reduces the patient's out-of-pocket cost. However, if an individual anticipates a future hospitalization and the factors influencing this expectation are unobserved by the modeler, a feedback effect may arise from hospitalization to enrollment in the subsidized health care program.

---

**Algorithm A20** Multivariate probit model

---

- 1: Select *Multivariate Models* on the top panel
  - 2: Select *Multivariate Probit* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Write down the number of cross-sectional units in the **Number of individuals: n** box
  - 6: Write down the number of exogenous variables in the **Number of exogenous variables: k** box
  - 7: Write down the number of choices in the **Number of choices: l** box
  - 8: Set the hyperparameters: mean vectors, covariance matrix, degrees of freedom, and the scale matrix. This step is not necessary as by default our GUI uses non-informative priors
  - 9: Click the *Go!* button
  - 10: Analyze results
  - 11: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

We consider seven regressors: a constant, gender (female), age, self-perceived health status (with categories fair, good, and excellent, using bad as

the reference category), and the proportion of the individual's age spent living in their neighborhood. The latter variable serves as a proxy for social capital, which may influence enrollment in the subsidized health insurance program, as the target population is identified by local government authorities [355]. The dataset comprises 12,975 individuals who can "choose" between two outcomes: hospitalization and enrollment in the subsidized health insurance regime.

Algorithm A20 shows how to run a multivariate probit model using our GUI.

### **R code. Self selection in hospitalization**

```

1 rm(list = ls()); set.seed(10101)
2 Data <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTAApp/refs/heads/master/DataApp/7HealthMed.
  csv", sep = ",", header = TRUE, quote = "")
3 attach(Data); str(Data)
4 p <- 2; nd <- 7; N <- length(y)/p; y <- y
5 Xd <- as.matrix(Data[seq(1, p*N, 2),3:9])
6 XcreateMP<-function(p,nxs,nind,Data){
7   pandterm = function(message) {
8     stop(message, call. = FALSE)
9   }
10  if (missing(nxs))
11    pandterm("requires number of regressors: include intercept
      if required")
12  if (missing(nind))
13    pandterm("requires number of units (individuals)")
14  if (missing(Data))
15    pandterm("requires dataset")
16  if (nrow(Data)!=nind*2)
17    pandterm("check dataset! number of units times number
      alternatives should be equal to dataset rows")
18  XXDat<-array(0,c(p,1+nxs,nind))
19  XX<-array(0,c(p,nxs*p,nind))
20  YY<-array(0,c(p,1,nind))
21  is<- seq(p,nind*p,p)
22  cis<- seq(nxs,nxs*p+1,nxs)
23  for(i in is){
24    j<-which(i==is)
25    XXDat[,j]<-as.matrix(Data[c((i-(p-1)):i),-1])
26    YY[,j]<-XXDat[,1,j]
27    for(l in 1:p{
28      XX[l,((cis[l]-nxs+1):cis[l]),j]<-XXDat[l,-1,j]
29    }
30  }
31  return(list(y=YY,X=XX))
32 }
```

### R code. Self selection in hospitalization

```

1 Dat <- XcreateMP(p = p, nxs = nd, nind = N, Data = Data)
2 y<-NULL; X<-NULL
3 for(i in 1:dim(Dat$y)[3]){
4   y<-c(y,Dat$y[, ,i])
5   X<-rbind(X,Dat$X[, ,i])
6 }
7 DataMP = list(p=p, y=y, X=X)
8 # Hyperparameters
9 k <- dim(X)[2]; b0 <- rep(0, k); c0 <- 1000
10 B0 <- c0*diag(k); B0i <- solve(B0)
11 a0 <- p - 1 + 3; Psi0 <- a0*diag(p)
12 Prior <- list(betabar = b0, A = B0i, nu = a0, V = Psi0)
13 # MCMC parameters
14 mcmc <- 20000; thin <- 5;
15 Mcmc <- list(R = mcmc, keep = thin)
16 Results <- bayesm::rmvpGibbs(Data = DataMP, Mcmc = Mcmc,
17   Prior = Prior)
17 betatilde1 <- Results$betadraw[,1:7] / sqrt(Results$
18   sigmadraw[,1])
18 summary(coda::mcmc(betatilde1))
19 betatilde2 <- Results$betadraw[,8:14] / sqrt(Results$
19   sigmadraw[,4])
20 summary(coda::mcmc(betatilde2))
21 sigmadraw12 <- Results$sigmadraw[,3] / (Results$sigmadraw
21   [,1]*Results$sigmadraw[,4])^0.5
22 summary(coda::mcmc(sigmadraw12))

```

We set 20,000 MCMC iterations with a thinning parameter of 5. The hyperparameters are specified as  $\beta_0 = \mathbf{0}_{14}$ ,  $B_0 = 100I_{14}$ ,  $\alpha_0 = 4$ , and  $\Psi_0 = 4I_2$ .<sup>10</sup>

The previous **R** code illustrates how to obtain posterior draws using the *rmvpGibbs* command from the *bayesm* package. The results suggest that females, older individuals, and those who perceive their health as poor are more likely to be hospitalized. Moreover, females, older individuals, and individuals with poor or fair self-perceived health who have spent a larger proportion of their lives in their current neighborhoods are more likely to be enrolled in the subsidized health care system. Finally, the results indicate no evidence of unobserved correlation between the two equations, as the 95% credible interval for the correlation parameter is (-0.07, 0.06).

---

<sup>10</sup>Note that the order of the location coefficients in our GUI follows the order of the equations rather than the regressors, as in the theoretical framework presented in this section. This distinction is important for correctly specifying the hyperparameters and interpreting the estimated location parameters.

---

## 7.5 Summary

In this chapter, we present the setting and posterior distributions of several common multivariate models. The multivariate framework enables us to address *endogeneity* issues by exploiting the conditional distribution of a multivariate normal vector. Moreover, in these models, the posterior conditional distributions belong to standard families, multivariate normal, Wishart, and truncated normal. This property allows the straightforward implementation of the Gibbs sampling algorithm across all these models.

---

## 7.6 Exercises

1. Show that  $\mathbb{E}[u_1 \text{PAER}] = \frac{\alpha_1}{1 - \beta_1 \alpha_1} \sigma_1^2$ , assuming that  $\mathbb{E}[u_1 u_2] = 0$  and  $\text{Var}(u_1) = \sigma_1^2$ , in the example on the effect of institutions on per capita GDP.
2. Show that  $\beta_1 = \pi_1 / \gamma_1$  in the example on the effect of institutions on per capita GDP.

3. **The effect of institutions on per capita GDP continues I**

Use the *rmultireg* command from the *bayesm* package to perform inference in the example of the effect of institutions on per capita GDP.

4. **Demand and supply simulation**

Consider the structural demand–supply model:

$$\begin{aligned} q_i^d &= \beta_1 + \beta_2 p_i + \beta_3 y_i + \beta_4 pc_i + \beta_5 ps_i + u_{i1}, \\ q_i^s &= \alpha_1 + \alpha_2 p_i + \alpha_3 er_i + u_{i2}, \end{aligned}$$

where  $q^d$  and  $q^s$  represent demand and supply, respectively;  $p$ ,  $y$ ,  $pc$ ,  $ps$ , and  $er$  denote price, income, complementary price, substitute price, and exchange rate. The variables  $pc$  and  $ps$  refer to the prices of complementary and substitute goods for  $q$ . Assume that  $\boldsymbol{\beta} = [5 \ -0.5 \ 0.8 \ -0.4 \ 0.7]^\top$ ,  $\boldsymbol{\alpha} = [-2 \ 0.5 \ -0.4]^\top$ ,  $u_1 \sim N(0, 0.5^2)$ , and  $u_2 \sim N(0, 0.5^2)$ . Additionally, assume that  $y \sim N(10, 1)$ ,  $pc \sim N(5, 1)$ ,  $ps \sim N(5, 1)$ , and  $er \sim N(15, 1)$ .

- Derive the *reduced-form* model under the equilibrium condition  $q^d = q^s$ , which defines the observable quantity  $q$ .
- Simulate  $p$  and  $q$  from the *reduced-form* equations.

- Perform inference for the *reduced-form* model using the *rmultireg* command from the *bayesm* package.
- Use the posterior draws of the *reduced-form* parameters to perform inference for the *structural* parameters. Discuss any potential issues. *Hint:* Are all structural parameters exactly identified?

#### 5. Utility demand continues

- Run the **Utility demand** application using our GUI and the dataset *Utilities.csv*. *Hint:* This file must be modified to match the structure required by our GUI (see the dataset *5Institutions.csv* in the *DataApp* folder of our GitHub repository—<https://github.com/besmarter/BSTAApp>—as a template).
- Program the Gibbs sampling algorithm for this application from scratch.

#### 6. Simulation exercise of instrumental variables continues I

- Use the setting of the simulation exercise with instrumental variables to analyze the impact of a weak instrument. For example, set  $\gamma_2 = 0.2$  and compare the posterior mean estimates of the ordinary and instrumental variable models.
- Perform a simulation to analyze how the degree of exogeneity of the instrument affects the performance of the posterior mean in the instrumental variable model.

#### 7. Simulation exercise of instrumental variables continues II

Program the Gibbs sampling algorithm for the instrumental variable model from scratch.

#### 8. The effect of institutions on per capita GDP continues II

Estimate the structural Equation 7.1 using the instrumental variable model, where  $\log(Mort)$  serves as the instrument for PAER. Compare the estimated effect of property rights on per capita GDP with that obtained in the example on the effect of institutions on per capita GDP. Use the file *6Institutions.csv* to conduct this exercise in our GUI, setting  $B_0 = 100I_5$ ,  $\beta_0 = \mathbf{0}_5$ ,  $\gamma_0 = \mathbf{0}_2$ ,  $G_0 = 100I_2$ ,  $\alpha_0 = 3$ , and  $\Psi_0 = 3I_2$ . Set the MCMC parameters as follows: 50,000 iterations, 1,000 burn-in, and a thinning interval of 5.

#### 9. Multivariate probit with different regressors

Simulate the model

$$y_{i1}^* = 0.5 - 1.2x_{i11} + 0.7x_{i12} + 0.8x_{i13} + \mu_{i1}, \quad y_{i2}^* = 1.5 - 0.8x_{i21} + 0.5x_{i22} + \mu_{i2},$$

with

$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix},$$

where all regressors follow a standard normal distribution and  $N = 5000$ . Use  $\beta_0 = \mathbf{0}$ ,  $B_0 = 1000\mathbf{I}$ ,  $\alpha_0 = 4$ , and  $\Psi_0 = 4\mathbf{I}_2$ . Set the number of iterations to 2,000 and the thinning parameter to 5.

- Perform inference under the setting of Section 7.4, assuming that  $x_{i3}$  may affect  $y_{i2}$ .
- Program a Gibbs sampling algorithm that accounts for different regressors in each binary decision, that is, where  $x_{i3}$  does not affect  $y_{i2}$ .



# 8

## Time series models

In this chapter, we provide a brief introduction to performing inference in time series models using a Bayesian framework. There is a large literature in time series in statistics and econometrics, and it would be impossible to present a good treatment in a few pages of an introductory book. However, there are excellent books in Bayesian inference in time series, see for instance, [450, 336, 340].

A time series is a sequence of observations collected in chronological order, allowing us to track how variables change over time. However, it also introduces technical challenges, as we must account for statistical features such as autocorrelation and stationarity. Since time series data is time-dependent, we adjust our notation. Specifically, we use  $t$  and  $T$  instead of  $i$  and  $N$  to explicitly indicate time.

Our starting point in this chapter is the *state-space representation* of time series models. Much of the Bayesian inference literature in time series adopts this approach, as it allows dynamic systems to be modeled in a structured way. This representation provides modularity, flexibility, efficiency, and interpretability in complex models where the state evolves over time. It also enables the use of recursive estimation methods, such as the *Kalman filter* for dynamic Gaussian linear models and the *particle filter* (also known as *sequential Monte Carlo*) for non-Gaussian and nonlinear state-space models. The latter method is especially useful for *online* predictions or when there are data storage limitations. These inferential tools are based on the sequential updating process of Bayes' rule, where the posterior at time  $t$  becomes the prior at time  $t + 1$  (see Equation 1.14).

Remember that we can run our GUI typing

**R code. How to display our graphical user interface**

```
1 shiny::runGitHub("besmarter/BSTApp", launch.browser = T)
```

in the **R** package console or any **R** code editor, and once our GUI is

deployed, select *Time series Models*. However, users should see Chapter 5 for details.

## 8.1 State-space representation

A *state-space model* is composed by of an *unobservable state vector*  $\beta_t \in \mathbb{R}^K$ , and an *observed measure*  $\mathbf{Y}_t \in \mathbb{R}^M$ ,  $t = 1, 2, \dots$  such that (i)  $\beta_t$  is a *Markov process*, this is,  $\pi(\beta_t | \beta_{1:t-1}) = \pi(\beta_t | \beta_{t-1})$ , all the information regarding  $\beta_t$  based on all its history up to  $t-1$  is carried by  $\beta_{t-1}$ , and (ii)  $\mathbf{Y}_t$  is independent of  $\mathbf{Y}_s$  conditional on  $\beta_t$ ,  $s < t$  [336, Chap. 2].

These assumptions imply that  $\pi(\beta_{0:t}, \mathbf{Y}_{1:t}) = \pi(\beta_0) \prod_{s=1}^t \pi(\beta_s | \beta_{s-1}) \pi(\mathbf{Y}_s | \beta_s)$ .<sup>1</sup>

There are three key aims in *state-space models*: *filtering*, *smoothing*, and *forecasting*. In *filtering*, we aim to estimate the current state given observations up to time  $t$ , specifically obtaining the density  $\pi(\beta_s | \mathbf{y}_{1:t})$  for  $s = t$ . In *smoothing*, we conduct a retrospective analysis of the system, obtaining  $\pi(\beta_s | \mathbf{y}_{1:t})$  for  $s < t$ . In *forecasting*, we forecast future observations by first obtaining  $\pi(\beta_s | \mathbf{y}_{1:t})$  as an intermediate step to compute  $\pi(\mathbf{Y}_s | \mathbf{y}_{1:t})$  for  $s > t$ . A valuable feature of these methods is that all these densities can be calculated recursively. [336] show the recursive equations in Propositions 2.1 (filtering), 2.3 (smoothing) and 2.5 (forecasting).

An important class of *state-space models* is the *Gaussian linear state-space model*, also known as, *dynamic linear model*:

$$\begin{aligned} \mathbf{Y}_t &= \mathbf{X}_t \beta_t + \boldsymbol{\mu}_t && \text{(Observation equations)} \\ \beta_t &= \mathbf{G}_t \beta_{t-1} + \mathbf{w}_t && \text{(States equations)}, \end{aligned}$$

where  $\beta_0 \sim N(\mathbf{b}_0, \mathbf{B}_0)$ ,  $\boldsymbol{\mu}_t \sim N(\mathbf{0}, \boldsymbol{\Sigma}_t)$ ,  $\mathbf{w}_t \sim N(\mathbf{0}, \boldsymbol{\Omega}_t)$ ,  $\beta_0$ ,  $\boldsymbol{\mu}_t$  and  $\mathbf{w}_t$  are independent,  $\mathbf{X}_t$  and  $\mathbf{G}_t$  are  $M \times K$  and  $K \times K$  known matrices. Observe that this assumption implies that  $\mathbf{Y}_t | \beta_t \sim N(\mathbf{X}_t \beta_t, \boldsymbol{\Sigma}_t)$ , and  $\beta_t | \beta_{t-1} \sim N(\mathbf{G}_t \beta_{t-1}, \boldsymbol{\Omega}_t)$ .<sup>2</sup>

Let  $\beta_{t-1} | \mathbf{y}_{1:t-1} \sim N(\mathbf{b}_{t-1}, \mathbf{B}_{t-1})$ , then, we can get the *Kalman filter* by getting

1. The one-step-ahead predictive distribution of  $\beta_t$  given  $\mathbf{y}_{1:t-1}$  is  $\beta_t | \mathbf{y}_{1:t-1} \sim N(\mathbf{a}_t, \mathbf{R}_t)$ , where  $\mathbf{a}_t = \mathbf{G}_t \mathbf{b}_{t-1}$  and  $\mathbf{R}_t = \mathbf{G}_t \mathbf{B}_{t-1} \mathbf{G}_t^\top + \boldsymbol{\Omega}_t$ .
2. The one-step-ahead predictive distribution of  $\mathbf{Y}_t$  given  $\mathbf{y}_{1:t-1}$  is  $\mathbf{Y}_t | \mathbf{y}_{1:t-1} \sim N(\mathbf{f}_t, \mathbf{Q}_t)$ , where  $\mathbf{f}_t = \mathbf{X}_t \mathbf{a}_t$  and  $\mathbf{Q}_t = \mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \boldsymbol{\Sigma}_t$ .

<sup>1</sup>A *state-space model* where the states are random variables taking discrete values is called *hidden Markov model*.

<sup>2</sup>A general *state-space model* is given by  $\mathbf{Y}_t = \mathbf{f}_t(\beta_t, \boldsymbol{\mu}_t)$ , and  $\beta_t = \mathbf{m}_t(\beta_{t-1}, \mathbf{w}_t)$  for arbitrary functions  $\mathbf{f}_t$  and  $\mathbf{m}_t$ , and distributions for  $\boldsymbol{\mu}_t$  and  $\mathbf{w}_t$ , and a prior distribution for  $\beta_0$ .

3. The distribution of the one-step-ahead prediction error  $\mathbf{e}_t = \mathbf{Y}_t - \mathbb{E}[\mathbf{Y}_t | \mathbf{y}_{1:t-1}] = \mathbf{Y}_t - \mathbf{f}_t$  is  $N(\mathbf{0}, \mathbf{Q}_t)$  [410, Chap. 6].
4. The filtering distribution of  $\beta_t$  given  $\mathbf{y}_{1:t}$  is  $\beta_t | \mathbf{y}_{1:t} \sim N(\mathbf{b}_t, \mathbf{B}_t)$ , where  $\mathbf{b}_t = \mathbf{a}_t + \mathbf{K}_t \mathbf{e}_t$ ,  $\mathbf{K}_t = \mathbf{R}_t \mathbf{X}_t^\top \mathbf{Q}_t^{-1}$  is the *Kalman gain*, and  $\mathbf{B}_t = \mathbf{R}_t - \mathbf{R}_t \mathbf{X}_t^\top \mathbf{Q}_t^{-1} \mathbf{X}_t \mathbf{R}_t$ .

The formal proofs of these results can be found in [336, Chap 2]. Just take into account that the logic of these results follow the results of the Seemingly unrelated regression (SUR) model in Section 7.2 for a particular time period. In addition, we know that the posterior distribution using information up to  $t-1$  becomes the prior in  $t$  (see Equation 1.14,  $\pi(\boldsymbol{\theta} | \mathbf{y}_{1:t}) \propto p(y_t | \mathbf{y}_{1:t-1}, \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta} | \mathbf{y}_{1:t-1})$ ). This is the updating process from  $\beta_t | \mathbf{y}_{1:t-1} \sim N(\mathbf{a}_t, \mathbf{R}_t)$  to  $\beta_t | \mathbf{y}_{1:t} \sim N(\mathbf{b}_t, \mathbf{B}_t)$ . Moreover, the posterior mean and variance of the SUR model with independent conjugate priors for a particular time period can be written as  $\mathbf{a}_t + \mathbf{R}_t \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \Sigma_t)^{-1} (\mathbf{y}_t - \mathbf{X}_t \mathbf{a}_t)$  and  $\mathbf{R}_t - \mathbf{R}_t \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \Sigma_t)^{-1} \mathbf{X}_t \mathbf{R}_t^\top$ , respectively. Let's see this, we know from Section 7.2 that  $\mathbf{B}_t = (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1}$  and  $\beta_t = \mathbf{B}_t (\mathbf{R}_t^{-1} \mathbf{a}_t + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{y}_t)$ . Thus, let's show that both conditional posterior distributions are the same. In particular, the posterior mean in the *state-space representation* is  $[\mathbf{I}_K - \mathbf{R}_t \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \Sigma_t)^{-1} \mathbf{X}_t] \mathbf{a}_t + \mathbf{R}_t \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \Sigma_t)^{-1} \mathbf{y}_t$ , where

$$\begin{aligned} \mathbf{R}_t \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \Sigma_t)^{-1} &= \mathbf{R}_t \mathbf{X}_t^\top [\Sigma_t^{-1} - \Sigma_t^{-1} \mathbf{X}_t (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{X}_t^\top \Sigma_t^{-1}] \\ &= \mathbf{R}_t [\mathbf{I}_K - \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1}] \mathbf{X}_t^\top \Sigma_t^{-1} \\ &= \mathbf{R}_t (\mathbf{I}_K - [\mathbf{I}_K - \mathbf{R}_t^{-1} (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1}]) \mathbf{X}_t^\top \Sigma_t^{-1} \\ &= (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{X}_t^\top \Sigma_t^{-1}, \end{aligned}$$

where the first equality uses the Woodbury matrix identity (matrix inversion lemma), and the third equality uses  $\mathbf{D}(\mathbf{D} + \mathbf{E})^{-1} = \mathbf{I} - \mathbf{E}(\mathbf{D} + \mathbf{E})^{-1}$ .

Thus, we have the following expression:

$$\begin{aligned} &[\mathbf{I}_K - \mathbf{R}_t \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \Sigma_t)^{-1} \mathbf{X}_t] \mathbf{a}_t + \mathbf{R}_t \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \Sigma_t)^{-1} \mathbf{y}_t \\ &= [\mathbf{I}_K - (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t] \mathbf{a}_t + (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{y}_t \\ &= (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{R}_t^{-1} \mathbf{a}_t + (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{y}_t \\ &= (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} (\mathbf{R}_t^{-1} \mathbf{a}_t + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{y}_t) \\ &= (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} (\mathbf{R}_t^{-1} \mathbf{a}_t + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t \hat{\beta}_t), \end{aligned}$$

where the second equality uses the identity:

$$\mathbf{I} - (\mathbf{D} + \mathbf{E})^{-1} \mathbf{D} = (\mathbf{D} + \mathbf{E})^{-1} \mathbf{E},$$

and the estimator  $\hat{\beta}_t$  is defined as:

$$\hat{\beta}_t = (\mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{y}_t.$$

This shows that the posterior mean is a weighted average of the prior mean

and the maximum likelihood estimator (which is the generalized least squares estimator).

The weights are linked to the signal-to-noise ratio, that is, the proportion of the total variability ( $\Omega_t + \Sigma_t$ ) due to the signal ( $\Omega_t$ ) versus the noise ( $\Sigma_t$ ). Note that in the simplest case where  $M = K = 1$ , and  $\mathbf{X}_t = \mathbf{G}_t = 1$ , then  $\mathbf{K}_t = \mathbf{R}_t \mathbf{Q}_t^{-1} = (B_{t-1} + \Omega_t)/(B_{t-1} + \Omega_t + \Sigma_t)$ . Thus, the weight associated with the observations is equal to 1 if  $\Sigma_t = 0$ , that is, the posterior mean is equal to the actual observation. On the other hand, if  $\Sigma_t$  increases compare to  $\Omega_t$ , there is more weight to the prior information, and consequently, the posterior mean is smoother as it heavily dependents on the history. We ask in Exercise 1 to perform simulations with different signal-to-noise ratios to see the effects on the system.

The equality of variances of both approaches is as follows:

$$\begin{aligned} \text{Var}[\beta_t | \mathbf{y}_{1:t}] &= \mathbf{R}_t - \mathbf{R}_t \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{R}_t \mathbf{X}_t^\top + \Sigma_t)^{-1} \mathbf{X}_t \mathbf{R}_t \\ &= \mathbf{R}_t - \mathbf{R}_t \mathbf{X}_t^\top (\Sigma_t^{-1} - \Sigma_t^{-1} \mathbf{X}_t (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{X}_t^\top \Sigma_t^{-1}) \mathbf{X}_t \mathbf{R}_t \\ &= \mathbf{R}_t - \mathbf{R}_t \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t \mathbf{R}_t \\ &\quad + \mathbf{R}_t \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t \mathbf{R}_t \\ &= \mathbf{R}_t - \mathbf{R}_t \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t \mathbf{R}_t \\ &\quad + \mathbf{R}_t \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t [\mathbf{I}_K - (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \mathbf{R}_t^{-1}] \mathbf{R}_t \\ &= \mathbf{R}_t - \mathbf{R}_t \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1} \\ &= \mathbf{R}_t [\mathbf{I}_K - \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1}] \\ &= \mathbf{R}_t [\mathbf{I}_K - (\mathbf{I}_K - \mathbf{R}_t^{-1} (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1})] \\ &= (\mathbf{R}_t^{-1} + \mathbf{X}_t^\top \Sigma_t^{-1} \mathbf{X}_t)^{-1}, \end{aligned}$$

where the second equality uses the Woodbury matrix identity, the fourth equality uses  $(\mathbf{D} + \mathbf{E})^{-1} \mathbf{D} = \mathbf{I} - (\mathbf{D} + \mathbf{E})^{-1} \mathbf{E}$ , and the seventh equality uses  $\mathbf{D}(\mathbf{D} + \mathbf{E})^{-1} = \mathbf{I} - \mathbf{E}(\mathbf{D} + \mathbf{E})^{-1}$ .

The *Kalman filter* allows calculating recursively in a forward way  $\pi(\beta_t | \mathbf{y}_{1:t})$  from  $\pi(\beta_{t-1} | \mathbf{y}_{1:t-1})$  starting from  $\pi(\beta_0)$ .

Let  $\beta_{t+1} | \mathbf{y}_{1:T} \sim N(\mathbf{s}_{t+1}, \mathbf{S}_{t+1})$ , then we can get the *Kalman smoother* by  $\beta_t | \mathbf{y}_{1:T} \sim N(\mathbf{s}_t, \mathbf{S}_t)$ , where  $\mathbf{s}_t = \mathbf{b}_t + \mathbf{B}_t \mathbf{G}_{t+1}^\top \mathbf{R}_{t+1}^{-1} (\mathbf{s}_{t+1} - \mathbf{a}_{t+1})$  and  $\mathbf{S}_t = \mathbf{B}_t - \mathbf{B}_t \mathbf{G}_{t+1}^\top \mathbf{R}_{t+1}^{-1} (\mathbf{R}_{t+1} - \mathbf{S}_{t+1}) \mathbf{R}_{t+1}^{-1} \mathbf{G}_{t+1} \mathbf{B}_t$ . The proof can be found in [336, Chap 2].

Thus, we can calculate the *Kalman smoother* starting from  $t = T - 1$ , that is,  $\beta_T | \mathbf{y}_{1:T} \sim N(\mathbf{s}_T, \mathbf{S}_T)$ . However, this is the filtering distribution at  $T$ , which means  $\mathbf{s}_T = \mathbf{b}_T$  and  $\mathbf{S}_T = \mathbf{B}_T$ , and then, we should proceed recursively in a backward way.

Finally, the forecasting recursion in the *dynamic linear model*, given  $\mathbf{a}_t(0) = \mathbf{b}_t$  and  $\mathbf{R}_t(0) = \mathbf{B}_t$ ,  $h \geq 1$ , is given by

1. The forecasting distribution of  $\beta_{t+h} | \mathbf{y}_{1:t}$  is  $N(\mathbf{a}_t(h), \mathbf{R}_t(h))$ , where  $\mathbf{a}_t(h) = \mathbf{G}_{t+h} \mathbf{a}_t(h-1)$  and  $\mathbf{R}_t(h) = \mathbf{G}_{t+h} \mathbf{R}_t(h-1) \mathbf{G}_{t+h}^\top + \Omega_{t+h}$ .

2. The forecasting distribution  $\mathbf{Y}_{t+h} \mid \mathbf{y}_{1:t}$  is  $N(\mathbf{f}_t(h), \mathbf{Q}_t(h))$ , where  $\mathbf{f}_t(h) = \mathbf{X}_{t+h}\mathbf{a}_t(h)$  and  $\mathbf{Q}_t(h) = \mathbf{X}_{t+h}\mathbf{R}_t(h)\mathbf{X}_{t+h}^\top + \boldsymbol{\Sigma}_{t+h}$ .

The proof can be found in [336, Chap 2].

These recursive equations allow us to perform probabilistic forecasting  $h$ -steps-ahead for the state and observation equations.

These results demonstrate how to use these recursive equations for filtering, smoothing, and forecasting in *dynamic linear models* (*Gaussian linear state-space models*). Although these algorithms appear simple, they suffer from numerical instability, which can lead to non-symmetric and negative-definite covariance matrices. Thus, special care must be taken when working with them.

In addition, this setup assumes that  $\boldsymbol{\Sigma}_t$  and  $\boldsymbol{\Omega}_t$  are known. However, this is rarely the case in most situations. Therefore, we need to estimate them. One option is to perform maximum likelihood estimation. However, this approach does not account for the uncertainty associated with the fact that  $\boldsymbol{\Sigma}_t$  and  $\boldsymbol{\Omega}_t$  are unknown when their estimates are *plugged into* the *state space* recursions. On the other hand, we can use a Bayesian approach and perform the recursions associated with each posterior draw of the unknown parameters, thus taking their uncertainty into account.

The point of departure is the posterior distribution, such that

$$\pi(\boldsymbol{\theta}, \boldsymbol{\beta}_0, \dots, \boldsymbol{\beta}_T \mid \mathbf{y}, \mathbf{X}, \mathbf{G}) \propto \pi(\boldsymbol{\beta}_0 \mid \boldsymbol{\theta})\pi(\boldsymbol{\theta}) \prod_{t=1}^T \pi(\boldsymbol{\beta}_t \mid \boldsymbol{\beta}_{t-1}, \boldsymbol{\theta})\pi(\mathbf{y}_t \mid \boldsymbol{\beta}_t, \boldsymbol{\theta}),$$

where  $\boldsymbol{\theta}$  is the vector of unknown parameters.

We can compute  $\pi(\boldsymbol{\beta}_s, \boldsymbol{\theta} \mid \mathbf{y}_{1:t}) = \pi(\boldsymbol{\beta}_s \mid \mathbf{y}_{1:t}, \boldsymbol{\theta})\pi(\boldsymbol{\theta} \mid \mathbf{y}_{1:t})$ , for  $s = t$  (*filtering*),  $s < t$  (*smoothing*), and  $s > t$  (*forecasting*). The marginal posterior distribution of the states is  $\pi(\boldsymbol{\beta}_s \mid \mathbf{y}_{1:t}) = \int_{\boldsymbol{\Theta}} \pi(\boldsymbol{\beta}_s \mid \mathbf{y}_{1:t}, \boldsymbol{\theta})\pi(\boldsymbol{\theta} \mid \mathbf{y}_{1:t})d\boldsymbol{\theta}$ .

We can use the Gibbs sampling algorithm to get the posterior draws in the *dynamic linear model* assuming conjugate families. In particular, let's see the univariate case with *random walk states*,

$$y_t = \mathbf{x}_t^\top \boldsymbol{\beta}_t + \mu_t \tag{8.1}$$

$$\boldsymbol{\beta}_t = \boldsymbol{\beta}_{t-1} + \mathbf{w}_t, \tag{8.2}$$

where  $\mu_t \sim N(0, \sigma^2)$  and  $\mathbf{w}_t \sim N(\mathbf{0}, \text{diag}\{\omega_1^2, \dots, \omega_K^2\})$ . We assume that  $\pi(\sigma^2, \omega_1^2, \dots, \omega_K^2, \boldsymbol{\beta}_0) = \pi(\sigma^2)\pi(\omega_1^2), \dots, \pi(\omega_K^2)\pi(\boldsymbol{\beta}_0)$  where  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ ,  $\omega_k^2 \sim IG(\alpha_{k0}/2, \delta_{k0}/2)$ ,  $k = 1, \dots, K$ , and  $\boldsymbol{\beta}_0 \sim N(\mathbf{b}_0, \mathbf{B}_0)$ . Thus, the conditional posterior distributions are  $\sigma^2 \mid \mathbf{y}, \mathbf{X}, \boldsymbol{\beta}_{1:T} \sim IG(\alpha_n/2, \delta_n/2)$ , where  $\alpha_n = T + \alpha_0$  and  $\delta_n = \sum_{t=1}^T (y_t - \mathbf{x}_t^\top \boldsymbol{\beta}_t)^2 + \delta_0$ , and  $\omega_k^2 \mid \mathbf{y}, \mathbf{X}, \boldsymbol{\beta}_{0:T} \sim IG(\alpha_{kn}/2, \delta_{kn}/2)$ , where  $\alpha_{kn} = T + \alpha_{k0}$  and  $\delta_{kn} = \sum_{t=1}^T (\boldsymbol{\beta}_{t,k} - \boldsymbol{\beta}_{t-1,k})^2 + \delta_{k0}$ . The vector of the dependent variable is  $\mathbf{y}$ , and all regressors are in  $\mathbf{X}$ .

We also need to sample the states from  $\pi(\boldsymbol{\beta}_{1:T} \mid \mathbf{y}, \mathbf{X}, \sigma^2, \omega_1^2, \dots, \omega_K^2)$ . This can be done using the *forward filtering backward sampling* (FFBS) algorithm

[71, 149, 409]. This algorithm is basically a simulation version of the *smoothing* recursion, which allows getting draws of the states, even if we do not have analytical solutions, for instance, in non-linear settings. See below and [336, Chap. 3] for details. A word of caution here, users should be careful to set non-informative priors in this setting, and in general, settings where there are a large number of parameters (see [250, Chap. 8] for details). Thus, it is useful to use empirical Bayes methods focusing on relevant hyperparameters, for instance, the hyperparameters of the inverse-gamma distributions which define the signal-to-noise ratio.

We use the command *dlmGibbsDIG* from the *dlm* package in our GUI to perform Bayesian inference in the univariate *dynamic linear model* with *random walk states*. This function uses the FFBS algorithm, and assumes independent gamma priors for the precision (inverse of variance) parameters. In addition, this package uses the singular value decomposition to calculate the covariance matrices to avoid numerical instability.

Algorithm A21 shows how to perform inference in the univariate *dynamic linear model* with random walk states in our GUI. See also Chapter 5 for details regarding the dataset structure.

---

**Algorithm A21** Dynamic linear models

---

- 1: Select *Time series Model* on the top panel
  - 2: Select *Dynamic linear model* using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Set the hyperparameters of the *precision of the observation equation*: prior mean and variance.
  - 6: Set the hyperparameters of the *precision of the state equations*: just one set of prior mean and variance parameters.
  - 7: Click the *Go!* button
  - 8: Analyze results
  - 9: Download posterior chains of variances of observation and state equations, and posterior chains of states using the *Download Results* button
- 

**Example: Simulation exercise of the dynamic linear model**

We simulate the process  $y_t = \beta_{t1} + x_t\beta_{t2} + \mu_t$  and  $\beta_t = \beta_{t-1} + \mathbf{w}_t$ ,  $t = 1, 2, \dots, 200$ , where  $\beta_t = [\beta_{t1} \ \beta_{t2}]^\top$ ,  $\mu_t \sim N(0, 0.5^2)$ ,  $\mathbf{w}_t \sim N(\mathbf{0}, \text{diag}\{0.2, 0.1\})$ ,  $x_t \sim N(1, 1)$ ,  $\beta_0$  and  $\mathbf{B}_0$  are the OLS estimates and variance of the recursive OLS estimates (see below), respectively.

The following algorithm demonstrates how to perform inference using *dlmGibbsDIG* and compares the results to those of the maximum likelihood estimator, which is based on the *dlmMLE* function. We also use the *dlmSvd2var* function, which is based on singular value decomposition, to calculate the vari-

ance of the smoothing states. All these functions are from the *dlm* package in **R**.

Users can observe that we employ a straightforward strategy for setting the hyperparameters. First, we recursively estimate the model using ordinary least squares (OLS), progressively increasing the sample size, and save the location parameters. Next, we compute the covariance matrix of this sequence and use it to set the priors: the prior mean of the precision of the state vector is set equal to the inverse of the maximum element of the main diagonal of this covariance matrix (*a.theta*), and the prior variance is set equal to ten times this value (*b.theta*). For the observation equation, the prior mean of the precision is set equal to the inverse of the OLS variance estimate (*a.y*), and the prior variance is set equal to ten times this value (*b.y*). We perform some sensitivity analysis of the results regarding the hyperparameters, and it seems that the results are robust. However, we encourage giving more consideration to empirical Bayes methods for setting hyperparameters in *state-space models*.

### **R code. Simulation: Dynamic linear model**

```

1 rm(list = ls()); set.seed(010101)
2 T <- 200; sig2 <- 0.5^2
3 x <- rnorm(T, mean = 1, sd = 1); X <- cbind(1, x); B0 <- c
  (1, 0.5)
4 K <- length(B0); e <- rnorm(T, mean = 0, sd = sig2^0.5)
5 Omega <- diag(c(0.2, 0.1))
6 w <- MASS::mvrnorm(T, c(0, 0), Omega)
7 Bt <- matrix(NA, T, K); Bt[1,] <- B0
8 yt <- rep(NA, T)
9 yt[1] <- X[1,] %*% B0 + e[1]
10 for(t in 1:T){
11   if(t == 1){
12     Bt[t,] <- w[t,]
13   }else{
14     Bt[t,] <- Bt[t-1,] + w[t,]
15   }
16   yt[t] <- X[t,] %*% Bt[t,] + e[t]
17 }
18 RegLS <- lm(yt ~ x)
19 SumRegLS <- summary(RegLS)
20 SumRegLS; SumRegLS$sigma^2
21 Bp <- matrix(RegLS$coefficients, T, K, byrow = TRUE)
22 S <- 20
23 for(t in S:T){
24   RegLSt <- lm(yt[1:t] ~ x[1:t])
25   Bp[t,] <- RegLSt$coefficients
26 }
```

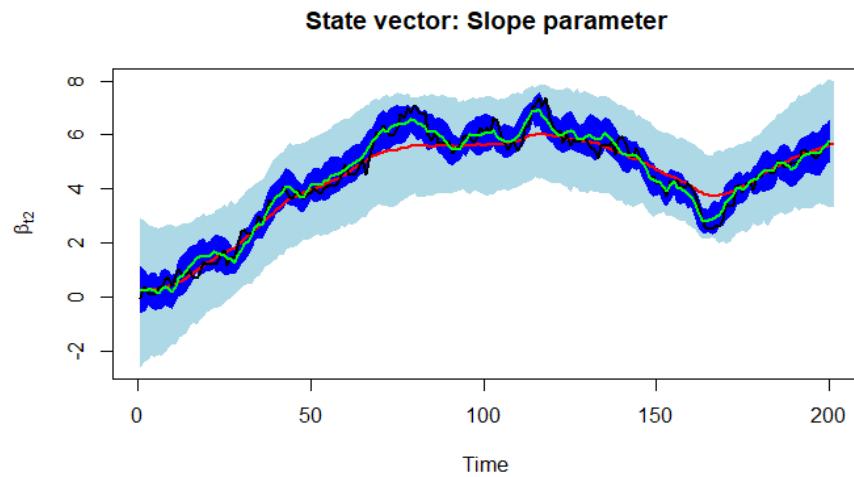
### R code. Simulation: Dynamic linear model

```

1 VarBp <- var(Bp)
2 # State space model
3 ModelReg <- function(par){
4   Mod <- dlm::dlmModReg(x, dV = exp(par[1]), dW = exp(par
5     [2:3]), m0 = RegLS$coefficients, CO = VarBp)
6   return(Mod)
7 }
8 outMLEReg <- dlm::dlmMLE(yt, parm = rep(0, K+1), ModelReg)
9 exp(outMLEReg$par)
10 RegFilter <- dlm::dlmFilter(yt, ModelReg(outMLEReg$par))
11 RegSmooth <- dlm::dlmSmooth(yt, ModelReg(outMLEReg$par))
12 SmoothB2 <- RegSmooth$s[-1,2]
13 VarSmooth <- dlm::dlmSvd2var(u = RegSmooth[["U.S"]], RegSmooth
14   [ ["D.S"]])
15 SDVarSmoothB2 <- sapply(2:(T+1), function(t){VarSmooth[[t]][
16   K,K]^0.5})
17 LimInfB2 <- SmoothB2 - qnorm(0.975)*SDVarSmoothB2
18 LimSupB2 <- SmoothB2 + qnorm(0.975)*SDVarSmoothB2
19 # Gibbs
20 MCMC <- 2000; burnin <- 1000
21 a.y <- (SumRegLS$sigma^2)^(-1); b.y <- 10*a.y; a.theta <- (
22   max(diag(VarBp)))^(-1); b.theta <- 10*a.theta
23 gibbsOut <- dlm::dlmGibbsDIG(yt, mod = dlm::dlmModReg(x), a.
24   y = a.y, b.y = b.y, a.theta = a.theta, b.theta = b.theta
25   , n.sample = MCMC, thin = 5, save.states = TRUE)
26 B2t <- matrix(0, MCMC - burnin, T + 1)
27 for(t in 1:(T+1)){
28   B2t[,t] <- gibbsOut[["theta"]][t,2,-c(1:burnin)]
29 }
30 Lims <- apply(B2t, 2, function(x){quantile(x, c(0.025,
31   0.975))})
32 summary(coda::mcmc(gibbsOut[["dV"]]))
33 summary(coda::mcmc(gibbsOut[["dW"]]))
34 # Figure
35 require(latex2exp) # LaTeX equations in figures
36 xx <- c(1:(T+1), (T+1):1)
37 yy <- c(Lims[1,], rev(Lims[2,]))
38 plot(xx, yy, type = "n", xlab = "Time", ylab = TeX("$\\beta_1$"))
39 polygon(xx, yy, col = "lightblue", border = "lightblue")
40 xxML <- c(1:T, T:1)
41 yyML <- c(LimInfB2, rev(LimSupB2))
42 polygon(xxML, yyML, col = "blue", border = "blue")
43 lines(colMeans(B2t), col = "red", lw = 2)
44 lines(Bt[,2], col = "black", lw = 2)
45 lines(SmoothB2, col = "green", lw = 2)
46 title("State vector: Slope parameter")

```

Figure 8.1 shows the comparison between maximum likelihood (ML) and Bayesian inference. The light blue (Bayesian) and dark blue (maximum likelihood) shadows show the credible and confidence intervals at 95% for the state slope parameter ( $\beta_{t2}$ ). We see that the Bayesian interval encompass the ML interval. This is a reflection of the extra uncertainty of the unknown variances. The black line is the actual trajectory of  $\beta_{t2}$ , the green and red lines are the *smoothing* recursions using the ML and Bayesian estimates (posterior mean), respectively.



**FIGURE 8.1**  
Simulation: Dynamic linear model.

#### Example: Effects of inflation on interest rate I

We use the dataset *16INTDEF.csv* provided by [457, Chaps. 10] to study the effects of inflation on the interest rate. The specification is

$$\Delta i_t = \beta_{t1} + \beta_{t2} \Delta \text{inf}_t + \beta_{t3} \Delta \text{def}_t + \mu_t$$

and

$$\boldsymbol{\beta}_t = \boldsymbol{\beta}_{t-1} + \boldsymbol{w}_t,$$

where  $\Delta z_t = z_t - z_{t-1}$  is the difference operator,  $i_t$  is the three-month T-bill rate,  $\text{inf}_t$  is the annual inflation rate based on the consumer price index (CPI), and  $\text{def}_t$  is the federal budget deficit as a percentage of gross domestic product (GDP) from 1948 to 2003 in the USA. In addition,  $\mu_t \sim N(0, \sigma^2)$  and  $\boldsymbol{w}_t \sim N(\mathbf{0}, \text{diag}\{\omega_1^2, \omega_2^2\})$ . We assume inverse-gamma distributions for the priors of the scale parameters and set 12,000 MCMC iterations, 2,000 as burn-in, and 10 as the thinning parameter.

The following code shows how to perform this application. We use the variance of the recursive estimation of OLS to set the hyperparameters of the inverse-gamma distribution for the variances of  $w_t$ , and the OLS estimate of the variance of the model to set the hyperparameters of the distribution of  $\sigma^2$ . Note that, as we are using the function *dlmGibbsDIG* from the *dlm* package, the hyperparameters are set in terms of precision parameters.

Figure 8.2 shows the posterior results of the effect of inflation on the interest rate. This is a fan chart indicating deciles from 10% to 90%. The red shaded area shows the range around the median value, and the black line represents the mean value of the state associated with the annual change in inflation. We see that the annual changes in interest rates are weakly positively related to annual changes in inflation.

### R code. Dynamic linear model: Effects of inflation on interest rate

```

1 rm(list = ls()); set.seed(10101)
2 DataIntRate <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/16INTDEF.csv"
  , sep = ",", header = TRUE, quote = "")
3 attach(DataIntRate); Xt <- cbind(diff(inf), diff(def))
4 K <- dim(Xt)[2] + 1; yt <- diff(i3)
5 T <- length(yt); RegLS <- lm(yt ~ Xt)
6 SumRegLS <- summary(RegLS); SumRegLS; SumRegLS$sigma^2
7 # Recursive OLS
8 Bp <- matrix(RegLS$coefficients, T, K, byrow = TRUE)
9 S <- 20
10 for(t in S:T){
11   RegLSt <- lm(yt[1:t] ~ Xt[1:t,])
12   Bp[t,] <- RegLSt$coefficients
13 }
14 VarBp <- var(Bp)
15 # State space model
16 ModelReg <- function(par){
17   Mod <- dlm::dlmModReg(Xt, dV = exp(par[1]), dW = exp(par
  [2:(K+1)]), m0 = RegLS$coefficients,
18   C0 = diag(VarBp))
19   return(Mod)
20 }
21 MCMC <- 12000; burnin <- 2000; thin <- 10
22 a.y <- (SumRegLS$sigma^2)^(-1); b.y <- 10*a.y; a.theta <- (
  max(diag(VarBp)))^(-1); b.theta <- 10*a.theta
23 gibbsOut <- dlm::dlmGibbsDIG(yt, mod = dlm::dlmModReg(Xt), a.
  .y = a.y, b.y = b.y, a.theta = a.theta, b.theta = b.
  theta, n.sample = MCMC, thin = 5, save.states = TRUE)

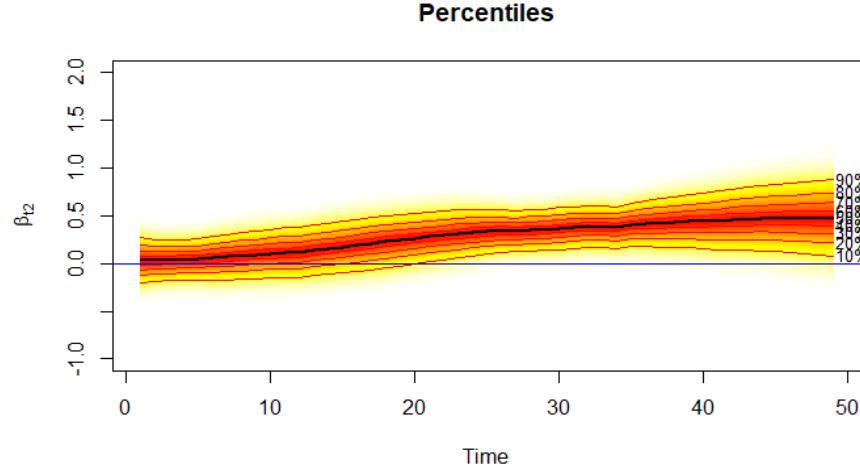
```

**R code. Dynamic linear model: Effects of inflation on interest rate**

```

1 B2t <- matrix(0, MCMC - burnin, T + 1)
2 for(t in 1:(T+1)){
3   B2t[,t] <- gibbsOut[["theta"]][t,2,-c(1:burnin)]
4 }
5 dV <- coda::mcmc(gibbsOut[["dV"]][-c(1:burnin)])
6 dW <- coda::mcmc(gibbsOut[["dW"]][-c(1:burnin),])
7 summary(dV); summary(dW)
8 plot(dV); plot(dW)
9 library(fanplot); library(latex2exp)
10 df <- as.data.frame(B2t)
11 plot(NULL, main="Percentiles", xlim = c(1, T+1), ylim = c
12   (-1, 2), xlab = "Time", ylab = TeX("$\\beta_{t1}$"))
13 fan(data = df); lines(colMeans(B2t), col = "black", lw = 2)
14 abline(h=0, col = "blue")

```



**FIGURE 8.2**  
Effects of inflation on interest rate: Dynamic linear model.

We can extend the *dynamic linear model* with *random walk states* to take into account time invariant location parameters. In particular, we follow [113], who propose the *simulation smoother*. This algorithm overcomes some shortcomings of the FFBS algorithm, such as slow convergence and computational

overhead. We focus on the case  $M = 1$ ,

$$y_t = \mathbf{z}_t^\top \boldsymbol{\alpha} + \mathbf{x}_t^\top \boldsymbol{\beta}_t + \mathbf{h}_t^\top \boldsymbol{\epsilon}_t, \quad t = 1, 2, \dots, T. \text{ (Observation equation)} \quad (8.3)$$

$$\boldsymbol{\beta}_t = \boldsymbol{\beta}_{t-1} + \mathbf{H}_t \boldsymbol{\epsilon}_t, \quad t = 1, 2, \dots, T. \text{ (States equations)}, \quad (8.4)$$

where  $\mathbf{z}_t$  and  $\mathbf{x}_t$  are  $L$ -dimensional and  $K$ -dimensional vectors of regressors associated with time-invariant and time-varying parameters, respectively,  $\mathbf{h}_t$  is a vector of dimension  $1 + K$ ,  $\mathbf{H}_t$  is a matrix of dimension  $K \times (1 + K)$ ,  $\boldsymbol{\beta}_0 = \mathbf{0}$  and  $\boldsymbol{\epsilon}_t \sim N(\mathbf{0}_{1+K}, \sigma^2 \mathbf{I}_{1+K})$ .

Observe that this specification encompasses Equations 8.1 and 8.2 setting  $\boldsymbol{\epsilon}_t = [\mu_t \ \mathbf{w}_t^\top]^\top$ ,  $\mathbf{h}_t = [1 \ 0 \ \dots \ 0]$ ,  $\mathbf{H}_t = [\mathbf{0}_K \ \mathbf{U}_{K \times K}]$  such that  $\text{diag}\{\omega_1^2 \ \dots \ \omega_K^2\} = \sigma^2 \mathbf{U} \mathbf{U}^\top$ ,  $\boldsymbol{\alpha} = \mathbf{0}$ , and  $\mathbf{h}_t \mathbf{H}_t^\top = \mathbf{0}_K$ .

The nice idea of [113] was to propose an efficient algorithm to get draws from  $\boldsymbol{\eta}_t = \mathbf{F}_t \boldsymbol{\epsilon}_t$ , where the most common choice is  $\mathbf{F}_t = \mathbf{H}_t$ , which means drawing samples from the perturbations of the states, and then, recovering the states from Equation 8.4 and  $\boldsymbol{\beta}_0 = \mathbf{0}$ . [113] present a more general version of the *state space model* that the one presented here.

Using the system given by Equations 8.3 and 8.4,  $\mathbf{F}_t = \mathbf{H}_t$  and  $\mathbf{h}_t \mathbf{H}_t^\top = \mathbf{0}_K$ , the *filtering* recursions are given by  $e_t = Y_t - \mathbf{z}_t^\top \boldsymbol{\alpha} - \mathbf{x}_t^\top \mathbf{b}_{t-1}$ ,  $q_t = \mathbf{x}_t^\top \mathbf{B}_{t-1} \mathbf{x}_t + \mathbf{h}_t^\top \mathbf{h}_t$ ,  $\mathbf{K}_t = \mathbf{B}_{t-1} \mathbf{x}_t q_t^{-1}$ ,  $\mathbf{b}_t = \mathbf{b}_{t-1} + \mathbf{K}_t e_t$ , and  $\mathbf{B}_t = \mathbf{B}_{t-1} - \mathbf{B}_{t-1} \mathbf{x}_t \mathbf{K}_t^\top + \mathbf{H}_t \mathbf{H}_t^\top$ , where  $\mathbf{b}_0 = \mathbf{0}$  and  $\mathbf{B}_0 = \mathbf{H}_0 \mathbf{H}_0^\top$ . See system 2 in [113] for a more general case. We should save  $e_t$  (innovation vector),  $q_t$  (scale innovation variance) and  $\mathbf{K}_t$  (*Kalman gain*) from this recursion.

Then, setting  $\mathbf{r}_T = 0$  and  $\mathbf{M}_T = \mathbf{0}$ , we run backwards from  $t = T-1, T-2, \dots, 1$ , the following recursions:  $\boldsymbol{\Lambda}_{t+1} = \mathbf{H}_{t+1} \mathbf{H}_{t+1}^\top$ ,  $\mathbf{C}_{t+1} = \boldsymbol{\Lambda}_{t+1} - \boldsymbol{\Lambda}_{t+1} \mathbf{M}_{t+1} \boldsymbol{\Lambda}_{t+1}^\top$ ,  $\boldsymbol{\xi}_{t+1} \sim N(\mathbf{0}_K, \sigma^2 \mathbf{C}_{t+1})$ ,  $\mathbf{L}_{t+1} = \mathbf{I}_K - \mathbf{K}_{t+1} \mathbf{x}_{t+1}^\top$ ,  $\mathbf{V}_{t+1} = \boldsymbol{\Lambda}_{t+1} \mathbf{M}_{t+1} \mathbf{L}_{t+1}$ ,  $\mathbf{r}_t = \mathbf{x}_{t+1} e_{t+1} / q_{t+1} + \mathbf{L}_{t+1}^\top \mathbf{r}_{t+1} - \mathbf{V}_{t+1}^\top \mathbf{C}_{t+1}^{-1} \boldsymbol{\xi}_{t+1}$ ,  $\mathbf{M}_t = \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top / q_{t+1} + \mathbf{L}_{t+1}^\top \mathbf{M}_{t+1} \mathbf{L}_{t+1} + \mathbf{V}_{t+1}^\top \mathbf{C}_{t+1}^{-1} \mathbf{V}_{t+1}$ , and  $\boldsymbol{\eta}_{t+1} = \boldsymbol{\Lambda}_{t+1} \mathbf{r}_{t+1} + \boldsymbol{\xi}_{t+1}$ . [113] show that  $\boldsymbol{\eta} = [\boldsymbol{\eta}_1^\top \ \dots \ \boldsymbol{\eta}_T^\top]^\top$  is drawn from  $p(\mathbf{H}_t \boldsymbol{\epsilon}_t \mid y_t, \mathbf{x}_t, \mathbf{z}_t, \mathbf{h}_t, \mathbf{H}_t, \boldsymbol{\alpha}, \sigma^2, t = 1, 2, \dots, T)$ . Thus, we can recover  $\boldsymbol{\beta}_t$  using 8.4 and  $\boldsymbol{\beta}_0 = \mathbf{0}_K$ .

We assume in the model given by Equations 8.3 and 8.4 that  $\mathbf{h}_t = [1 \ 0 \ \dots \ 0]^\top$  and  $\mathbf{H}_t = [\mathbf{0}_K \ \text{diag}\{1/\tau_1 \dots 1/\tau_K\}]$ , and then perform Bayesian inference assuming independent priors, that is,  $\pi(\boldsymbol{\beta}_0, \boldsymbol{\alpha}, \sigma^2, \boldsymbol{\tau}) = \pi(\boldsymbol{\beta}_0)\pi(\boldsymbol{\alpha})\pi(\sigma^2)\prod_{k=1}^K \pi(\tau_k^2)$  where  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ ,  $\tau_k^2 \sim G(v_0/2, v_0/2)$ ,  $k = 1, \dots, K$ ,  $\boldsymbol{\alpha} \sim N(\mathbf{a}_0, \mathbf{A}_0)$  and  $\boldsymbol{\beta}_0 \sim N(\mathbf{b}_0, \mathbf{B}_0)$ . The conditional posterior distributions are  $\sigma^2 \mid \mathbf{y}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\beta}_{0:T}, \boldsymbol{\alpha}, \boldsymbol{\tau} \sim IG(\alpha_n/2, \delta_n/2)$ , where  $\delta_n = \sum_{t=1}^T [(\boldsymbol{\beta}_t - \boldsymbol{\beta}_{t-1})^\top \boldsymbol{\Psi} (\boldsymbol{\beta}_t - \boldsymbol{\beta}_{t-1}) + (y_t - \mathbf{z}_t^\top \boldsymbol{\alpha} - \mathbf{x}_t^\top \boldsymbol{\beta}_t)^\top (y_t - \mathbf{z}_t^\top \boldsymbol{\alpha} - \mathbf{x}_t^\top \boldsymbol{\beta}_t)] + \delta_0$  and  $\alpha_n = T(K+1) + \alpha_0$ ,  $\boldsymbol{\tau} = [\tau_1 \ \dots \ \tau_K]$ ,  $\boldsymbol{\Psi} = \text{diag}\{\tau_1^2, \dots, \tau_K^2\}$ , and  $\tau_k^2 \mid \mathbf{y}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\beta}_{0:T}, \sigma^2 \sim G(v_{1n}/2, v_{2kn}/2)$ , where  $v_{1n} = T + v_0$  and  $v_{2kn} = \sigma^{-2} \sum_{t=1}^T (\boldsymbol{\beta}_{t,k} - \boldsymbol{\beta}_{t-1,k})^2 + v_0$ , and  $\boldsymbol{\alpha} \mid \mathbf{y}, \mathbf{X}, \mathbf{Z}, \sigma^2, \boldsymbol{\beta}_{1:T}, \boldsymbol{\tau} \sim N(\mathbf{a}_n, \mathbf{A}_n)$ , where  $\mathbf{A}_n = (\mathbf{A}_0^{-1} + \sigma^{-2} \sum_{t=1}^T \mathbf{z}_t \mathbf{z}_t^\top)^{-1}$  and  $\mathbf{a}_n = \mathbf{A}_n (\mathbf{A}_0^{-1} \mathbf{a}_0 + \sigma^{-2} \sum_{t=1}^T \mathbf{z}_t (y_t - \mathbf{x}_t^\top \boldsymbol{\beta}_t))$ . The vector of the dependent variable is  $\mathbf{y}$ , and all regressors are in  $\mathbf{X}$  and  $\mathbf{Z}$ .

We can see that all the previous posterior distributions are conditional on

the state vector  $\beta_{0:T}$ , which can be sampled using the *simulation smoother* algorithm, conditional on draws of the time-invariant parameters. Thus, the *state space model* provides an excellent illustration of the modular nature of the Bayesian framework, where performing inference on more complex models often simply involves adding new blocks to an MCMC algorithm. This means we can break down a complex inferential problem into smaller, more manageable parts, which is a “divide and conquer” approach. This is possible due to the structure of the conditional posterior distributions. Exercise 3 asks you to perform a simulation of the model given by Equations 8.3 and 8.4, and to program the MCMC algorithm, including the *simulation smoother*.

## 8.2 ARMA processes

Since the seminal work of [56], autoregressive moving average (ARMA) models have become ubiquitous in time series analysis. Thus, we present a brief introduction to these models in this section.

Let’s start with the linear Gaussian model with autoregressive errors,

$$y_t = \mathbf{x}_t^\top \boldsymbol{\beta} + \mu_t \quad (8.5)$$

$$\phi(L)\mu_t = \epsilon_t, \quad (8.6)$$

where  $\mathbf{x}_t$  is a  $K$ -dimensional vector of regressors,  $\epsilon_t \stackrel{iid}{\sim} N(0, \sigma^2)$ ,  $\phi(L) = 1 - \phi_1 L - \phi_2 L^2 - \dots - \phi_p L^p$  is a polynomial in the lag operator ( $L$ ), where  $Lz_t = z_{t-1}$ , and in general,  $L^r z_t = z_{t-r}$ .

Thus, we see that stochastic error  $\mu_t$  follows an *autoregressive process of order  $p$* , that is,  $\mu_t \sim AR(p)$ . It is standard practice to assume that  $\mu_t$  is second-order stationary, this implies that the mean, variance and autocovariance of  $\mu_t$  are finite and independent of  $t$  and  $s$ , although  $\mathbb{E}[\mu_t \mu_s]$  may depend on  $|t-s|$ . Then, all roots of  $\phi(L)$  lie outside the unit circle, for instance, given an  $AR(1)$ , then,  $1 - \phi_1 L = 0$ , implies,  $L = 1/\phi_1$  such that  $|\phi_1| < 1$  for the process being second-order stationary.

The likelihood function conditional on the first  $p$  observations is

$$\begin{aligned} p(y_{p+1}, \dots, y_T | y_p, \dots, y_1, \boldsymbol{\theta}) &= \prod_{t=p+1}^T p(y_t | \mathcal{I}_{t-1}, \boldsymbol{\theta}) \\ &\propto \sigma^{-(T-p)} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{t=p+1}^T (y_t - \hat{y}_{t|t-1, \boldsymbol{\theta}})^2 \right\}, \end{aligned}$$

where  $\mathcal{I}_{t-1}$  is the past information,  $\boldsymbol{\theta}$  collects all parameters  $(\boldsymbol{\beta}, \phi_1, \dots, \phi_p, \sigma^2)$ , and  $\hat{y}_{t|t-1, \boldsymbol{\theta}} = (1 - \phi(L))y_t + \phi(L)\mathbf{x}_t^\top \boldsymbol{\beta}$ .

We can see that multiplying the first expression in Equation 8.5 by  $\phi(L)$ , we can express the model as

$$y_t^* = \mathbf{x}_t^{*\top} \boldsymbol{\beta} + \epsilon_t \quad (8.7)$$

where  $y_t^* = \phi(L)y_t$  and  $\mathbf{x}_t^* = \phi(L)\mathbf{x}_t$ .

Thus, collecting all observations  $t = p+1, p+2, \dots, T$ , we have  $\mathbf{y}^* = \mathbf{X}^* \boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_{T-p})$ ,  $\mathbf{y}^*$  is a  $T-p$  dimensional vector, and  $\mathbf{X}^*$  is a  $(T-p) \times K$  dimensional matrix.

Assuming that  $\boldsymbol{\beta} \mid \sigma \sim N(\boldsymbol{\beta}_0, \sigma^2 \mathbf{B}_0)$ ,  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$  and  $\boldsymbol{\phi} \sim N(\boldsymbol{\phi}_0, \boldsymbol{\Phi}_0) \mathbb{1}(\boldsymbol{\phi} \in S_{\boldsymbol{\phi}})$ , where  $S_{\boldsymbol{\phi}}$  is the stationary region of  $\boldsymbol{\phi} = [\phi_1 \dots \phi_p]^\top$ . Then, Equation 8.7 implies that  $\boldsymbol{\beta} \mid \sigma^2, \boldsymbol{\phi}, \mathbf{y}, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \sigma^2 \mathbf{B}_n)$ , where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \mathbf{X}^{*\top} \mathbf{X}^*)^{-1}$  and  $\boldsymbol{\beta}_n = \mathbf{B}_n(\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{X}^{*\top} \mathbf{y}^*)$ . In addition,  $\sigma^2 \mid \boldsymbol{\beta}, \boldsymbol{\phi}, \mathbf{y}, \mathbf{X} \sim IG(\alpha_n/2, \delta_n/2)$  where  $\alpha_n = \alpha_0 + T - p$  and  $\delta_n = \delta_0 + (\mathbf{y}^* - \mathbf{X}^* \boldsymbol{\beta})^\top (\mathbf{y}^* - \mathbf{X}^* \boldsymbol{\beta}) + (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \mathbf{B}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)$ . Thus, the previous conditional posterior distributions imply that we can use a Gibbs sampling algorithm to perform inference of these parameters [88].

We know from Equation 8.5 that  $\mu_t = y_t - \mathbf{x}_t^\top \boldsymbol{\beta}$ , and from Equation 8.6 that  $\mu_t = \phi_1 \mu_{t-1} + \dots + \phi_p \mu_{t-p} + \epsilon_t$ ,  $t = p+1, \dots, T$ . In matrix notation  $\boldsymbol{\mu} = \mathbf{U} \boldsymbol{\phi} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\mu}$  is a  $T-p$  dimensional vector,  $\mathbf{U}$  is a  $(T-p) \times p$  matrix whose  $t$ -th row is  $[\mu_{t-1} \dots \mu_{t-p}]$ . Thus, the posterior distribution of  $\boldsymbol{\phi} \mid \boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X}$  is  $N(\boldsymbol{\phi}_n, \boldsymbol{\Phi}_n) \mathbb{1}(\boldsymbol{\phi} \in S_{\boldsymbol{\phi}})$ , where  $\boldsymbol{\Phi}_n = (\boldsymbol{\Phi}_0^{-1} + \sigma^{-2} \mathbf{U}^\top \mathbf{U})$  and  $\boldsymbol{\phi}_n = \boldsymbol{\Phi}_n(\boldsymbol{\Phi}_0^{-1} \boldsymbol{\phi}_0 + \sigma^{-2} \mathbf{U}^\top \boldsymbol{\mu})$  (see Exercise 4).

Drawing from the model under the stationarity restriction is straightforward: we simply sample from the multivariate normal distribution and discard draws that do not satisfy the stationarity condition. The proportion of draws that meet this restriction represents the conditional probability that the process is stationary.

### Example: Effects of inflation on interest rate II

We specify a *dynamic linear model* in the example of the effects of inflation on interest rates to account for a potential dynamic relationship. However, we can introduce dynamics in this model by assuming

$$\Delta i_t = \beta_1 + \beta_2 \Delta inf_t + \beta_3 \Delta def_t + \mu_t,$$

where  $\mu_t = \phi \mu_{t-1} + \epsilon_t$ . This leads to the model:

$$\Delta i_t = \beta_1(1-\phi_1) + \phi_1 \Delta i_{t-1} + \beta_2(\Delta inf_t - \phi_1 \Delta inf_{t-1}) + \beta_3(\Delta def_t - \phi_1 \Delta def_{t-1}) + \epsilon_t.$$

Thus, we again use the dataset *16INTDEF.csv* provided by [457, Chaps. 10] to illustrate linear regressions with *AR(1)* errors.

The following code demonstrates how to implement this application using vague priors, assuming  $\alpha_0 = \delta_0 = 0.01$ ,  $\boldsymbol{\beta}_0 = \mathbf{0}$ ,  $\mathbf{B}_0 = \mathbf{I}$ ,  $\boldsymbol{\phi}_0 = \mathbf{0}$ , and  $\boldsymbol{\Phi}_0 = \mathbf{I}$ . We use 15,000 MCMC iterations, with a burn-in of 5,000 and a thinning parameter of 5.

**R code. AR(1) model: Effects of inflation on interest rate**

```

1 rm(list = ls()); set.seed(10101)
2 DataIntRate <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/16INTDEF.csv"
  , sep = ",", header = TRUE, quote = "")
3 attach(DataIntRate)
4 yt <- diff(i3); ytlag <- dplyr::lag(yt, n = 1)
5 T <- length(yt)
6 Xt <- cbind(diff(inf), diff(def)); Xtag <- dplyr::lag(Xt, n
  = 1)
7 K <- dim(Xt)[2] + 1
8 Reg <- lm(yt ~ ytlag + I(Xt[,-1] - Xtag))
9 SumReg <- summary(Reg); SumReg
10 PostSig2 <- function(Beta, Phi){
11   Xstar<- matrix(NA, T-1, K - 1)
12   ystar <- matrix(NA, T-1, 1)
13   for(t in 2:T){
14     Xstar[t-1,] <- Xt[t,] - Phi*Xt[t-1,]
15     ystar[t-1,] <- yt[t] - Phi*yt[t-1]
16   }
17   Xstar <- cbind(1, Xstar)
18   an <- T - 1 + a0
19   dn <- d0 + t(ystar - Xstar%*%Beta)%*%(ystar - Xstar%*%Beta
    ) + t(Beta - b0)%*%B0i%*%(Beta - b0)
20   sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
21   return(sig2)
22 }
23 PostBeta <- function(sig2, Phi){
24   Xstar<- matrix(NA, T-1, K - 1)
25   ystar <- matrix(NA, T-1, 1)
26   for(t in 2:T){
27     Xstar[t-1,] <- Xt[t,] - Phi*Xt[t-1,]
28     ystar[t-1,] <- yt[t] - Phi*yt[t-1]
29   }
30   Xstar <- cbind(1, Xstar)
31   XtXstar <- t(Xstar)%*%Xstar
32   Xtystar <- t(Xstar)%*%ystar
33   Bn <- solve(B0i + XtXstar)
34   bn <- Bn%*%(B0i%*%b0 + Xtystar)
35   Beta <- MASS::mvrnorm(1, bn, sig2*Bn)
36   return(Beta)
37 }
```

**R code. AR(1) model: Effects of inflation on interest rate**

```

1 PostPhi <- function(sig2, Beta){
2   u <- yt - cbind(1,Xt)%*%Beta
3   U <- u[-T]
4   ustар <- u[-1]
5   UtU <- t(U)%*%U
6   UtU <- t(U)%*%ustар
7   Phin <- solve(Phi0i + sig2^(-1)*UtU)
8   phin <- Phin%*%(Phi0i%*%phi0 + sig2^(-1)*UtU)
9   Phi <- truncnorm::rtruncnorm(1, a = -1, b = 1, mean = phin
10   , sd = Phi^0.5)
11   return(Phi)
12 }
13 # Hyperparameters
14 d0 <- 0.01; a0 <- 0.01
15 b0 <- rep(0, K); c0 <- 1;
16 B0 <- c0*diag(K); B0i <- solve(B0)
17 phi0 <- 0; Phi0 <- 1; Phi0i <- 1/Phi0
18 # MCMC parameters
19 mcmc <- 15000; burnin <- 5000
20 tot <- mcmc + burnin; thin <- 1
21 PostBetas <- matrix(0, mcmc+burnin, K)
22 PostSigma2s <- rep(0, mcmc+burnin)
23 PostPhis <- rep(0, mcmc+burnin)
24 Beta <- rep(0, K); Phi <- 0
25 sig2 <- SumReg$sigma^2; Phi <- SumReg$coefficients[2,1]
26 Beta <- SumReg$coefficients[c(1,3,4),1]
27 pb <- txtProgressBar(min = 0, max = tot, style = 3)
28 for(s in 1:tot){
29   sig2 <- PostSig2(Beta = Beta, Phi = Phi)
30   PostSigma2s[s] <- sig2
31   Beta <- PostBeta(sig2 = sig2, Phi = Phi)
32   PostBetas[s,] <- Beta
33   Phi <- PostPhi(sig2 = sig2, Beta = Beta)
34   PostPhis[s] <- Phi
35   setTxtProgressBar(pb, s)
36 }
37 close(pb)
38 keep <- seq((burnin+1), tot, thin)
39 PosteriorBetas <- coda::mcmc(PostBetas[keep,])
40 summary(PosteriorBetas)
41 PosteriorSigma2 <- coda::mcmc(PostSigma2s[keep])
42 summary(PosteriorSigma2)
43 PosteriorPhi <- coda::mcmc(PostPhis[keep])
44 summary(PosteriorPhi)
45 dfBinf <- as.data.frame(PosteriorBetas[,2])

```

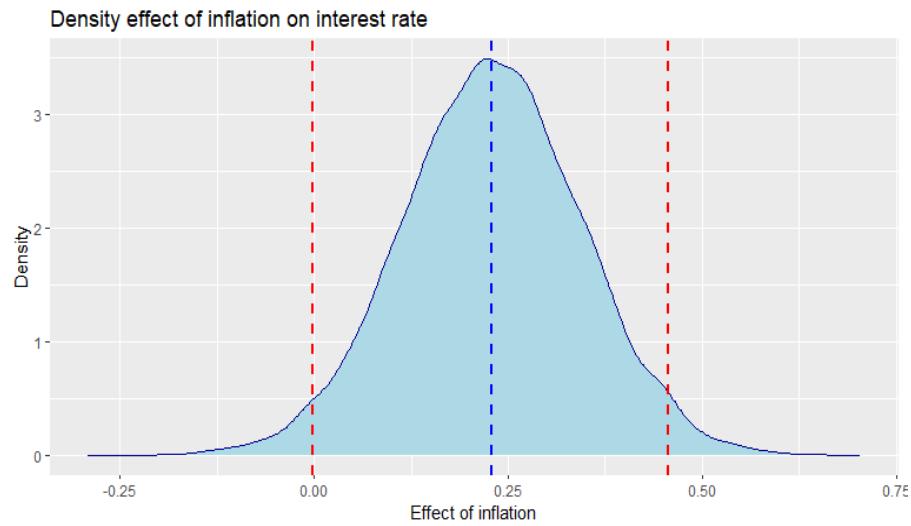
**R code. AR(1) model: Effects of inflation on interest rate**

```

1 # Basic density
2 library(ggplot2)
3 p <- ggplot(dfBinf, aes(x=var1)) + geom_density(color="darkblue",
4   fill="lightblue") + geom_vline(aes(xintercept = mean(var1)), color="blue", linetype="dashed", linewidth =1) + geom_vline(aes(xintercept=quantile(var1, 0.025)), color="red", linetype="dashed", linewidth=1) + geom_vline(aes(xintercept=quantile(var1, 0.975)), color="red", linetype="dashed", linewidth=1) + labs(title="Density effect of inflation on interest rate", x="Effect of inflation", y = "Density")
4

```

Figure 8.3 shows the posterior density plot of the effects of inflation rate on interest rate. The posterior mean of this coefficient is approximately 0.25, and the credible interval at 95% is (0, 0.46), which indicates again that the annual changes in interest rate are weakly positive related to annual changes in inflation (see Figure 8.2 as reference).



**FIGURE 8.3**  
Density: Effects of inflation on interest rate.

Observe that the previous setting encompasses the particular relevant case  $y_t \sim AR(p)$ , it is just omitting the covariates such that  $y_t = \mu_t$ . [91] extend the Bayesian inference of linear regression with  $AR(p)$  errors to  $ARMA(p, q)$  errors using a *state-space* representation.

Setting  $y_t = \mu_t$  such that  $y_t = \sum_{s=1}^p \phi_j y_{t-s} + \sum_{s=1}^q \theta_s \epsilon_{t-s} + \epsilon_t$ , letting  $r = \max\{p, q+1\}$ ,  $\phi_s = 0$  for  $s > p$  and  $\theta_s = 0$  for  $s > q$ , and defining  $\mathbf{x}^\top = [1 \ 0 \ \dots \ 0]$ , and  $\mathbf{H} = [1 \ \theta_1 \ \dots \ \theta_{r-1}]^\top$   $r$ -dimensional vectors, and

$$\mathbf{G} = \begin{bmatrix} \phi_1 & 1 & 0 & \dots & 0 \\ \phi_2 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & & \\ \phi_{r-1} & 0 & 0 & \dots & 1 \\ \phi_r & 0 & 0 & \dots & 0 \end{bmatrix} = \begin{bmatrix} \phi_1 & & & & \\ \phi_2 & & & & \\ \vdots & & & & \\ \phi_r & & & & \\ & 0 & 0 & \dots & 0 \end{bmatrix} \mathbf{I}_{r-1},$$

which is a  $r \times r$  dimensional matrix, and give the *state* vector  $\boldsymbol{\beta}_t = [\beta_{1,t} \ \beta_{2,t} \ \dots \ \beta_{r,t}]^\top$ , the  $ARMA$  model has the following representation:

$$\begin{aligned} y_t &= \mathbf{x}^\top \boldsymbol{\beta}_t \\ \boldsymbol{\beta}_t &= \mathbf{G} \boldsymbol{\beta}_{t-1} + \mathbf{H} \epsilon_t. \end{aligned}$$

This is a *dynamic linear model* where  $\boldsymbol{\Sigma}_t = 0$ , and  $\boldsymbol{\Omega}_t = \sigma^2 \mathbf{H} \mathbf{H}^\top$  (see [336, 91]).

A notable advantage of the *state-space* representation of the  $ARMA$  model is that the evaluation of the likelihood can be performed efficiently using the recursive laws. Extensions to autoregressive integrated moving average  $ARIMA(p, d, q)$  models are discussed in [336, Chap. 3]. In  $ARIMA(p, d, q)$  models,  $d$  refers to the level of integration (or differencing) required to eliminate the stochastic trend in a time series (see [136, Chap. 4] for details).

#### Example: $AR(2)$ process

Let's see the *state-space* representation of a stationary  $AR(2)$  process with intercept, that is,  $y_t = \mu + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \epsilon_t$ , where  $\epsilon_t \sim N(0, \sigma^2)$ . Thus,  $\mathbb{E}[y_t] = \frac{\mu}{1-\phi_1-\phi_2}$ , and variance  $Var[y_t] = \frac{\sigma^2(1-\phi_2)}{1-\phi_2-\phi_1^2-\phi_1^2\phi_2-\phi_2^2+\phi_2^3}$ .

In addition, we can proof that setting  $z_t = Y_t - \bar{\mu}$ , we have  $z_t = \phi_1 z_{t-1} + \phi_2 z_{t-2} + \epsilon_t$  where  $\mathbb{E}[z_t] = 0$ , and these are equivalent representations (see Exercise 5). Then, setting  $\mathbf{x}^\top = [1 \ 0]$ ,  $\mathbf{H} = [1 \ 0]^\top$ ,  $\mathbf{G} = \begin{bmatrix} \phi_1 & 1 \\ \phi_2 & 0 \end{bmatrix}$ ,  $\boldsymbol{\beta}_t = [\beta_{t1} \ \beta_{t2}]^\top$ ,  $\boldsymbol{\Sigma}_t = 0$  and  $\boldsymbol{\Omega}_t = \sigma^2$  we have

$$\begin{aligned} z_t &= \mathbf{x}^\top \boldsymbol{\beta}_t && \text{(Observation equation)} \\ \boldsymbol{\beta}_t &= \mathbf{G} \boldsymbol{\beta}_{t-1} + \mathbf{H} \epsilon_t && \text{(States equations).} \end{aligned}$$

We use the function `stan_sarima` from the package `bayesforecast` to perform Bayesian inference in  $ARMA$  models in our GUI. The following code shows

how to simulate an  $AR(2)$  process, and perform Bayesian inference using this function.

### R code. Simulation and inference: $AR(2)$ model

```

1 rm(list = ls()); set.seed(10101)
2 T <- 200; mu <- 0.5
3 phi1 <- 0.5; phi2 <- 0.3; sig <- 0.5
4 Ey <- mu/(1-phi1-phi2); Sigy <- sig*((1-phi2)/(1-phi1-phi2-
   ^2-phi2*phi1^2-phi2^2+phi2^3))^0.5
5 y <- rnorm(T, mean = Ey, sd = Sigy)
6 e <- rnorm(T, mean = 0, sd = sig)
7 for(t in 3:T){
8   y[t] <- mu + phi1*y[t-1] + phi2*y[t-2] + e[t]
9 }
10 mean(y); sd(y)
11 y <- ts(y, start=c(1820, 1), frequency=1)
12 plot(y)
13 iter <- 10000; burnin <- 5000; thin <- 1; tot <- iter +
   burnin
14 library(bayesforecast)
15 sf1 <- bayesforecast::stan_sarima(y, order = c(2, 0, 0),
   prior_mu0 = normal(0, 1),
16 prior_ar = normal(0, 1), prior_sigma0 = inverse.gamma(0.01/
   2, 0.01/2),
17 seasonal = c(0, 0, 0), iter = tot, warmup = burnin, chains =
   1)
18 keep <- seq(burnin+1, tot, thin)
19 Postmu <- sf1[["stanfit"]]\$sim[["samples"]][[1]][["mu0"]][
   keep]
20 Postsig <- sf1[["stanfit"]]\$sim[["samples"]][[1]][["sigma0"]]
   ][keep]
21 Postphi1 <- sf1[["stanfit"]]\$sim[["samples"]][[1]][["ar0[1]"]]
   ][keep]
22 Postphi2 <- sf1[["stanfit"]]\$sim[["samples"]][[1]][["ar0[2]"]]
   ][keep]
23 Postdraws <- coda::mcmc(cbind(Postmu, Postsig, Postphi1,
   Postphi2))
24 summary(Postdraws)
25 plot(Postdraws)

```

We perform 10,000 MCMC iterations plus a burn-in equal 5,000 assuming  $\sigma^2 \sim IG(0.01/2, 0.01/2)$ ,  $\mu \sim N(0, 1)$  and  $\phi_k \sim N(0, 1)$ ,  $k = 1, 2$ . The trace plots look well, and all 95% credible intervals encompass the population values.

Algorithm A22 shows how to perform inference in  $ARMA(p, q)$  models using our GUI. See also Chapter 5 for details regarding the dataset structure.

The function *stan\_sarima* uses the *Stan* software [423], which in turn employs *Hamiltonian Monte Carlo* (HMC). The following code illustrates how

**Algorithm A22** Autoregressive Moving Average (*ARMA*) models

- 
- 1: Select *Time series Model* on the top panel
  - 2: Select *ARMA* using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Set the order of the ARMA model,  $p$  and  $q$  parameters
  - 6: Set the frequency: annual (1), quarterly (4), monthly (12)
  - 7: Set the location and scale hyperparameters of the *intercept*, autoregressive (*AR*), moving average (*MA*) and standard deviation. Take into account that there is just one set of hyperparameters for *AR* and *MA* coefficients. This step is not necessary as by default our GUI uses non-informative priors
  - 8: Click the *Go!* button
  - 9: Analyze results
  - 10: Download posterior chains and figures (density, autocorrelation and trace plots) using the *Download Results* button
- 

to perform Bayesian inference in the *AR*(2) model by programming the HMC from scratch. It is important to note that this is only an illustration, as HMC is less efficient than the Gibbs sampler in this example. However, HMC can outperform traditional MCMC algorithms in more complex models, particularly when dealing with high-dimensional probability distributions or when MCMC struggles with poor mixing due to posterior correlation.

In the first block, we perform the simulation by setting  $\mu = 0.5$ ,  $\phi_1 = 0.5$ ,  $\phi_2 = 0.3$ ,  $\sigma = 0.25$ , and a sample size of 200. We then set the hyperparameters and define the function to calculate the logarithm of the posterior distribution. The model is parametrized using  $\tau = \log(\sigma^2)$ , such that  $\sigma^2 = \exp(\tau)$ , which avoids issues related to the non-negativity restriction of  $\sigma^2$ . As a result, we need to account for the Jacobian due to this transformation, specifically  $d\sigma^2/d\tau = \exp(\tau)$ .

Next, we define the function to compute the gradient vector of the log posterior distribution. It is preferable to calculate the gradient vector analytically, as using finite differences can be computationally expensive. However, it is a good practice to check the analytical calculations by evaluating the function at the maximum posterior estimate, where the function should return values close to 0, or by comparing the results with finite differences at a few evaluation points.

The posterior distribution is given by<sup>3</sup>

$$\begin{aligned}\pi(\mu, \phi_1, \phi_2, \tau \mid \mathbf{y}) &\propto \prod_{t=3}^T \exp \left\{ -\frac{1}{2 \exp(\tau)} (y_t - \mu - \phi_1 y_{t-1} - \phi_2 y_{t-2})^2 \right\} \\ &\quad \times (\exp(\tau))^{-1/2} \exp \left\{ -\frac{1}{2\sigma_\mu^2} (\mu - \mu_0)^2 \right\} \\ &\quad \times \exp \left\{ -\frac{1}{2\sigma_{\phi_1}^2} (\phi_1 - \phi_{10})^2 \right\} \\ &\quad \times \exp \left\{ -\frac{1}{2\sigma_{\phi_2}^2} (\phi_2 - \phi_{20})^2 \right\} \times \exp \{-(\alpha_0/2 + 1)\tau\} \\ &\quad \times \exp \{-\delta_0/(2 \exp(\tau))\} \exp(\tau).\end{aligned}$$

The components of the gradient vector of the log posterior distribution are given by

$$\begin{aligned}\frac{\partial \log(\pi(\mu, \phi_1, \phi_2, \tau \mid \mathbf{y}))}{\partial \mu} &= \frac{\sum_{t=3}^T (y_t - \mu - \phi_1 y_{t-1} - \phi_2 y_{t-2})}{\exp(\tau)} - \frac{1}{\sigma_\mu^2} (\mu - \mu_0) \\ \frac{\partial \log(\pi(\mu, \phi_1, \phi_2, \tau \mid \mathbf{y}))}{\partial \phi_1} &= \frac{\sum_{t=3}^T (y_t - \mu - \phi_1 y_{t-1} - \phi_2 y_{t-2}) y_{t-1}}{\exp(\tau)} - \frac{1}{\sigma_{\phi_1}^2} (\phi_1 - \phi_{10}) \\ \frac{\partial \log(\pi(\mu, \phi_1, \phi_2, \tau \mid \mathbf{y}))}{\partial \phi_2} &= \frac{\sum_{t=3}^T (y_t - \mu - \phi_1 y_{t-1} - \phi_2 y_{t-2}) y_{t-2}}{\exp(\tau)} - \frac{1}{\sigma_{\phi_2}^2} (\phi_2 - \phi_{20}) \\ \frac{\partial \log(\pi(\mu, \phi_1, \phi_2, \tau \mid \mathbf{y}))}{\partial \tau} &= -\frac{(T-2)}{2} + \frac{\sum_{t=3}^T (y_t - \mu - \phi_1 y_{t-1} - \phi_2 y_{t-2})^2}{2 \exp(\tau)} \\ &\quad - (\alpha_0/2 + 1) + \delta_0/(2 \exp(\tau)) + 1.\end{aligned}$$

Next, we provide the code for the Hamiltonian Monte Carlo, as outlined in Chapter 4. The initial values are set as follows:  $\mu = \bar{y} = \frac{1}{T-2} \sum_{t=3}^T y_t$ ,  $\phi_1 = \phi_2 = 0$ , and  $\tau = \exp \left( \frac{1}{T-2} \sum_{t=3}^T (y_t - \bar{y})^2 \right)$ , with  $M$  being the inverse covariance matrix of the posterior distribution evaluated at its maximum. Additionally,  $\epsilon$  is randomly drawn from a uniform distribution between 0 and  $2\epsilon_0$ , and  $L$  is set to the highest integer near  $1/\epsilon$ , in order to approximately satisfy  $L\epsilon = 1$ , where  $\epsilon_0 = 0.1$ .

We can verify that all 95% credible intervals encompass the population values, and the posterior means are close to the population values. The acceptance rate averages above 65%, so we should consider increasing the base step ( $\epsilon_0$ ). Furthermore, we do not impose the stationarity conditions on  $\phi_1$

---

<sup>3</sup>Take into account that we do not consider the first two observations when present the likelihood, this is no an issue when there is a large sample size.

and  $\phi_2$ . Exercise 6 asks to program an HMC that takes these requirements into account.

### R code. Simulation and inference: AR(2) model using Hamiltonian Monte Carlo

```

1 # Simulation AR(2)
2 rm(list = ls()); set.seed(10101); T <- 1000; K <- 4
3 mu <- 0.5; phi1 <- 0.5; phi2 <- 0.3; sig <- 0.5
4 Ey <- mu/(1-phi1-phi2); Sigy <- sig*((1-phi2)/(1-phi1
  ^2-phi2*phi1^2-phi2^2+phi2^3))^0.5
5 y <- rnorm(T, mean = Ey, sd = Sigy); e <- rnorm(T, mean = 0,
  sd = sig)
6 for(t in 3:T){
7   y[t] <- mu + phi1*y[t-1] + phi2*y[t-2] + e[t]
8 }
9 # Hyperparameters
10 d0 <- 0.01; a0 <- 0.01; mu0 <- 0; MU0 <- 1
11 phio <- c(0, 0); Phio <- diag(2)
12 # Log posterior multiply by -1 to use optim
13 LogPost <- function(theta, y){
14   mu <- theta[1]; phi1 <- theta[2]; phi2 <- theta[3]
15   tau <- theta[4]; sig2 <- exp(tau); logLik <- NULL
16   for(t in 3:T){
17     logLikt <- dnorm(y[t], mean = mu + phi1*y[t-1] + phi2*y[
      t-2], sd = sig2^0.5, log = TRUE)
18     logLik <- c(logLik, logLikt)
19   }
20   logLik <- sum(logLik)
21   logPrior <- dnorm(mu, mean = mu0, sd = MU0^0.5, log = TRUE
    ) + dnorm(phi1, mean = phio[1], sd = Phio[1,1]^0.5, log
    = TRUE) + dnorm(phi2, mean = phio[2], sd = Phio
    [2,2]^0.5, log = TRUE) + invgamma::dinvgamma(sig2, shape
    = a0/2, rate = d0/2, log = TRUE)
22   logPosterior <- logLik + logPrior + tau
23   return(-logPosterior) # Multiply by -1 to minimize using
    optim
24 }
25 theta0 <- c(mean(y), 0, 0, var(y))
26 Opt <- optim(theta0, LogPost, y = y, hessian = TRUE)
27 theta0 <- Opt$par; VarPost <- solve(Opt$hessian)
```

**R code. Simulation and inference: AR(2) model using Hamiltonian Monte Carlo**

```

1 # Gradient log posterior
2 GradientTheta <- function(theta, y){
3   mu <- theta[1]; phi1 <- theta[2]; phi2 <- theta[3]
4   tau <- theta[4]; sig2 <- exp(tau); SumLik <- matrix(0, 3,
5     1)
6   SumLik2 <- NULL
7   for(t in 3:T){
8     xt <- matrix(c(1, y[t-1], y[t-2]), 3, 1)
9     SumLik <- (y[t] - (mu + phi1*y[t-1] + phi2*y[t-2]))*xt
10    SumLik2t <- (y[t] - (mu + phi1*y[t-1] + phi2*y[t-2]))^2
11    SumLik <- rowSums(cbind(SumLik, SumLik2t))
12    SumLik2 <- sum(SumLik2, SumLik2t)
13  }
14  Grad_mu <- SumLik[1]/sig2 - (1/MU0)*(mu - mu0)
15  Grad_phi1 <- SumLik[2]/exp(tau) - 1/Phi0[1,1]*(phi1 - phi0
16    [1])
17  Grad_phi2 <- SumLik[3]/exp(tau) - 1/Phi0[2,2]*(phi2 - phi0
18    [2])
19  Grad_tau <- -(T-2)/2 + SumLik2/(2*exp(tau)) - (a0/2 + 1) +
20    d0/(2*exp(tau)) + 1
21  Grad <- c(Grad_mu, Grad_phi1, Grad_phi2, Grad_tau)
22  return(Grad)
23 }
24 # Hamiltonian Monte Carlo function
25 HMC <- function(theta, y, epsilon, M){
26   L <- ceiling(1/epsilon)
27   Minv <- solve(M); thetat <- theta
28   K <- length(theta)
29   mom <- t(mvtnorm::rmvnorm(1, rep(0, K), M))
30   logPost_Mom_t <- -LogPost(theta, y) + mvtnorm::dmvnorm(t
31     (mom), rep(0, K), M, log = TRUE)
32   for(l in 1:L){
33     if(l == 1 | l == L){
34       mom <- mom + 0.5*epsilon*GradientTheta(theta, y)
35       theta <- theta + epsilon*Minv%*%mom
36     }else{
37       mom <- mom + epsilon*GradientTheta(theta, y)
38       theta <- theta + epsilon*Minv%*%mom
39     }
40   }

```

**R code. Simulation and inference: AR(2) model using Hamiltonian Monte Carlo**

```

1 logPost_Mom_star <- -LogPost(theta, y) + mvtnorm::dmvnorm
2   (t(mom), rep(0, K), M, log = TRUE)
3 alpha <- min(1, exp(logPost_Mom_star-logPost_Mom_t))
4 u <- runif(1)
5 if(u <= alpha){
6   thetaNew <- c(theta)
7 }else{
8   thetaNew <- thetat
9 }
10 rest <- list(theta = thetaNew, Prob = alpha)
11 return(rest)
12 }
13 # Posterior draws
14 S <- 1000; burnin <- 1000; thin <- 2; tot <- S + burnin
15 thetaPost <- matrix(NA, tot, K)
16 ProbAccept <- rep(NA, tot)
17 theta0 <- c(mean(y), 0, 0, exp(var(y)))
18 M <- solve(VarPost); epsilon0 <- 0.1
19 pb <- txtProgressBar(min = 0, max = tot, style = 3)
20 for(s in 1:tot){
21   epsilon <- runif(1, 0, 2*epsilon0)
22   L <- ceiling(1/epsilon)
23   HMCs <- HMC(theta = theta0, y, epsilon, M)
24   theta0 <- HMCs$theta
25   thetaPost[s,] <- HMCs$theta
26   ProbAccept[s] <- HMCs$Prob
27   setTxtProgressBar(pb, s)
28 }
29 close(pb)
30 keep <- seq((burnin+1), tot, thin)
31 thetaF <- coda::mcmc(thetaPost[keep,])
32 summary(thetaF)
33 ProbAcceptF <- coda::mcmc(ProbAccept[keep])
34 summary(ProbAcceptF)

```

### 8.3 Stochastic volatility models

A notable example of non-linear and non-Gaussian *state-space models* is stochastic volatility models (SVMs), which are widely used to model the

volatility of financial returns. SVMs have gained significant attention due to their flexibility, ability to capture complex dynamics such as asymmetries, and ease of generalization to simultaneously model multiple returns, making them advantageous over generalized autoregressive conditional heteroskedasticity (GARCH) models proposed by [52]. However, estimating SVMs is more challenging than estimating GARCH models. This is because GARCH models set variance in a deterministic manner, whereas SVMs do so stochastically. Consequently, GARCH models are typically estimated using maximum likelihood methods, while SVMs require Bayesian approaches, adding complexity to the estimation process.

The specification of the stochastic volatility model is given by

$$y_t = \mathbf{x}_t^\top \boldsymbol{\beta} + \exp\{0.5h_t\} \mu_t \quad (\text{Observation equation}) \quad (8.8)$$

$$h_t = \mu + \phi(h_{t-1} - \mu) + \sigma w_t \quad (\text{State equation}), \quad (8.9)$$

where  $y_t$  are the log-returns,  $\mathbf{x}_t$  are controls,  $\boldsymbol{\beta}$  are time-invariant location parameters,  $\mu_t \sim N(0, 1)$ ,  $w_t \sim N(0, 1)$ ,  $\mu_t \perp w_t$ , the initial log-variance process  $h_0 \sim N(\mu, \sigma^2/(1 - \phi^2))$ ,  $\mu$ ,  $\phi$  and  $\sigma$  are the level, persistence and standard deviation of the log-variance, respectively.

Given the specification in Equations 8.8 and 8.9, we can write the observation equation as

$$\log\{(y_t - \mathbf{x}_t^\top \boldsymbol{\beta})^2\} = h_t + \log(\mu_t^2),$$

which leads to a linear, but non-Gaussian, *state-space model*. [242] approximate the distribution of  $\log(\mu_t^2)$  by a mixture of normal distributions, that is,

$$\log(\mu_t^2) | l_t \sim N(m_{l_t}, s_{l_t}^2),$$

where  $l_t \in \{1, 2, \dots, 10\}$  defines the mixture component indicator at time  $t$ . Thus, the model can be written as

$$\log\{(y_t - \mathbf{x}_t^\top \boldsymbol{\beta})^2\} = h_t + \log(\mu_t^2),$$

and

$$h_t = \mu + \phi(h_{t-1} - \mu) + \sigma w_t.$$

This forms a linear and conditionally Gaussian *state-space model*, where

$$\log\{(y_t - \mathbf{x}_t^\top \boldsymbol{\beta})^2\} = m_{l_t} + h_t + \mu_t,$$

and

$$\mu_t \sim N(0, s_{l_t}^2).$$

We use the *stochvol* package in our GUI to perform MCMC inference in the SVMs [211]; this package is based on the MCMC algorithms proposed by [242]. The default prior distributions in the *stochvol* package are:

$$\boldsymbol{\beta} \sim N(\mathbf{b}_0, \mathbf{B}_0), \quad \mu \sim N(\mu_0, \sigma_{\mu 0}^2), \quad \frac{\phi + 1}{2} \sim B(\alpha_0, \beta_0), \quad \sigma^2 \sim G\left(\frac{1}{2}, \frac{1}{2\sigma_{\sigma^2}^2}\right).$$

The prior distribution for  $\phi$  is set to ensure stationarity of the process ( $\phi \in (-1, 1)$ ). In most applications,  $\phi \approx 1$ , so the authors of the package recommend setting  $\alpha_0 \gtrsim 5$  and  $\beta_0 \approx 1.5$ . The prior distribution for  $\sigma$  is  $|N(0, \sigma_{\sigma^2}^2)|$  (a half-normal distribution). This is recommended by the authors since the conjugate inverse-gamma distribution does not work well in this case, as it bounds  $\sigma$  away from 0, which is undesirable when modeling the log-variance of log-returns.

Algorithm A23 shows how to perform inference in stochastic volatility models using our GUI. See also Chapter 5 for details regarding the dataset structure.

---

**Algorithm A23** Stochastic volatility models

---

- 1: Select *Time series Model* on the top panel
  - 2: Select *Stochastic volatility* using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Set the hyperparameters: the mean and standard deviation of the Gaussian prior for the regression parameters, mean and standard deviation for the Gaussian prior distribution of the level of the log-volatility, shape parameters for the Beta prior distribution of the transformed persistence parameter, and the positive real number, which stands for the scaling of the transformed volatility of log-volatility. This step is not necessary as by default our GUI uses default values in *stochvol* package
  - 6: Click the *Go!* button
  - 7: Analyze results
  - 8: Download posterior chains of the fixed coefficients, and the states using the *Download Results* button
- 

**Example: Simulation exercise of the stochastic volatility model**

The following code shows how to simulate and perform Bayesian inference in the stochastic volatility model using the function *svsample* from the *stochvol* package. We set the stochastic volatility parameters to  $\mu = -10$ ,  $\phi = 0.95$ , and  $\sigma = 0.3$ . We assume two regressors, which are distributed as standard normal, with  $\boldsymbol{\beta} = [0.5 \ 0.3]^\top$ , and the sample size is 1,250, which corresponds to approximately 5 years of daily returns. We use the default hyperparameters: 10,000 MCMC iterations, a burn-in of 5,000, and a thinning parameter of 5.

The summary statistics of the posterior draws show that all 95% credible intervals encompass the population parameters, and the posterior chains appear to have converged. Figure 8.4 displays the posterior results for the volatility ( $h_t$ ). The posterior mean (blue) follows the “observed” series (black), and the 95% credible intervals (light blue) typically encompass the “observed” series.

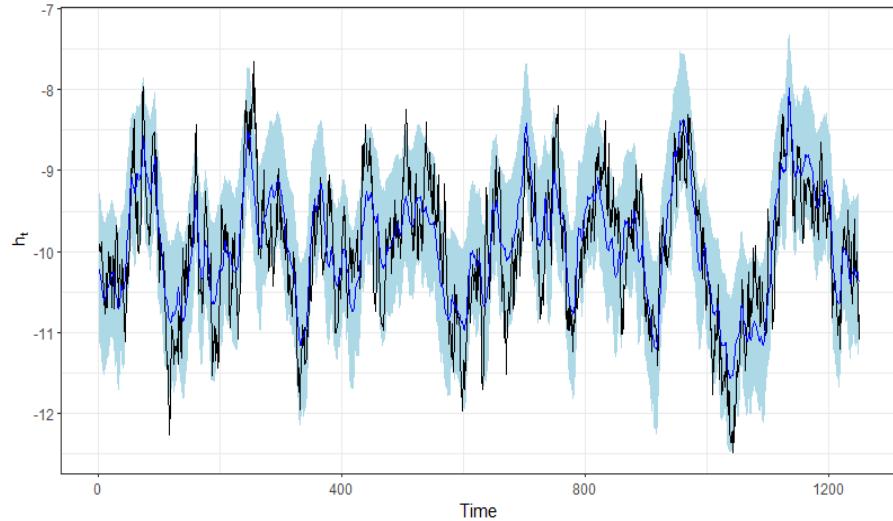
### R code. Simulation and inference: Stochastic volatility model

```

1 rm(list = ls()); set.seed(10101)
2 T <- 1250; K <- 2
3 X <- matrix(rnorm(T*K), T, K)
4 B <- c(0.5, 0.3); mu <- -10; phi <- 0.95; sigma <- 0.3
5 h <- numeric(T); y <- numeric(T)
6 h[1] <- rnorm(1, mu, sigma / sqrt(1 - phi^2)) # Initial
    state
7 y[1] <- X[,1]*%*%B + rnorm(1, 0, exp(h[1] / 2))           #
    Initial observation
8 for (t in 2:T) {
9   h[t] <- mu + phi*(h[t-1]-mu) + rnorm(1, 0, sigma)
10  y[t] <- X[,t]*%*%B + rnorm(1, 0, sd = exp(0.5*h[t]))
11 }
12 df <- as.data.frame(cbind(y, X))
13 colnames(df) <- c("y", "x1", "x2")
14 MCMC <- 10000; burnin <- 10000; thin <- 5
15 res <- stochvol::svsample(y, designmatrix = X, draws = MCMC,
    burnin = burnin, thin = thin, priormu = c(0, 100),
    priorsigma = c(1), priorphi = c(5, 1.5), priorbeta = c
    (0, 10000))
16 summary(res[["para"]][[1]][,-c(4,5)])
17 summary(res[["beta"]])
18 ht <- res[["latent"]][[1]]
19 library(dplyr)
20 library(ggplot2)
21 library(latex2exp)
22 ggplot2::theme_set(theme_bw())
23 x_means <- colMeans(ht)
24 x_quantiles <- apply(ht, 2, function(x) quantile(x, probs =
    c(0.025, 0.975)))
25 df <- tibble(t = seq(1, T), mean = x_means, lower = x_
    quantiles[1, ], upper = x_quantiles[2, ], x_true = h,
    observations = y)
26 plot_filtering_estimates <- function(df) {
27   p <- ggplot(data = df, aes(x = t)) + geom_ribbon(aes(ymin =
      lower, ymax = upper), alpha = 1, fill = "lightblue") +
      geom_line(aes(y = x_true), colour = "black", alpha = 1,
      linewidth = 0.5) + geom_line(aes(y = mean), colour =
      "blue", linewidth = 0.5) + ylab(TeX("\$h_{\$t\$}")) + xlab("Time")
28   print(p)
29 }
30 plot_filtering_estimates(df)

```

So far, we have used MCMC algorithms to perform inference in *state-space models*. These algorithms require all observations to estimate the unknown



**FIGURE 8.4**  
Stochastic volatility model.

parameters, a process referred to as offline or batch inference. However, this approach has limitations when online inference is needed, as every new observation requires simulating a new posterior chain. This is because MCMC algorithms do not naturally adapt to sequential updates. In contrast, particle filter algorithms, which are a subset of sequential Monte Carlo (SMC) methods, are specifically designed for sequential use, making them suitable for online inference.

Remember from Chapter 4 that particle filters (sequential Monte Carlo) are algorithms that allow computing a numerical approximation to the filtering distribution  $\pi(\theta_{1:t} | \mathbf{y}_{1:t})$  sequentially in time. This is particularly relevant in non-linear and non-Gaussian models where there is no analytical solution for the filtering distribution.

The following code shows how to perform particle filtering in the vanilla stochastic volatility model assuming that the proposal distribution is the conditional prior distribution, that is,  $q(h_t | h_{t-1}, y_t) = \pi(h_t | h_{t-1})$ , which is normal with mean  $\mu + \phi(h_{t-1} - \mu)$  and variance  $\sigma^2$ . This choice implies that the incremental importance weights are equal to  $p(y_t | h_t)$ , which is  $N(0, \exp(h_t))$ . Therefore, the weights are proportional to the likelihood function. We perform multinomial resampling every time period in the code, and start the algorithm in the stationary distribution of  $h_t$ . Remember that there are other resampling approaches that are more efficient, for instance, residual resampling (see Section 4.3). We ask in Exercise 7 to modify this code to perform resampling when the effective sample size is lower than 50% of the initial number of particles. In addition, we ask to program a sequential importance

sampling, and check why is important to perform resampling in this simple example.

Figure 8.5 illustrates the filtering recursion using SMC with uneven weights (blue line), even weights (purple line), bands corresponding to plus/minus two standard deviations (light blue shaded area), and the true state (black line).<sup>4</sup> The results indicate that SMC performs well even with a simple implementation, with no significant differences between using even and uneven weights (see Chapter 4).

In this example, we use the population parameters to perform the filtering recursion. However, this is not the case in practice, as we must estimate the time-invariant parameters. Therefore, more elaborate algorithms are required to achieve this (see Chapter 4). For instance, [14] propose particle Markov chain Monte Carlo, a family of methods that combines MCMC and SMC. See [110] for a tutorial on particle Metropolis-Hastings in **R**. A potential practical solution for applications that require sequential updating of a posterior distribution over an unbounded time horizon is to estimate the time-invariant parameters offline using MCMC algorithms up to a specific time period, and then update the state vector sequentially online during subsequent time periods, iterating this process. This is not optimal, but it can be practical.

***R code. Simulation and inference: Stochastic volatility model programming sequential Monte Carlo from scratch***

```

1 rm(list = ls()); set.seed(10101)
2 T <- 1250; mu <- -10; phi <- 0.95; sigma <- 0.3
3 h <- numeric(T); y <- numeric(T)
4 h[1] <- rnorm(1, mu, sigma / sqrt(1 - phi^2))
5 y[1] <- rnorm(1, 0, exp(h[1] / 2))
6 for (t in 2:T) {
7   h[t] <- mu + phi*(h[t-1]-mu) + rnorm(1, 0, sigma)
8   y[t] <- rnorm(1, 0, sd = exp(0.5*h[t]))
9 }
10 N <- 10000
11 log_Weights <- matrix(NA, N, T) # Log weights
12 Weights <- matrix(NA, N, T) # Weights
13 WeightsST <- matrix(NA, N, T) # Normalized weights
14 WeightsSTT <- matrix(1/N, N, T) # Normalized weights bar
15 particles <- matrix(NA, N, T) # Particles
16 particlesT <- matrix(NA, N, T) # Particles bar

```

---

<sup>4</sup>This standard deviation estimates the conditional posterior's standard deviation derived from the particles, not the estimator's standard deviation. The latter requires several independent particle runs on the same data.

**R code. Simulation and inference: Stochastic volatility model programming sequential Monte Carlo from scratch**

```

1 logalphas <- matrix(NA, N, T)    # Incremental importance
2 particles[, 1] <- rnorm(N, mu, sigma / sqrt(1 - phi^2))  #
   Stationary prior
3 log_Weights[, 1] <- dnorm(y[1], 0, sd = exp(0.5*particles
   [,1]), log = TRUE)  # Likelihood
4 Weights[, 1] <- exp(log_Weights[, 1])
5 WeightsST[, 1] <- Weights[, 1] / sum(Weights[, 1])
6 ind <- sample(1:N, size = N, replace = TRUE, prob =
   WeightsST[, 1])  # Resample
7 particles[, 1] <- particles[ind, 1] # Resampled particles
8 particlesT[, 1] <- particles[, 1] # Resampled particles
9 WeightsST[, 1] <- rep(1/N, N) # Resampled weights
10 pb <- txtProgressBar(min = 0, max = T, style = 3)
11 for (t in 2:T) {
12   particles[, t] <- rnorm(N, mu + phi*(particles[, t - 1] -
     mu), sigma)  # Sample from proposal
13   logalphas[, t] <- dnorm(y[t], 0, sd = exp(0.5*particles[,t
     ]), log = TRUE)
14   Weights[, t] <- exp(logalphas[, t])
15   WeightsST[, t] <- Weights[, t] / sum(Weights[, t])
16   if(t < T){
17     ind <- sample(1:N, size = N, replace = TRUE, prob =
     WeightsST[, t])
18     particles[, 1:t] <- particles[ind, 1:t]
19   }else{
20     ind <- sample(1:N, size = N, replace = TRUE, prob =
     WeightsST[, t])
21     particlesT[, 1:t] <- particles[ind, 1:t]
22   }
23   setTxtProgressBar(pb, t)
24 }
25 close(pb)
26 FilterDist <- colSums(particles * WeightsST)
27 SDFilterDist <- (colSums(particles^2 * WeightsST) -
   FilterDist^2)^0.5
28 FilterDistT <- colSums(particlesT * WeightsSTT)
29 SDFilterDistT <- (colSums(particlesT^2 * WeightsSTT) -
   FilterDistT^2)^0.5
30 MargLik <- colMeans(Weights)
31 plot(MargLik, type = "l")
32 library(dplyr); library(ggplot2); library(latex2exp)

```

**R code. Simulation and inference: Stochastic volatility model programming sequential Monte Carlo from scratch**

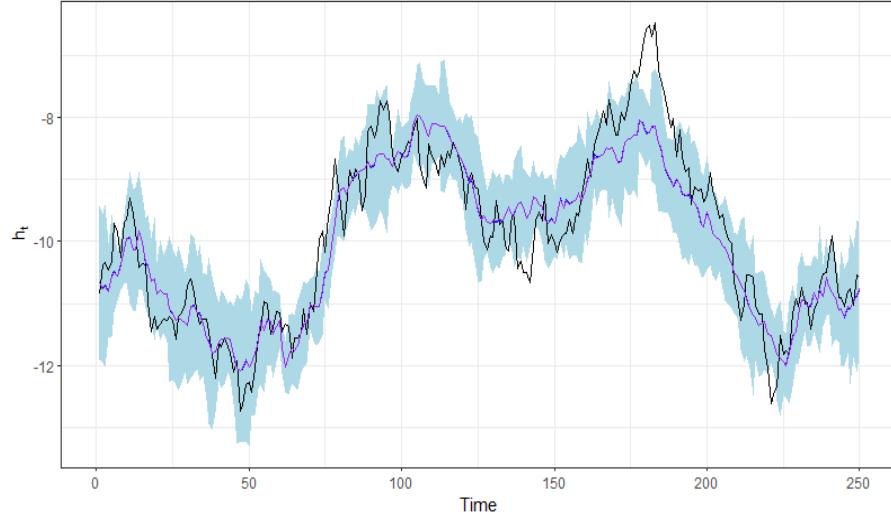
```

1 ggplot2:::theme_set(theme_bw())
2 Tfig <- 250; keepFig <- 1:Tfig
3 df <- tibble(t = keepFig,
4   mean = FilterDist[keepFig],
5   lower = FilterDist[keepFig] - 2*SDFilterDist[keepFig],
6   upper = FilterDist[keepFig] + 2*SDFilterDist[keepFig],
7   meanT = FilterDistT[keepFig],
8   lowerT = FilterDistT[keepFig] - 2*SDFilterDistT[keepFig],
9   upperT = FilterDistT[keepFig] + 2*SDFilterDistT[keepFig],
10  x_true = h[keepFig])
11 plot_filtering_estimates <- function(df) {
12   p <- ggplot(data = df, aes(x = t)) +
13     geom_ribbon(aes(ymin = lower, ymax = upper), alpha = 1,
14     fill = "lightblue") +
15     geom_line(aes(y = x_true), colour = "black", alpha = 1,
16     linewidth = 0.5) +
17     geom_line(aes(y = mean), colour = "blue", linewidth = 0.5)
18     +
19     geom_line(aes(y = meanT), colour = "purple", linewidth =
20     0.5) +
21     ylab(TeX("$h_{\{t\}}$")) + xlab("Time")
22   print(p)
23 }
24 plot_filtering_estimates(df)

```

## 8.4 Vector Autoregressive models

Another widely used methodological approach in time series analysis is the vector autoregressive (VAR) model, which extends AR( $p$ ) models to the multivariate case. Since the seminal work by Sims (1980) [412], these models have become a cornerstone of macroeconomic research to perform forecasts, and impulse-response (structural) analysis. This chapter provides an introduction to Bayesian inference in VAR models, with detailed discussions available in [248, 117, 458, 76].

**FIGURE 8.5**

Stochastic volatility model: Sequential Monte Carlo (SMC).

The *reduced-form* VAR(p) model can be written as

$$\mathbf{y}_t = \mathbf{v} + \sum_{j=1}^p \mathbf{A}_j \mathbf{y}_{t-j} + \boldsymbol{\mu}_t, \quad (8.10)$$

where  $\mathbf{y}_t$  is a  $M$ -dimensional vector having information of  $M$  time series variables,  $\mathbf{v}$  is a  $M$ -dimensional vector of intercepts,  $\mathbf{A}_j$  are  $M \times M$  matrices of coefficients, and  $\boldsymbol{\mu}_t \stackrel{iid}{\sim} N_M(\mathbf{0}, \boldsymbol{\Sigma})$  are stochastic errors,  $t = 1, 2, \dots, T$  and  $j = 1, 2, \dots, p$ . Other deterministic terms and exogenous variables can be added to the specification without main difficulty, we do not do this to keep simply the notation. In addition, we assume that the stability condition is satisfied such that the stochastic process is stationary (see [203, Chap. 2] for details), and we have available  $p$  presample values for each variable.

Following the matrix-form notation of the multivariate regression model (see sections 3.4 and 7.1), we can set  $\mathbf{Y} = [\mathbf{y}_1 \mathbf{y}_2 \dots \mathbf{y}_M]$ , which is an  $T \times M$  matrix,  $\mathbf{x}_t = [1 \mathbf{y}_{t-1}^\top \dots \mathbf{y}_{t-p}^\top]$  is a  $(1+Mp)$ -dimensional row vector, we define  $k = 1 + Mp$  to facilitate notation, and set

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix},$$

which is a  $T \times k$  matrix,  $\mathbf{B} = [\mathbf{v} \ \mathbf{A}_1 \ \mathbf{A}_2 \dots \mathbf{A}_P]^\top$  is a  $k \times M$  matrix of param-

eters, and  $\mathbf{U} = [\boldsymbol{\mu}_1 \ \boldsymbol{\mu}_2 \dots \boldsymbol{\mu}_M]$  is a  $T \times M$ -dimensional matrix of stochastic random errors such that  $\mathbf{U} \sim N_{T \times M}(\mathbf{0}_{T \times M}, \boldsymbol{\Sigma} \otimes \mathbf{I}_T)$ . Thus, we can express the VAR(p) model in the form of a multivariate regression model,

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{U}.$$

We can assume conjugate priors to facilitate computation, that is,

$$\pi(\mathbf{B}, \boldsymbol{\Sigma}) = \pi(\mathbf{B} | \boldsymbol{\Sigma})\pi(\boldsymbol{\Sigma}),$$

where  $\mathbf{B} | \boldsymbol{\Sigma} \sim N_{k \times M}(\mathbf{B}_0, \mathbf{V}_0, \boldsymbol{\Sigma})$  and  $\boldsymbol{\Sigma} \sim IW(\boldsymbol{\Psi}_0, \alpha_0)$ . Thus,

$$\pi(\mathbf{B}, \boldsymbol{\Sigma} | \mathbf{Y}, \mathbf{X}) = \pi(\mathbf{B} | \boldsymbol{\Sigma}, \mathbf{Y}, \mathbf{X})\pi(\boldsymbol{\Sigma} | \mathbf{Y}, \mathbf{X}),$$

where  $\mathbf{B} | \boldsymbol{\Sigma}, \mathbf{Y}, \mathbf{X} \sim N_{k \times M}(\mathbf{B}_n, \mathbf{V}_n, \boldsymbol{\Sigma})$  and  $\boldsymbol{\Sigma} | \mathbf{Y}, \mathbf{X} \sim IW(\boldsymbol{\Psi}_n, \alpha_n)$ . The quantities  $\mathbf{B}_n$ ,  $\mathbf{V}_n$ ,  $\boldsymbol{\Psi}_n$ , and  $\alpha_n$  are given by the following expressions:

$$\begin{aligned} \mathbf{B}_n &= (\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}(\mathbf{V}_0^{-1} \mathbf{B}_0 + \mathbf{X}^\top \mathbf{X} \hat{\mathbf{B}}), \\ \mathbf{V}_n &= (\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}, \\ \boldsymbol{\Psi}_n &= \boldsymbol{\Psi}_0 + \mathbf{S} + \mathbf{B}_0^\top \mathbf{V}_0^{-1} \mathbf{B}_0 + \hat{\mathbf{B}}^\top \mathbf{X}^\top \mathbf{X} \hat{\mathbf{B}} - \mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n, \\ \mathbf{S} &= (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}})^\top (\mathbf{Y} - \mathbf{X} \hat{\mathbf{B}}), \\ \hat{\mathbf{B}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}, \end{aligned}$$

and

$$\alpha_n = T + \alpha_0.$$

Thus, we see that once we express a VAR(p) model in the correct form, we can perform Bayesian inference as we did in the multivariate regression model. However, assuming conjugate priors has some limitations. First, VAR(p) models have many parameters. For instance, with 4 lags and 6 variables, we would have 150 location parameters  $((1 + (6 \times 4)) \times 6)$  and 21 scale parameters  $(6 \times (6 + 1)/2)$  for the covariance matrix. This can lead to a loss of precision, especially when using macroeconomic data, due to the typical lack of large sample sizes. Therefore, it is desirable to impose prior restrictions on the model specification, which cannot be achieved using conjugate priors.

Second, natural conjugate priors do not allow for flexible extensions, such as having different regressors in different equations. Third, the prior structure implies that the prior covariance of the coefficients in any two equations must be proportional to each other. This is because the prior covariance form is  $\boldsymbol{\Sigma} \otimes \mathbf{V}_0$ . However, this does not always make sense in certain applications. For example, imposing zero prior restrictions on some coefficients would imply that the prior variance of these coefficients should be near zero, but this does not need to be true for all coefficients in the model.

To address the first issue, we can think of the VAR(p) specification in a similar way to the seemingly unrelated regression (SUR) model, where we

have different regressors in different equations and account for unobserved dependence. This approach allows us to impose zero restrictions on the VAR(p) model, thereby improving its parsimony. Following the setup in Section 7.2, we have

$$\mathbf{y}_m = \mathbf{Z}_m \boldsymbol{\beta}_m + \boldsymbol{\mu}_m,$$

where  $\mathbf{y}_m$  is a  $T$ -dimensional vector corresponding to the  $m$ -th time series variable,  $\mathbf{Z}_m$  is a  $T \times K_m$  matrix of regressors,  $\boldsymbol{\beta}_m$  is a  $K_m$ -dimensional vector of location parameters, and  $\boldsymbol{\mu}_m$  is a  $T$ -dimensional vector of stochastic errors, for  $m = 1, 2, \dots, M$ .

Stacking the  $M$  equations, we can write  $\mathbf{y} = \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\mu}$  where  $\mathbf{y} = [\mathbf{y}_1^\top \mathbf{y}_2^\top \dots \mathbf{y}_M^\top]^\top$  is a  $MT$ -dimensional vector,  $\boldsymbol{\beta} = [\boldsymbol{\beta}_1^\top \boldsymbol{\beta}_2^\top \dots \boldsymbol{\beta}_M^\top]^\top$  is a  $K$  dimensional vector,  $K = \sum_{m=1}^M K_m$ , note that having the same number of regressors implies  $K = M \cdot k$  coefficients,  $\mathbf{Z}$  is an  $MT \times K$  block diagonal matrix composed of  $\mathbf{Z}_m$ , that is,

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{Z}_M \end{bmatrix},$$

and  $\boldsymbol{\mu} = [\boldsymbol{\mu}_1^\top \boldsymbol{\mu}_2^\top \dots \boldsymbol{\mu}_M^\top]^\top$  is a  $MT$ -dimensional vector of stochastic errors such that  $\boldsymbol{\mu} \sim N(\mathbf{0}, \boldsymbol{\Sigma} \otimes \mathbf{I}_T)$ .

We can use independent priors in this model to overcome the limitations of the conjugate prior, that is,  $\pi(\boldsymbol{\beta}) \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$  and  $\pi(\boldsymbol{\Sigma}^{-1}) \sim W(\alpha_0, \boldsymbol{\Psi}_0)$ . Thus, we know from Section 7.2 that the posterior distributions are

$$\boldsymbol{\beta} | \boldsymbol{\Sigma}, \mathbf{y}, \mathbf{Z} \sim N(\boldsymbol{\beta}_n, \mathbf{B}_n),$$

$$\boldsymbol{\Sigma}^{-1} | \boldsymbol{\beta}, \mathbf{y}, \mathbf{Z} \sim W(\alpha_n, \boldsymbol{\Psi}_n),$$

where  $\mathbf{B}_n = (\mathbf{Z}^\top (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I}_T) \mathbf{Z} + \mathbf{B}_0^{-1})^{-1}$ ,  $\boldsymbol{\beta}_n = \mathbf{B}_n (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{Z}^\top (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I}_T) \mathbf{y})$ ,  $\alpha_n = \alpha_0 + T$  and  $\boldsymbol{\Psi}_n = (\boldsymbol{\Psi}_0^{-1} + \mathbf{U}^\top \mathbf{U})^{-1}$ , where  $\mathbf{U}$  is an  $T \times M$  matrix whose columns are  $\mathbf{y}_m - \mathbf{Z}_m \boldsymbol{\beta}_m$ .<sup>5</sup>

Observe that we have standard conditional posteriors, thus, we can employ a Gibbs sampling algorithm to get the posterior draws. We can calculate the prediction  $\mathbf{y}_{T+1} = [y_{1T+1} \ y_{2T+1} \ \dots \ y_{MT+1}]^\top$  knowing that  $\mathbf{y}_{T+1} \sim N(\mathbf{Z}_T \boldsymbol{\beta}, \boldsymbol{\Sigma})$ , where

$$\mathbf{Z}_T = \begin{bmatrix} \mathbf{z}_{1T}^\top & 0 & \dots & 0 \\ 0 & \mathbf{z}_{2T}^\top & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{z}_{MT}^\top \end{bmatrix},$$

---

<sup>5</sup>We can also use the alternative representation presented in Section 7.2.

and using the posterior draws of  $\beta^{(s)}$  and  $\Sigma^{(s)}$ ,  $s = 1, 2, \dots, S$ . We can also perform inference of functions of the parameters that are of main interest when using VAR models.

Note that independent priors offer more flexibility regarding prior information. For instance, we can set  $\Psi_0 = \mathbf{S}^{-1}$ ,  $\alpha_0 = T$ ,  $\beta_0 = \mathbf{0}$ , and  $\mathbf{B}_0$  as a diagonal matrix, where the variance of the components associated with the coefficients in the  $m$ -th equation is such that the prior variance of the coefficients for the own lags is  $a_1/l^2$ , the variances for lag  $l$  of variable  $m \neq j$  are  $a_2 s_m^2/(l^2 s_j^2)$ , and the variance of the intercepts is set to  $a_3 s_m^2$ , with  $l = 1, 2, \dots, p$ , where  $s_m$  is the estimated standard error of the residuals from an unrestricted univariate autoregression of variable  $m$  against a constant and its  $p$  lags [274, 248].

Note that setting  $a_1 > a_2$  implies that own lags are more important as predictors than lags of other variables, and dividing by  $l^2$  implies that more recent lags are more relevant than those further in the past. The specific choices of  $a_1$ ,  $a_2$ , and  $a_3$  ( $a_k > 0$ ,  $k = 1, 2, 3$ ) depend on the specific application, but it is generally easier to elicit these parameters rather than the  $K(K+1)/2$  different components of  $\mathbf{B}_0$ .<sup>6</sup>

This setting is known as the *Minnesota prior*, as it is based on the seminal proposals for Bayesian VAR models by researchers at the University of Minnesota and the Federal Reserve Bank of Minneapolis [123, 274].<sup>7</sup>

An important non-linear function of parameters when performing VAR analysis is the *impulse response* function, which is, the response of one variable to an impulse in another variable in the model. The *impulse response* function can be deduced using the *MA* representation of the VAR model. In particular, we can write Equation 8.10 using the lag operator (see Section 8.2),

$$\mathbf{y}_t = \mathbf{v} + (\mathbf{A}_1 L + \mathbf{A}_2 L^2 + \dots + \mathbf{A}_p L^p) \mathbf{y}_t + \boldsymbol{\mu}_t, \quad (8.11)$$

thus  $\mathbf{A}(L)\mathbf{y}_t = \mathbf{v} + \boldsymbol{\mu}_t$ , where  $\mathbf{A}(L) = \mathbf{I}_M - \mathbf{A}_1 L - \mathbf{A}_2 L^2 - \dots - \mathbf{A}_p L^p$ . Let  $\Phi(L) := \sum_{s=0}^{\infty} \Phi_s L^s$  an operator such that  $\Phi(L)\mathbf{A}(L) = \mathbf{I}_M$ . Thus, we have that  $\Phi(L)\mathbf{A}(L)\mathbf{y}_t = (\sum_{s=0}^{\infty} \Phi_s L^s) \mathbf{v} + (\sum_{s=0}^{\infty} \Phi_s L^s) \boldsymbol{\mu}_t = \boldsymbol{\mu} + \sum_{s=0}^{\infty} \Phi_s \boldsymbol{\mu}_{t-s}$ . Note that  $L^s \mathbf{v} = \mathbf{v}$  because  $\mathbf{v}$  is constant, thus we set  $\sum_{s=0}^{\infty} \Phi_s L^s \mathbf{v} = \sum_{s=0}^{\infty} \Phi_s \mathbf{v} = \Phi(1) \mathbf{v} = (\mathbf{I}_M - \mathbf{A}_1 - \mathbf{A}_2 - \dots - \mathbf{A}_p)^{-1} \mathbf{v} := \boldsymbol{\mu}$ , which is the mean of the process [203, Chap. 2]. Therefore, the MA representation of the VAR is

$$\mathbf{y}_t = \boldsymbol{\mu} + \sum_{s=0}^{\infty} \Phi_s \boldsymbol{\mu}_{t-s}, \quad (8.12)$$

where  $\Phi_0 = \mathbf{I}_M$ , and we can get the coefficients in  $\Phi_s$  by the recursion  $\Phi_s =$

<sup>6</sup>In our GUI, we use the *bvar tools* package, which adopts a slightly different notation such that  $a_2 = a_1 \kappa_2$  and  $a_3 = a_1 \kappa_3$ . Thus, we need to set  $a_1$ ,  $\kappa_2$ , and  $\kappa_3$ .

<sup>7</sup>In the case that the variables are not stationary, which is more likely when using variables in levels (e.g., gross domestic product), we set  $\beta_0 = \mathbf{0}$ , except for the elements associated with the first own lags of the dependent variables in each equation, where the prior mean is set to 1. Additionally, the original proposal of the Minnesota prior set  $\Sigma = \mathbf{S}/T$ , meaning it did not account for uncertainty regarding  $\Sigma$ .

$\sum_{l=1}^s \Phi_{s-l} \mathbf{A}_l$ ,  $\mathbf{A}_l = \mathbf{0}$ ,  $l > p$  and  $s = 1, 2, \dots$  [203, Chap. 2]. This impulse response function is called *forecast error impulse response function*.

The MA coefficients contain the impulse responses of the system. In particular,  $\phi_{mj,s}$ , which is the  $mj$ -th element of the matrix  $\Phi_s$ , represents the response of the  $m$ -th variable to a unit shock of the variable  $j$  in the system,  $s$  periods ago, provided that the effect is not contaminated by other shocks in the system. The long-term effects (total multipliers) are given by  $\Psi_\infty := \sum_{s=1}^\infty \Phi_s = (\mathbf{I}_M - \mathbf{A}_1 - \mathbf{A}_2 - \dots - \mathbf{A}_p)^{-1}$ .

An assumption in these *impulse response* functions is that a shock occurs in only one variable at a time. This can be questionable as different shocks may be correlated, consequently, occurring simultaneously. Thus, the *impulse response* analysis can be performed based on the alternative MA representation,  $\mathbf{y}_t = \boldsymbol{\mu} + \sum_{s=0}^\infty \Phi_s \mathbf{P} \mathbf{P}^{-1} \boldsymbol{\mu}_{t-s} = \boldsymbol{\mu} + \sum_{s=0}^\infty \Theta_s \mathbf{w}_{t-s}$ , where  $\Theta_s = \Phi_s \mathbf{P}$  and  $\mathbf{w}_t = \mathbf{P}^{-1} \boldsymbol{\mu}_t$ ,  $\mathbf{P}$  is a lower triangular matrix such that  $\Sigma = \mathbf{P} \mathbf{P}^\top$  (Cholesky factorization/decomposition). Note that the covariance matrix of  $\mathbf{w}_t$  is  $\mathbf{I}_M$  due to  $\mathbb{E}[\mathbf{w}_t \mathbf{w}_t^\top] = \mathbb{E}[\mathbf{P}^{-1} \boldsymbol{\mu}_t \boldsymbol{\mu}_t^\top (\mathbf{P}^{-1})^\top] = \mathbf{P}^{-1} \Sigma (\mathbf{P}^{-1})^\top = \mathbf{P}^{-1} \mathbf{P} \mathbf{P}^\top (\mathbf{P}^{-1})^\top = \mathbf{I}_M$ .

In this representation is sensible to assume that each shock occurs independently due to the covariance matrix of  $\mathbf{w}_t$  being an identity. In addition, a unit shock is a shock of size one standard deviation due the result of the covariance matrix. This is named the *orthogonalized impulse response*, where  $\theta_{mj,s}$ , which is the  $mj$ -th element of the matrix  $\Theta_s$ , represents the response of the  $m$ -th variable to a standard deviation shock of the variable  $j$  in the system,  $s$  periods ago. The critical point with the orthogonalized impulse responses is that the order of the variables in the VAR is really important because implicitly establishes a recursive model, that is, the  $m$ -th equation in the system may contain  $y_{1t}, y_{2t}, \dots, y_{m-1t}$ , but not  $y_{mt}, y_{m+1t}, \dots, y_{Mt}$  on the right-hand side of its equation. Thus,  $y_{mt}$  cannot have an instantaneous impact on  $y_{jt}$  for  $j < m$  [203, Chap. 2].

Beyond the fascinating macroeconomic implications embedded in the specification of VAR models, the key point for this section is that we can infer impulse response functions using the posterior draws.

### Example: US fiscal system

Let's use the dataset provided by [459] of the US fiscal system, where  $ttr$  is the quarterly total tax revenue,  $gs$  is the quarterly total government spending, and  $gdp$  is the quarterly gross domestic product, all expressed in log, real, per person terms, and the period is 1948q1 to 2024q2. This dataset is the *18USAfiscal.csv* file. [307] analyze the US fiscal policy shocks using these variables.

Let's estimate a VAR model where  $\mathbf{y}_t = [\Delta(ttr_t) \Delta(gs_t) \Delta(gdp_t)]^\top$ , that is, we work with the log differences (variation rates), and we set  $p = 1$ . We use the package *bvar tools* to estimate the *forecast error* and *orthogonalized* impulse response functions. We use vague independent priors setting  $\beta_0 = \mathbf{0}$ ,

$\mathbf{B}_0 = 100\mathbf{I}$ ,  $\mathbf{V}_0 = 5^{-1}\mathbf{I}$  and  $\alpha_0 = 3$ , and the Minnesota prior setting  $a_1 = 2$ ,  $\kappa_2 = 0.5$  and  $\kappa_3 = 5$  (default values).<sup>8</sup>

The following code shows how to do this, take into account that we use the first 301 observations to estimate the model, and keep the last 4 observations to check the forecasting performance. Figures 8.6 and 8.7 show the impulse response functions of  $gs$  with respect to  $gs$ , the *forecast error impulse response* using vague independent priors, and the *orthogonalized impulse response* using the Minnesota prior, respectively. We see that the effect of the Minnesota prior is to decrease uncertainty. In addition, the forecasting exercise results indicate that these assumptions have same effects in this example. In particular, Figure 8.8 shows that the mean forecasts using the vague prior (green line) and the Minnesota prior (red line) are indistinguishable from the true observations (black line). However, the Minnesota prior enhances forecast precision, as its 95% predictive interval (blue shaded area) is narrower and fully contained within the 95% predictive interval obtained using vague priors (light blue shaded area). This improvement is attributable to the shrinkage properties of the Minnesota prior.

### R code. VAR model: US fiscal shocks

```

1 rm(list = ls()); set.seed(010101)
2 DataUSfilcal <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTAApp/refs/heads/master/DataApp/18USAfiscal.
  csv", sep = ",", header = TRUE, quote = "")
3 attach(DataUSfilcal) # upload data
4 Y <- cbind(diff(as.matrix(DataUSfilcal[,-c(1:2)])))
5 T <- dim(Y)[1]-1; K <- dim(Y)[2]
6 Ynew <- Y[-c((T-2):(T+1)), ] # Use 4 last observations to
  check forecast
7 y1 <- Ynew[-1, 1]; y2 <- Ynew[-1, 2]; y3 <- Ynew[-1, 3]
8 X1 <- cbind(1, lag(Ynew)); X1 <- X1[-1,]
9 X2 <- cbind(1, lag(Ynew)); X2 <- X2[-1,]
10 X3 <- cbind(1, lag(Ynew)); X3 <- X3[-1,]
11 M <- dim(Y)[2]; K1 <- dim(X1)[2]; K2 <- dim(X2)[2]; K3 <-
  dim(X3)[2]
12 K <- K1 + K2 + K3
13 # Hyperparameters
14 b0 <- 0; c0 <- 100; V0 <- 5^(-1); a0 <- M
15 #Posterior draws
16 MCMC <- 10000; burnin <- 1000; H <- 10; YnewPack <- ts(Ynew)

```

<sup>8</sup>The *bvarTools* package uses the inverse Wishart distribution as prior for  $\Sigma$ , where the hyperparameters are the degrees of freedom of the error term, and the prior error variance of endogenous variables.

### R code. VAR model: US fiscal shocks

```

1 model <- bvartools::gen_var(YnewPack, p = 1, deterministic =
2   "const", iterations = MCMC, burnin = burnin) # Create
3   model
4 model <- bvartools::add_priors(model, coef = list(v_i = c0
5     ^-1, v_i_det = c0^-1, const = b0), sigma = list(df = a0,
6     scale = V0/a0), coint_var = FALSE) # Add priors
7 object <- bvartools::draw_posterior(model) # Posterior draws
8 ir <- bvartools::irf.bvar(object, impulse = "gs", response =
9   "gs", n.ahead = H, type = "feir", cumulative = FALSE) # Calculate IR
10 # Plot IR
11 plot_IR <- function(df) {
12   p <- ggplot(data = df, aes(x = t)) + geom_ribbon(aes(ymin =
13     lower, ymax = upper), alpha = 1, fill = "lightblue") +
14     geom_line(aes(y = mean), colour = "blue", linewidth =
15     0.5) + ylab("Impulse response") + xlab("Time") + xlim(0,
16     H)
17   print(p)
18 }
19 dfNew <- tibble(t = 0:H, mean = as.numeric(ir[,2]), lower =
20   as.numeric(ir[,1]), upper = as.numeric(ir[,3]))
21 FigNew <- plot_IR(dfNew)
22 # Using Minnesota prior
23 library(tibble)
24 modelMin <- bvartools::gen_var(YnewPack, p = 1,
25   deterministic = "const", iterations = MCMC, burnin =
26   burnin)
27 modelMin <- bvartools::add_priors(modelMin, minnesota = list
28   (kappa0 = 2, kappa1 = 0.5, kappa3 = 5), coint_var =
29   FALSE) # Minnesota prior
30 objectMin <- bvartools::draw_posterior(modelMin) # Posterior
31   draws
32 irMin <- bvartools::irf.bvar(objectMin, impulse = "gs",
33   response = "gs", n.ahead = H, type = "feir", cumulative =
34   FALSE) # Calculate IR
35 dfNewMin <- tibble(t = 0:H, mean = as.numeric(irMin[,2]),
36   lower = as.numeric(irMin[,1]), upper = as.numeric(irMin
37   [,3]))
38 FigNewMin <- plot_IR(dfNewMin)
39 #### Forecasting
40 bvar_pred <- predict(object, n.ahead = 4, new_d = rep(1, 4))
41 bvar_predOR <- predict(objectMin, n.ahead = 4, new_d = rep
42   (1, 4))
43 dfFore <- tibble(t = c((T-2):(T+1)), mean = as.numeric(bvar_
44   pred[["fcst"]][["gs"]][,2]), lower = as.numeric(bvar_
45   pred[["fcst"]][["gs"]][,1]), upper = as.numeric(bvar_
46   pred[["fcst"]][["gs"]][,3]), mean1 = as.numeric(bvar_
47   predOR[["fcst"]][["gs"]][,2]), lower1 = as.numeric(bvar_
48   predOR[["fcst"]][["gs"]][,1]), upper1 = as.numeric(bvar_
49   predOR[["fcst"]][["gs"]][,3]), true = as.numeric(Y[c((T
50   -2):(T+1)),2]))

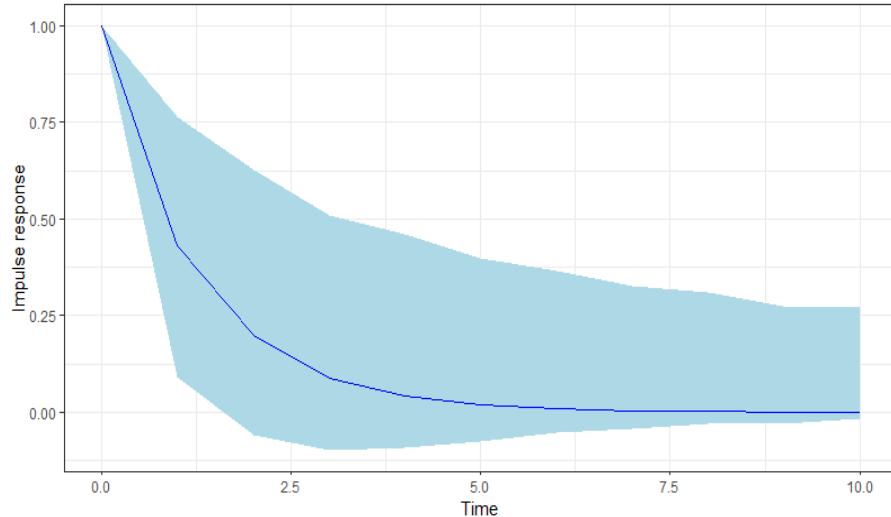
```

### R code. VAR model: US fiscal shocks

```

1 plot_FORE <- function(df) {
2   p <- ggplot(data = dfFore, aes(x = t)) + geom_ribbon(aes(
3     ymin = lower1, ymax = upper1), alpha = 1, fill = "blue")
4   + geom_ribbon(aes(ymin = lower, ymax = upper), alpha =
5     1, fill = "lightblue") + geom_line(aes(y = mean), colour =
6     "green", linewidth = 0.5) + geom_line(aes(y = mean1),
7     colour = "red", linewidth = 0.5) + geom_line(aes(y =
8     true), colour = "black", linewidth = 0.5) + ylab("Forecast")
9   + xlab("Time") + xlim(c((T-2),(T+1)))
10  print(p)
11 }
12 FigFore <- plot_FORE(dfFore)

```

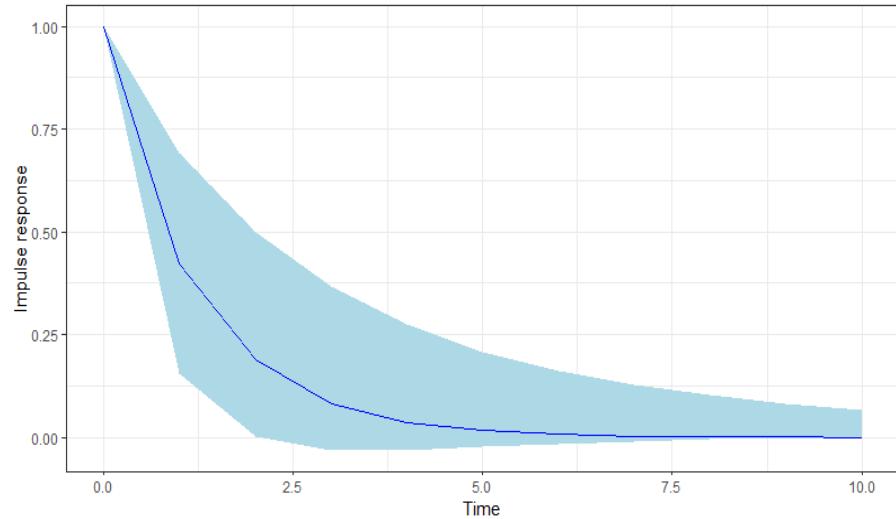


**FIGURE 8.6**

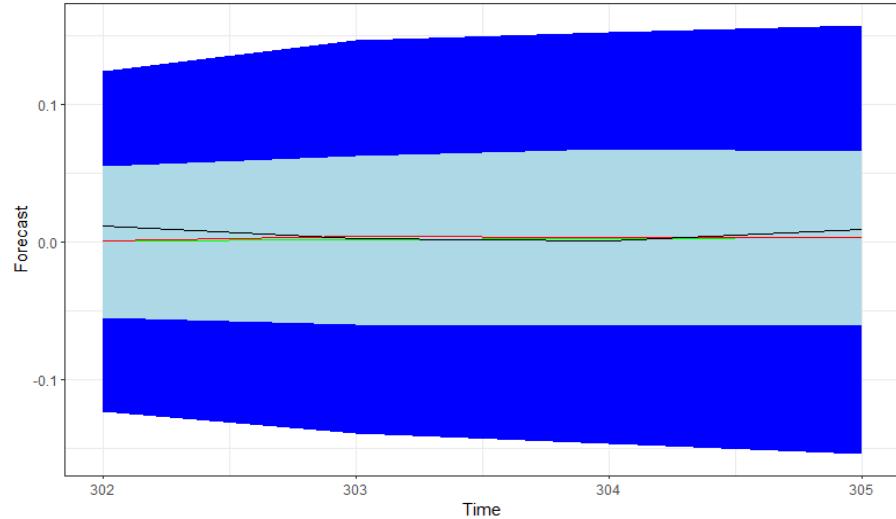
Forecasting error impulse response: the government spending ( $gs$ ) with respect to the government spending ( $gs$ ).

Algorithm A24 shows how to do perform inference in VAR models using our GUI. See also Chapter 5 for details regarding the dataset structure.

There are other good packages in **R** to perform Bayesian inference in VAR models. For instance, *bayesianVARs* package implements inference of reduced-form VARs with stochastic volatility [281], *BVAR* package performs inference

**FIGURE 8.7**

Ortogonalized impulse response: the government spending ( $gs$ ) with respect to the government spending ( $gs$ ).

**FIGURE 8.8**

Forecast performance: the government spending ( $gs$ ).

using hierarchical priors [324], *bvarsrv* implements time-varying parameters

**Algorithm A24** Vector Autoregressive models

- 
- 1: Select *Time series Model* on the top panel
  - 2: Select *VAR models* using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Set the number of lags ( $p$ )
  - 6: Set the hyperparameters for the Minnesota prior:  $a_1$ ,  $\kappa_2$  and  $\kappa_3$ . This step is not necessary as by default our GUI uses default values in *bvar tools* package
  - 7: Select the type of *impulse response functions*: forecast error or orthogonalized, and ordinary or cumulative.
  - 8: Set the time horizon for the impulse response functions and the forecasts
  - 9: Click the *Go!* button
  - 10: Analyze results
  - 11: Download impulse responses and forecasts using the *Download Results* button
- 

models [252], *bsvars* performs estimation of structural VAR models [459], and *bsvarSIGNs* to estimating structural VAR models with sign restrictions [461].

---

## 8.5 Summary

We present a brief review of Bayesian inference in time series models. In particular, we introduce the *state-space* representation and demonstrate how to perform inferential analysis for these models, focusing on the dynamic linear model and the stochastic volatility model. Additionally, we show how  $ARMA(p, q)$  processes can be expressed in *state-space* form and provide methods for estimating such models.

We also include code for implementing computational inference algorithms, such as sequential Monte Carlo (SMC), Hamiltonian Monte Carlo (HMC), and various Markov chain Monte Carlo (MCMC) methods. Finally, we introduce VAR( $p$ ) models, detailing how to perform impulse-response analysis and forecasting within this framework.

Time series analysis is a highly active research area with remarkable methodological developments and applications. Interested readers can refer to excellent materials in chapters 7 and 9 of [174], and chapters 17 to 20 of [76], along with the references therein.

---

## 8.6 Exercises

1. Simulate the *dynamic linear model* assuming  $X_t \sim N(1, 0.1\sigma^2)$ ,  $w_t \sim N(0, 0.5\sigma^2)$ ,  $\mu_t \sim N(0, \sigma^2)$ ,  $\beta_0 = 1$ ,  $B_0 = 0.5\sigma^2$ ,  $\sigma^2 = 0.25$ , and  $G_t = 1$ ,  $t = 1, \dots, 100$ . Then, perform the filtering recursion fixing  $\Sigma = 25 \times 0.25$ ,  $\Omega_1 = 0.5\Sigma$  (high signal-to-noise ratio) and  $\Omega_2 = 0.1\Sigma$  (low signal-to-noise ratio). Plot and compare the results.
2. Simulate the *dynamic linear model*  $y_t = \beta_t x_t + \mu_t$ ,  $\beta_t = \beta_{t-1} + w_t$ , where  $x_t \sim N(1, 0.1\sigma^2)$ ,  $w_t \sim N(0, 0.5\sigma^2)$ ,  $\mu_t \sim N(0, \sigma^2)$ ,  $\beta_0 = 0$ ,  $B_0 = 0.5\sigma^2$ , and  $\sigma^2 = 1$ ,  $t = 1, \dots, 100$ . Perform the filtering and smoothing recursions from scratch.
3. Simulate the process  $y_t = \alpha z_t + \beta_t x_t + \mathbf{h}^\top \boldsymbol{\epsilon}_t$ ,  $\beta_t = \beta_{t-1} + \mathbf{H}^\top \boldsymbol{\epsilon}_t$ , where  $\mathbf{h}^\top = [1 \ 0]$ ,  $\mathbf{H}^\top = [0 \ 1/\tau]$ ,  $\boldsymbol{\epsilon}_t \sim N(\mathbf{0}_2, \sigma^2 \mathbf{I}_2)$ ,  $x_t \sim N(1, 2\sigma^2)$ ,  $z_t \sim N(0, 2\sigma^2)$ ,  $\alpha = 2$ ,  $\tau^2 = 5$  and  $\sigma^2 = 0.1$ ,  $t = 1, \dots, 200$ . Assume  $\pi(\beta_0, \alpha, \sigma^2, \tau) = \pi(\beta_0)\pi(\alpha)\pi(\sigma^2)\pi(\tau^2)$  where  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ ,  $\tau^2 \sim G(v_0/2, v_0/2)$ ,  $\alpha \sim N(a_0, A_0)$  and  $\beta_0 \sim N(b_0, B_0)$  such that  $\alpha_0 = \delta_0 = 1$ ,  $v_0 = 5$ ,  $a_0 = 0$ ,  $A_0 = 1$ ,  $b_0 = 0$ ,  $B_0 = \sigma^2/\tau^2$ . Program the MCMC algorithm including the *simulation smoother*.
4. Show that the posterior distribution of  $\phi \mid \beta, \sigma^2, \mathbf{y}, \mathbf{X}$  in the model  $y_t = \mathbf{x}_t^\top \beta + \mu_t$  where  $\phi(L)\mu_t = \epsilon_t$  and  $\epsilon_t \stackrel{iid}{\sim} N(0, \sigma^2)$  is  $N(\phi_n, \Phi_n) \mathbb{1}(\phi \in S_\phi)$ , where  $\Phi_n = (\Phi_0^{-1} + \sigma^{-2} \mathbf{U}^\top \mathbf{U})$ ,  $\phi_n = \Phi_n(\Phi_0^{-1} \phi_0 + \sigma^{-2} \mathbf{U}^\top \boldsymbol{\mu})$ , and  $S_\phi$  is the stationary region of  $\phi$ .
5. Show that in the *AR(2)* stationary process,  $y_t = \mu + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \epsilon_t$ , where  $\epsilon_t \sim N(0, \sigma^2)$ ,  $\mathbb{E}[y_t] = \frac{\mu}{1-\phi_1-\phi_2}$ , and  $Var[y_t] = \frac{\sigma^2(1-\phi_2)}{1-\phi_2-\phi_1^2-\phi_1\phi_2-\phi_2^2+\phi_2^3}$ .
6. Program a Hamiltonian Monte Carlo taking into account the stationary restrictions on  $\phi_1$  and  $\phi_2$ , and  $\epsilon_0$  such that the acceptance rate is near 65%.
7. **Stochastic volatility model**
  - Program a sequential importance sampling (SIS) from scratch in the vanilla stochastic volatility model setting  $\mu = -10$ ,  $\phi = 0.95$ ,  $\sigma = 0.3$  and  $T = 250$ . Check what happen with its performance.
  - Modify the sequential Monte Carlo (SMC) to perform multinomial resampling when the effective sample size is lower than 50% the initial number of particles.
8. Estimate the vanilla stochastic volatility model using the dataset *17ExcRate.csv*, provided by [358], which contains the exchange rate log daily returns for USD/EUR, USD/GBP, and GBP/EUR from

one year before and after the WHO declared the COVID-19 pandemic on 11 March 2020.

9. Simulate the VAR(1) process

$$\begin{bmatrix} y_{1t} \\ y_{2t} \\ y_{3t} \end{bmatrix} = \begin{bmatrix} 2.8 \\ 2.2 \\ 1.3 \end{bmatrix} + \begin{bmatrix} 0.5 & 0 & 0 \\ 0.1 & 0.1 & 0.3 \\ 0 & 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} y_{1t-1} \\ y_{2t-1} \\ y_{3t-1} \end{bmatrix} + \begin{bmatrix} \mu_{1t} \\ \mu_{2t} \\ \mu_{3t} \end{bmatrix},$$

where  $\Sigma = \begin{bmatrix} 2.25 & 0 & 0 \\ 0 & 1 & 0.5 \\ 0 & 0.5 & 0.74 \end{bmatrix}$ .

- Use vague independent priors setting  $\beta_0 = \mathbf{0}$ ,  $B_0 = 100\mathbf{I}$ ,  $V_0 = 5\mathbf{I}$  and  $\alpha_0 = 3$ , and estimate a VAR(1) model using the *rsurGibbs* function from the package *bayesm*. Then, program from scratch algorithms to perform inference of the *forecast error* and *orthogonalized* impulse response functions, and compare with the population impulse response functions, that is, using the population parameters.
- Using the previous setting perform inference of the *forecast error* and the *orthogonalized* impulse responses using the package *bvarTools*.



# 9

## *Longitudinal/Panel data models*

We describe how to perform inference in longitudinal/panel models using a Bayesian framework. In this context, multiple cross-sectional units are observed repeatedly over time, a structure referred to as panel data by econometricians and longitudinal data by statisticians. Specifically, we present models for continuous (normal), binary (logit), and count (Poisson) responses. Applications and exercises illustrate the potential of these models.

In longitudinal/panel data sets, we have  $y_{it}$  where  $i = 1, 2, \dots, N$  and  $t = 1, 2, \dots, T_i$ . If  $T_i = T$  for all  $i$ , the dataset is *balanced*; otherwise, it is *unbalanced*. Panel data typically involves by far more cross-sectional units than time periods, this is called typically a *short panel*. It assumes that cross-sectional units are independent, though serial correlation exists within each unit over time, and unobserved heterogeneity for each unit must be accounted for. We can treat this unobserved heterogeneity as random variables, assuming it is either independent or dependent on control variables. Econometricians refer to these cases as *random effects* and *fixed effects*, respectively. The Bayesian literature takes a different approach, modeling the panel structure hierarchically, where the unobserved heterogeneity may or may not depend on other controls.<sup>1</sup>

Remember that the easiest way to run our GUI is typing

### **R code. How to display our graphical user interface**

```
1 shiny::runGitHub("besmarter/BSTApp", launch.browser = T)
```

in the **R** package console or any **R** code editor, and once our GUI is deployed, select *Hierarchical Longitudinal Models*. However, users should see Chapter 5 for details.

<sup>1</sup>See [367] for a nice comparison of Frequentist and Bayesian treatments of panel data models.

## 9.1 Normal model

The longitudinal (panel) normal model is specified as

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{W}_i\mathbf{b}_i + \boldsymbol{\mu}_i,$$

where  $\mathbf{y}_i$  is a  $T_i$ -dimensional vector corresponding to unit  $i = 1, 2, \dots, N$ , and  $\mathbf{X}_i$  and  $\mathbf{W}_i$  are  $T_i \times K_1$  and  $T_i \times K_2$  design matrices, respectively. In the statistical literature,  $\boldsymbol{\beta}$  denotes a  $K_1$ -dimensional vector of *fixed effects*, whereas  $\mathbf{b}_i$  is a  $K_2$ -dimensional vector of unit-specific *random effects* that allow for unit-specific means and capture marginal dependence among observations within each cross-sectional unit. We assume normally distributed stochastic errors,  $\boldsymbol{\mu}_i \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_{T_i})$ , implying that the likelihood function is

$$\begin{aligned} p(\boldsymbol{\beta}, \mathbf{b}, \sigma^2 \mid \mathbf{y}, \mathbf{X}, \mathbf{W}) &\propto \prod_{i=1}^N |\sigma^2 \mathbf{I}_{T_i}|^{-1/2} \\ &\quad \times \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{W}_i\mathbf{b}_i)^\top (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{W}_i\mathbf{b}_i) \right\} \\ &= (\sigma^2)^{-\frac{1}{2} \sum_{i=1}^N T_i} \\ &\quad \times \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{W}_i\mathbf{b}_i)^\top (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{W}_i\mathbf{b}_i) \right\}, \end{aligned}$$

where  $\mathbf{b} = [\mathbf{b}_1^\top, \mathbf{b}_2^\top, \dots, \mathbf{b}_N^\top]^\top$ .

Panel data modeling in the Bayesian approach assumes a hierarchical structure in the *random effects*. Following [87], there is a first stage where  $\mathbf{b}_i \sim N(\mathbf{0}, \mathbf{D})$ ,  $\mathbf{D}$  allows serial correlation within each cross-sectional unit  $i$ , and then, there is a second stage where  $\mathbf{D} \sim IW(d_0, d_0 \mathbf{D}_0)$ . Thus, we can see that there is an additional layer of priors as there is a prior on the hyperparameter  $\mathbf{D}$ .

In addition, we have standard conjugate prior distributions for  $\boldsymbol{\beta}$  and  $\sigma^2$ ,  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$  and  $\sigma^2 \sim IG(\alpha_0, \delta_0)$ .

[87] propose a blocking algorithm to perform inference in longitudinal hierarchical models by considering the distribution of  $\mathbf{y}_i$  marginalized over the random effects. Given that  $\mathbf{y}_i \mid \boldsymbol{\beta}, \mathbf{b}_i, \sigma^2, \mathbf{X}_i, \mathbf{W}_i \sim N(\mathbf{X}_i\boldsymbol{\beta} + \mathbf{W}_i\mathbf{b}_i, \sigma^2 \mathbf{I}_{T_i})$ , we can see that  $\mathbf{y}_i \mid \boldsymbol{\beta}, \mathbf{D}, \sigma^2, \mathbf{X}_i, \mathbf{W}_i \sim N(\mathbf{X}_i\boldsymbol{\beta}, \mathbf{V}_i)$ , where  $\mathbf{V}_i = \sigma^2 \mathbf{I}_{T_i} + \mathbf{W}_i \mathbf{D} \mathbf{W}_i^\top$  given that  $\mathbb{E}[\mathbf{b}_i] = \mathbf{0}$  and  $Var[\mathbf{b}_i] = \mathbf{D}$ . If we have just random intercepts, then  $\mathbf{W}_i = \mathbf{i}_{T_i}$ , where  $\mathbf{i}_{T_i}$  is a  $T_i$ -dimensional vector of ones. Thus,  $\mathbf{V}_i = \sigma^2 \mathbf{I}_{T_i} + \sigma_b^2 \mathbf{i}_{T_i} \mathbf{i}_{T_i}^\top$ , the variance is  $\sigma^2 + \sigma_b^2$  and the covariance is  $\sigma_b^2$  within each cross-sectional unit through time.

We can deduce the posterior distribution of  $\beta$  given  $\sigma^2$  and  $D$ ,

$$\begin{aligned}\pi(\beta | \sigma^2, D, y, X, W) &\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \beta)^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \beta) \right\} \\ &\quad \times \exp \left\{ -\frac{1}{2} (\beta - \beta_0)^T \mathbf{B}_0^{-1} (\beta - \beta_0) \right\}.\end{aligned}$$

This implies that (see Exercise 1)

$$\beta | \sigma^2, D, y, X, W \sim N(\beta_n, \mathbf{B}_n),$$

where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \sum_{i=1}^N \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{X}_i)^{-1}$ ,  $\beta_n = \mathbf{B}_n (\mathbf{B}_0^{-1} \beta_0 + \sum_{i=1}^N \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{y}_i)$ .

We can use the likelihood  $p(\beta, b_i, \sigma^2 | y, X, W)$  to get the posterior distributions of  $b_i$ ,  $\sigma^2$  and  $D$ . In particular,

$$\begin{aligned}\pi(b_i | \beta, \sigma^2, D, y, X, W) &\propto \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \beta - \mathbf{W}_i b_i)^T (\mathbf{y}_i - \mathbf{X}_i \beta - \mathbf{W}_i b_i) \right\} \\ &\quad \times \exp \left\{ -\frac{1}{2} \sum_{i=1}^N b_i^T \mathbf{D}^{-1} b_i \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^N (-2b_i^T (\sigma^{-2} \mathbf{W}_i^T (\mathbf{y}_i - \mathbf{X}_i \beta)) \right. \\ &\quad \left. + b_i^T (\sigma^{-2} \mathbf{W}_i^T \mathbf{W}_i + \mathbf{D}^{-1}) b_i) \right\} \\ &\propto \exp \left\{ -\frac{1}{2} (-2b_i^T \mathbf{B}_{ni}^{-1} \mathbf{B}_{ni} (\sigma^{-2} \mathbf{W}_i^T (\mathbf{y}_i - \mathbf{X}_i \beta)) \right. \\ &\quad \left. + b_i^T \mathbf{B}_{ni}^{-1} b_i) \right\} \\ &= \exp \left\{ -\frac{1}{2} (-2b_i^T \mathbf{B}_{ni}^{-1} b_{ni} + b_i^T \mathbf{B}_{ni}^{-1} b_i) \right\},\end{aligned}$$

where  $\mathbf{B}_{ni} = (\sigma^{-2} \mathbf{W}_i^T \mathbf{W}_i + \mathbf{D}^{-1})^{-1}$  and  $b_{ni} = \mathbf{B}_{ni} (\sigma^{-2} \mathbf{W}_i^T (\mathbf{y}_i - \mathbf{X}_i \beta))$ .

We can complete the square in this expression by adding and subtracting  $b_{ni}^T \mathbf{B}_{ni}^{-1} b_{ni}$ . Thus,

$$\begin{aligned}\pi(b_i | \beta, \sigma^2, D, y, X, W) &\propto \exp \left\{ -\frac{1}{2} (-2b_i^T \mathbf{B}_{ni}^{-1} b_{ni} + b_i^T \mathbf{B}_{ni}^{-1} b_i + b_{ni}^T \mathbf{B}_{ni}^{-1} b_{ni} \right. \\ &\quad \left. - b_{ni}^T \mathbf{B}_{ni}^{-1} b_{ni}) \right\} \\ &\propto \exp \left\{ (b_i - b_{ni})^T \mathbf{B}_{ni}^{-1} (b_i - b_{ni}) \right\}.\end{aligned}$$

This is the kernel of a multivariate normal distribution with mean  $b_{ni}$  and variance  $\mathbf{B}_{ni}$ . Thus,

$$b_i | \beta, \sigma^2, D, y, X, W \sim N(b_{ni}, \mathbf{B}_{ni}),$$

Let's see the posterior distribution of  $\sigma^2$ ,

$$\begin{aligned}\pi(\sigma^2 | \boldsymbol{\beta}, \mathbf{b}, \mathbf{y}, \mathbf{X}, \mathbf{W}) &\propto (\sigma^2)^{-\frac{\sum_{i=1}^N T_i}{2}} \\ &\times \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{W}_i \mathbf{b}_i)^\top (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{W}_i \mathbf{b}_i) \right\} \\ &\times (\sigma^2)^{-\alpha_0 - 1} \exp \left\{ -\frac{\delta_0}{\sigma^2} \right\} \\ &= (\sigma^2)^{-\frac{\sum_{i=1}^N T_i}{2} - \alpha_0 - 1} \\ &\times \exp \left\{ -\left( \sum_{i=1}^N \frac{(\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{W}_i \mathbf{b}_i)^\top (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{W}_i \mathbf{b}_i)}{2} \right. \right. \\ &\quad \left. \left. + \delta_0 \right) \frac{1}{\sigma^2} \right\}.\end{aligned}$$

Thus,

$$\sigma^2 | \boldsymbol{\beta}, \mathbf{b}, \mathbf{y}, \mathbf{X}, \mathbf{W} \sim IG(\alpha_n, \delta_n),$$

where  $\alpha_n = \alpha_0 + \frac{1}{2} \sum_{i=1}^N T_i$  and  $\delta_n = \delta_0 + \frac{1}{2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{W}_i \mathbf{b}_i)^\top (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{W}_i \mathbf{b}_i)$ .

The posterior distribution of  $\mathbf{D}$  is the following,

$$\begin{aligned}\pi(\mathbf{D} | \mathbf{b}) &\propto |\mathbf{D}|^{-N/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \mathbf{b}_i^\top \mathbf{D}^{-1} \mathbf{b}_i \right\} \\ &\times |\mathbf{D}|^{-(d_0 + K_2 + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(d_0 \mathbf{D}_0 \mathbf{D}^{-1}) \right\} \\ &= |\mathbf{D}|^{-(d_0 + N + K_2 + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr} \left( \left( d_0 \mathbf{D}_0 + \sum_{i=1}^N \mathbf{b}_i \mathbf{b}_i^\top \right) \mathbf{D}^{-1} \right) \right\}.\end{aligned}$$

This is the kernel of an inverse Wishart distribution with degrees of freedom  $d_n = d_0 + N$  and scale matrix  $\mathbf{D}_n = d_0 \mathbf{D}_0 + \sum_{i=1}^N \mathbf{b}_i \mathbf{b}_i^\top$ . Thus,

$$\mathbf{D} | \mathbf{b} \sim IW(d_n, \mathbf{D}_n).$$

Observe that the posterior distribution of  $\mathbf{D}$  depends just on  $\mathbf{b}$ .

All the posterior conditional distributions belong to standard families, this implies that we can use a Gibbs sampling algorithm to perform inference in these hierarchical normal models.

#### **Example: The relation between productivity and public investment**

We used the dataset named *8PublicCap.csv* used by [353] to analyze the relation between public investment and gross state product in the setting

of a spatial panel dataset consisting of 48 US states from 1970 to 1986. In particular, we perform inference based on the following equation

$$\log(gsp_{it}) = b_i + \beta_1 + \beta_2 \log(pcap_{it}) + \beta_3 \log(pc_{it}) + \beta_4 \log(emp_{it}) + \beta_5 \text{unemp}_{it} + \mu_{it},$$

where gsp is the gross state product, pcap is public capital, and pc is private capital all in USD, emp is employment (people), and unemp is the unemployment rate in percentage.

Algorithm A25 shows how to perform inference in hierarchical longitudinal normal models in our GUI. See also Chapter 5 for details regarding the dataset structure.

---

**Algorithm A25** Hierarchical longitudinal normal models

---

- 1: Select *Hierarchical Longitudinal Model* on the top panel
  - 2: Select *Normal* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Write down the formula of the *fixed effects* equation in the **Main Equation: Fixed Effects** box. This formula must be written using the syntax of the *formula* command of **R** software. This equation includes intercept by default, do not include it in the equation
  - 6: Write down the formula of the *random effects* equation in the **Main Equation: Random Effects** box without writing the dependent variable, that is, starting the equation with the *tilde* (~) symbol. This formula must be written using the syntax of the *formula* command of **R** software. This equation includes intercept by default, do not include it in the equation. If there are just random intercepts do not write anything in this box
  - 7: Write down the name of the grouping variable, that is, the variable that indicates the cross-sectional units
  - 8: Set the hyperparameters of the *fixed effects*: mean vector, covariance matrix, shape and scale parameters. This step is not necessary as by default our GUI uses non-informative priors
  - 9: Set the hyperparameters of the *random effects*: degrees of freedom and scale matrix of the inverse Wishart distribution. This step is not necessary as by default our GUI uses non-informative priors
  - 10: Click the *Go!* button
  - 11: Analyze results
  - 12: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

We ask in Exercise 2 to run this application in our GUI using 10,000

MCMC iterations plus a burn-in equal to 5,000 iterations, and a thinning parameter equal to 1. We also used the default values for the hyperparameters of the prior distributions, that is,  $\beta_0 = \mathbf{0}_5$ ,  $B_0 = I_5$ ,  $\alpha_0 = \delta_0 = 0.001$ ,  $d_0 = 5$  and  $D_0 = I_1$ . It seems that all posterior draws come from stationary distributions, as suggested by the diagnostics and posterior plots (see Exercise 2).

The following code uses the command *MCMChregress* from the package *MCMCpack* to run this application. This command is also used by our GUI to perform inference in hierarchical longitudinal normal models.

We can see that the 95% symmetric credible intervals for public capital, private capital, employment, and unemployment are (-2.54e-02, -2.06e-02), (2.92e-01, 2.96e-01), (7.62e-01, 7.67e-01) and (-5.47e-03, -5.31e-03), respectively. The posterior mean elasticity estimate of public capital to gsp is -0.023, that is, an increase by 1% in public capital means a 0.023% decrease in gross state product. The posterior mean estimates of private capital and employment elasticities are 0.294 and 0.765, respectively. In addition, 1 percentage point increase in the unemployment rate means a decrease of 0.54% in gsp. It seems that all these variables are statistically relevant. In addition, the posterior mean estimates of the variance associated with the unobserved heterogeneity and stochastic errors are 1.06e-01 and 1.45e-03. We obtained the posterior chain of the proportion of the variance associated with the unobserved heterogeneity. The 95% symmetric credible interval is (0.98, 0.99) for this proportion, that is, unobserved heterogeneity is very important to explain the total variability.

### *R code. The relationship between productivity and public investment, MCMChregress command*

```

1 rm(list = ls()); set.seed(12345)
2 DataGSP <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/8PublicCap.
  csv", sep = ",", header = TRUE, quote = "")
3 attach(DataGSP)
4 K1 <- 5; K2 <- 1
5 b0 <- rep(0, K1); B0 <- diag(K1)
6 r0 <- 5; R0 <- diag(K2)
7 a0 <- 0.001; d0 <- 0.001
8 Resultshreg <- MCMCpack::MCMChregress(fixed = log(gsp)^log(
  pcap)+log(pc)+log(emp)+unemp, random = ~1, group = "id",
  data = DataGSP, burnin = 5000, mcmc = 10000, thin = 1,
  r = r0, R = R0, nu = a0, delta = d0)

```

**R code. The relationship between productivity and public investment, MCMChregress command**

```

1 Betas <- Resultshreg[["mcmc"]][,1:K1]
2 Sigma2RanEff <- Resultshreg[["mcmc"]][,54]
3 Sigma2 <- Resultshreg[["mcmc"]][,55]
4 summary(Betas)
5 summary(Sigma2RanEff)
6 summary(Sigma2RanEff/(Sigma2RanEff+Sigma2))
7

```

There are many extensions of this model, for instance, [87] propose to introduce heteroskedasticity in this model by assuming  $\mu_{it} | \tau_{it} \sim N(0, \sigma^2/\tau_{it})$ ,  $\tau_{it} \sim G(v/2, v/2)$ . We ask in Exercise 2 to perform inference in the relation between productivity and public investment example using this setting. Another potential extension is to allow dependence between  $\mathbf{b}_i$  and some controls, let's say  $\mathbf{z}_i$ , a  $K_3$ -dimensional vector, and assume  $\mathbf{b}_i \sim N(\mathbf{Z}_i\boldsymbol{\gamma}, \mathbf{D})$  where  $\mathbf{Z}_i = \mathbf{I}_{K_2} \otimes \mathbf{z}_i^\top$ , and complete the model using a prior for  $\boldsymbol{\gamma}$ ,  $\boldsymbol{\gamma} \sim N(\boldsymbol{\gamma}_0, \boldsymbol{\Gamma}_0)$ . We ask to perform a simulation using this setting in Exercise 3.

**Example: Simulation exercise of the longitudinal normal model with heteroskedasticity**

Let's perform a simulation exercise to assess some potential extensions of the longitudinal hierarchical normal model. The point of departure is to assume that

$$y_{it} = \beta_0 + \beta_1 x_{it1} + \beta_2 x_{it2} + \beta_3 x_{it3} + b_i + w_{it1}b_{i1} + \mu_{it},$$

where  $x_{itk} \sim N(0, 1)$ ,  $k = 1, 2, 3$ ,  $w_{it1} \sim N(0, 1)$ ,  $b_i \sim N(0, 0.7^{1/2})$ ,  $b_{i1} \sim N(0, 0.6^{1/2})$ ,  $\mu_{it} \sim N(0, (0.1/\tau_{it})^{1/2})$ ,  $\tau_{it} \sim G(v/2, v/2)$  and  $\boldsymbol{\beta} = [0.5 \ 0.4 \ 0.6 \ -0.6]^\top$ ,  $i = 1, 2, \dots, 50$ . The sample size is 2,000 in an *unbalanced panel structure*.

Following same stages as in this section and Exercise 1, the posterior conditional distributions assuming that  $\mu_{it} | \tau_{it} \sim N(0, \sigma^2/\tau_{it})$ ,  $\tau_{it} \sim G(v/2, v/2)$  are given by

$$\boldsymbol{\beta} | \sigma^2, \boldsymbol{\tau}, \mathbf{D}, \mathbf{y}, \mathbf{X}, \mathbf{W} \sim N(\boldsymbol{\beta}_n, \mathbf{B}_n),$$

where  $\boldsymbol{\tau} = [\tau_{it}]^\top$ ,  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \sum_{i=1}^N \mathbf{X}_i^\top \mathbf{V}_i^{-1} \mathbf{X}_i)^{-1}$ ,  $\boldsymbol{\beta}_n = \mathbf{B}_n(\mathbf{B}_0^{-1}\boldsymbol{\beta}_0 + \sum_{i=1}^N \mathbf{X}_i^\top \mathbf{V}_i^{-1} \mathbf{y}_i)$ ,  $\mathbf{V}_i = \sigma^2 \boldsymbol{\Psi}_i + \sigma_b^2 \mathbf{i}_{T_i} \mathbf{i}_{T_i}^\top$  and  $\boldsymbol{\Psi}_i = \text{diag}\{\tau_{it}^{-1}\}$ .

$$\mathbf{b}_i | \boldsymbol{\beta}, \sigma^2, \boldsymbol{\tau}, \mathbf{D}, \mathbf{y}, \mathbf{X}, \mathbf{W} \sim N(\mathbf{b}_{ni}, \mathbf{B}_{ni}),$$

where  $\mathbf{B}_{ni} = (\sigma^{-2} \mathbf{W}_i^\top \boldsymbol{\Psi}_i^{-1} \mathbf{W}_i + \mathbf{D}^{-1})^{-1}$  and  $\mathbf{b}_{ni} = \mathbf{B}_{ni}(\sigma^{-2} \mathbf{W}_i^\top \boldsymbol{\Psi}_i^{-1} (\mathbf{y}_i -$

$\mathbf{X}_i\boldsymbol{\beta})$ .

$$\sigma^2 \mid \boldsymbol{\beta}, \mathbf{b}, \boldsymbol{\tau}, \mathbf{y}, \mathbf{X}, \mathbf{W} \sim IG(\alpha_n, \delta_n),$$

where  $\alpha_n = \alpha_0 + \frac{1}{2} \sum_{i=1}^N T_i$  and  $\delta_n = \delta_0 + \frac{1}{2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{W}_i\mathbf{b}_i)^\top \boldsymbol{\Psi}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{W}_i\mathbf{b}_i)$ .

$$\mathbf{D} \mid \mathbf{b} \sim IW(d_n, \mathbf{D}_n),$$

where  $d_n = d_0 + N$  and  $\mathbf{D}_n = d_0 \mathbf{D}_0 + \sum_{i=1}^N \mathbf{b}_i \mathbf{b}_i^\top$ . And

$$\tau_{it} \mid \sigma^2, \boldsymbol{\beta}, \mathbf{b}, \mathbf{y}, \mathbf{X}, \mathbf{W} \sim G(v_{1n}/2, v_{2ni}/2),$$

where  $v_{1n} = v + 1$  and  $v_{2ni} = v + \sigma^{-2}(y_{it} - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i)^2$ .

The following code implements this simulation, and gets draws of the posterior distributions. We set MCMC iterations, burn-in and thinning parameters equal to 5,000, 1,000 and 1, respectively. In addition,  $\boldsymbol{\beta}_0 = \mathbf{0}_5$ ,  $\mathbf{B}_0 = \mathbf{I}_5$ ,  $\alpha_0 = \delta_0 = 0.001$ ,  $d_0 = 2$ ,  $\mathbf{D}_0 = \mathbf{I}_2$  and  $v = 5$ .

### R code. Simulation exercise: Longitudinal normal model with heteroskedasticity from scratch

```

1 rm(list = ls()); set.seed(10101)
2 NT <- 2000; N <- 50
3 id <- c(1:N, sample(1:N, NT - N, replace=TRUE))
4 table(id)
5 x1 <- rnorm(NT); x2 <- rnorm(NT); x3 <- rnorm(NT)
6 X <- cbind(1, x1, x2, x3); K1 <- dim(X)[2]
7 w1 <- rnorm(NT); W <- cbind(1, w1)
8 K2 <- dim(W)[2]; B <- c(0.5, 0.4, 0.6, -0.6)
9 D <- c(0.7, 0.6)
10 b1 <- rnorm(N, 0, sd = D[1]^0.5)
11 b2 <- rnorm(N, 0, sd = D[2]^0.5)
12 b <- cbind(b1, b2)
13 v <- 5; tau <- rgamma(NT, shape = v/2, rate = v/2)
14 sig2 <- 0.1; u <- rnorm(NT, 0, sd = (sig2/tau)^0.5)
15 y <- NULL
16 for(i in 1:NT){
17   yi <- X[i,] %*% B + W[i,] %*% b[id[i],] + u[i]
18   y <- c(y, yi)
19 }
20 Data <- as.data.frame(cbind(y, x1, x2, x3, w1, id))
21 mcmc <- 5000; burnin <- 1000; thin <- 1; tot <- mcmc +
  burnin
22 b0 <- rep(0, K1); B0 <- diag(K1); B0i <- solve(B0)
23 r0 <- K2; R0 <- diag(K2); a0 <- 0.001; d0 <- 0.001

```

**R code. Simulation exercise: Longitudinal normal model with heteroskedasticity from scratch**

```

1 PostBeta <- function(sig2, D, tau){
2   XVX <- matrix(0, K1, K1)
3   XVy <- matrix(0, K1, 1)
4   for(i in 1:N){
5     ids <- which(id == i)
6     Ti <- length(ids)
7     Wi <- W[ids, ]
8     taui <- tau[ids]
9     Vi <- sig2*solve(diag(1/taui)) + Wi%*%D%*%t(Wi)
10    ViInv <- solve(Vi)
11    Xi <- X[ids, ]
12    XVXi <- t(Xi)%*%ViInv%*%Xi
13    XVX <- XVX + XVXi
14    yi <- y[ids]
15    XVyi <- t(Xi)%*%ViInv%*%yi
16    XVy <- XVy + XVyi
17  }
18  Bn <- solve(B0i + XVX)
19  bn <- Bn%*%(B0i%*%b0 + XVy)
20  Beta <- MASS:::mvrnorm(1, bn, Bn)
21  return(Beta)
22 }
23 Postb <- function(Beta, sig2, D, tau){
24   Di <- solve(D); bis <- matrix(0, N, K2)
25   for(i in 1:N){
26     ids <- which(id == i)
27     Wi <- W[ids, ]; Xi <- X[ids, ]
28     yi <- y[ids]; taui <- tau[ids]
29     Taui <- solve(diag(1/taui))
30     Wtei <- sig2^(-1)*t(Wi)%*%Taui%*%(yi - Xi%*%Beta)
31     Bni <- solve(sig2^(-1)*t(Wi)%*%Taui%*%Wi + Di)
32     bni <- Bni%*%Wtei
33     bi <- MASS:::mvrnorm(1, bni, Bni)
34     bis[i, ] <- bi
35   }
36   return(bis)
37 }
```

**R code. Simulation exercise: Longitudinal normal model with heteroskedasticity from scratch**

```

1 PostSig2 <- function(Beta, bs, tau){
2   an <- a0 + 0.5*NT
3   ete <- 0
4   for(i in 1:N){
5     ids <- which(id == i)
6     Xi <- X[ids, ]; yi <- y[ids]
7     Wi <- W[ids, ]; taui <- tau[ids]
8     Taui <- solve(diag(1/taui))
9     ei <- yi - Xi%*%Beta - Wi%*%bs[i, ]
10    etei <- t(ei)%*%Taui%*%ei
11    ete <- ete + etei
12  }
13  dn <- d0 + 0.5*ete
14  sig2 <- MCMCpack::rinvgamma(1, shape = an, scale = dn)
15  return(sig2)
16 }
17 PostD <- function(bs){
18   rn <- r0 + N
19   btb <- matrix(0, K2, K2)
20   for(i in 1:N){
21     bsi <- bs[i, ]
22     btbi <- bsi%*%t(bsi)
23     btb <- btb + btbi
24   }
25   Rn <- d0*R0 + btb
26   Sigma <- MCMCpack::riwish(v = rn, S = Rn)
27   return(Sigma)
28 }
29 PostTau <- function(sig2, Beta, bs){
30   v1n <- v + 1
31   v2n <- NULL
32   for(i in 1:NT){
33     Xi <- X[i, ]; yi <- y[i]
34     Wi <- W[i, ]; bi <- bs[id[i], ]
35     v2ni <- v + sig2^(-1)*(yi - Xi%*%Beta - Wi%*%bi)^2
36     v2n <- c(v2n, v2ni)
37   }
38   tau <- rgamma(NT, shape = rep(v1n/2, NT), rate = v2n/2)
39   return(tau)
40 }
41 PostBetas <- matrix(0, tot, K1); PostDs <- matrix(0, tot, K2
  *(K2+1)/2)
42 PostSig2s <- rep(0, tot); Postbs <- array(0, c(N, K2, tot))
43 PostTaus <- matrix(0, tot, NT); RegLS <- lm(y ~ X - 1)
44 SumLS <- summary(RegLS)
45 Beta <- SumLS[["coefficients"]][,1]
46 sig2 <- SumLS[["sigma"]]^2; D <- diag(K2)
47 tau <- rgamma(NT, shape = v/2, rate = v/2)

```

**R code. Simulation exercise: Longitudinal normal model with heteroskedasticity from scratch**

```

1 pb <- txtProgressBar(min = 0, max = tot, style = 3)
2 for(s in 1:tot){
3   bs <- Postb(Beta = Beta, sig2 = sig2, D = D, tau = tau)
4   D <- PostD(bs = bs)
5   Beta <- PostBeta(sig2 = sig2, D = D, tau = tau)
6   sig2 <- PostSig2(Beta = Beta, bs = bs, tau = tau)
7   tau <- PostTau(sig2 = sig2, Beta = Beta, bs = bs)
8   PostBetas[s,] <- Beta
9   PostDs[s,] <- matrixcalc::vech(D)
10  PostSig2s[s] <- sig2
11  Postbs[, , s] <- bs
12  PostTaus[s,] <- tau
13  setTxtProgressBar(pb, s)
14 }
15 close(pb)
16 keep <- seq((burnin+1), tot, thin)
17 Bs <- PostBetas[keep,]; Ds <- PostDs[keep,]
18 bs <- Postbs[, , keep]; sig2s <- PostSig2s[keep]
19 taus <- PostTaus[keep,]
20 summary(coda::mcmc(Bs)); summary(coda::mcmc(Ds))
21 summary(coda::mcmc(sig2s))

```

We can see that all the 95% credible intervals encompass the population parameters, except for the second *fixed effect* and the variance of the model, but both for a tiny margin.

## 9.2 Logit model

We can use the framework of Section 9.1 to perform inference in models with longitudinal/panel data of binary response variables. In particular, let  $y_{it} \sim B(\pi_{it})$ , where  $\text{logit}(\pi_{it}) = \log\left(\frac{\pi_{it}}{1-\pi_{it}}\right) \equiv y_{it}^*$ , such that  $y_{it}^* \sim N(\mathbf{x}_{it}^\top \boldsymbol{\beta} + \mathbf{w}_{it}^\top \mathbf{b}_i, \sigma^2)$ . Thus, we can *augment* the model with the latent variable  $y_{it}^*$  and perform inference using a Metropolis-within-Gibbs sampling algorithm based on the posterior conditional distributions from the previous section.

We can implement a Gibbs sampling algorithm to sample draws from the posterior conditional distributions of  $\boldsymbol{\beta}$ ,  $\sigma^2$ ,  $\mathbf{b}_i$ , and  $\mathbf{D}$  using the equations in Section 9.1 conditional on  $\mathbf{y}_i^*$ . Then, we can use a random walk Metropolis-

Hastings algorithm to sample  $y_{it}^*$ , where the proposal distribution is Gaussian with mean  $y_{it}^*$  and variance  $v^2$ , that is,  $y_{it}^{*c} = y_{it}^* + \epsilon_{it}$ , where  $\epsilon_{it} \sim \mathcal{N}(0, v^2)$ , and  $v$  is a tuning parameter to achieve good acceptance rates.

Finally, for making predictions, we should take into account that

$$\mathbb{E}[\pi_{it}] = \frac{1}{1 + \exp \left\{ (\mathbf{x}_{it}^\top \boldsymbol{\beta} + \mathbf{w}_{it}^\top \mathbf{b}_i) / \sqrt{1 + \left( \frac{16\sqrt{3}}{15\pi} \right)^2 \sigma^2} \right\}}.$$

See [121, pag. 136].

The posterior distribution of this model is

$$\begin{aligned} \pi(\boldsymbol{\beta}, \sigma^2, \mathbf{b}_i, \mathbf{D}, \mathbf{y}^* \mid \mathbf{y}, \mathbf{X}, \mathbf{W}) &\propto \prod_{i=1}^N \prod_{t=1}^{T_i} \left\{ \pi_{it}^{y_{it}} (1 - \pi_{it})^{1-y_{it}} (\sigma^2)^{-1} \right. \\ &\quad \times \exp \left\{ -(y_{it}^* - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i)^\top (y_{it}^* - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i) \right. \\ &\quad \left. \left. \frac{1}{2\sigma^2} \right\} \right\} \\ &\quad \times \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \right\} \\ &\quad \times \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \mathbf{b}_i^\top \mathbf{D}^{-1} \mathbf{b}_i \right\} \\ &\quad \times (\sigma^2)^{-\alpha_0 - 1} \exp \left\{ -\frac{\delta_0}{\sigma^2} \right\} \\ &\quad \times |\mathbf{D}|^{-(d_0 + K_2 + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(d_0 \mathbf{D}_0 \mathbf{D}^{-1}) \right\}. \end{aligned}$$

We can get samples of  $y_{it}^*$  from a normal distribution with mean equal to  $\mathbf{x}_{it}^\top \boldsymbol{\beta} + \mathbf{w}_{it}^\top \mathbf{b}_i$  and variance  $\sigma^2$ , and use these samples to get  $\pi_{it} = \frac{1}{1+e^{-y_{it}^*}}$ ,  $y_{it}^{*c} = y_{it}^* + \epsilon_{it}$  and  $\pi_{it}^c = \frac{1}{1+e^{-y_{it}^{*c}}}$ , and calculate the acceptance rate of the Metropolis-Hastings algorithm,

$$\begin{aligned} \alpha &= \min \left( 1, \frac{\pi_{it}^{cy_{it}} (1 - \pi_{it}^c)^{(1-y_{it})}}{\pi_{it}^{y_{it}} (1 - \pi_{it})^{(1-y_{it})}} \right. \\ &\quad \times \left. \frac{\exp \left\{ -\frac{1}{2\sigma^2} (y_{it}^{*c} - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i)^\top (y_{it}^{*c} - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i) \right\}}{\exp \left\{ -\frac{1}{2\sigma^2} (y_{it}^* - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i)^\top (y_{it}^* - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i) \right\}} \right). \end{aligned}$$

### Example: Doctor visits in Germany

We used the dataset *9VisitDoc.csv* provided by [452].<sup>2</sup> We analyze the determinants of a binary variable (DocVis), which equals 1 if an individual

<sup>2</sup>See <http://qed.econ.queensu.ca/jae/2004-v19.4/winkelmann/> for details

visited a physician in the last three months and 0 otherwise. The dataset contains 32,837 observations of 9,197 individuals in an *unbalanced longitudinal/panel* dataset over the years 1995–1999 from the German Socioeconomic Panel Data.

The specification is given by

$$\begin{aligned}\text{logit}(\pi_{it}) = & \beta_1 + \beta_2 \text{Age} + \beta_3 \text{Male} + \beta_4 \text{Sport} + \beta_5 \text{LogInc} \\ & + \beta_6 \text{GoodHealth} + \beta_7 \text{BadHealth} + b_i + b_{i1} \text{Sozh},\end{aligned}$$

where  $\pi_{it} = p(\text{DocVis}_{it} = 1)$ .

This specification controls for age, a gender indicator (with 1 representing male), whether the individual practices any sport (with 1 for sport), the logarithm of monthly gross income, and self-perception of health status, where “good” and “bad” are compared to a baseline of “regular”. Additionally, we assume that unobserved heterogeneity is linked to whether the individual receives welfare payments (with Sozh equal to 1 for receiving welfare).

We set 10,000 MCMC iterations, plus 1,000 burn-in, and a thinning parameter equal to 10. In addition,  $\boldsymbol{\beta}_0 = \mathbf{0}_7$ ,  $\boldsymbol{B}_0 = \mathbf{I}_7$ ,  $\alpha_0 = \delta_0 = 0.001$ ,  $d_0 = 5$ , and  $\boldsymbol{D}_0 = \mathbf{I}_2$ .

The Algorithm A26 shows how to perform inference of the hierarchical longitudinal logit model using our GUI. We show in the following code how to perform inference of this example using the command *MCMChlogit* from the *MCMCpack* package. We fixed the variance for over-dispersion ( $\sigma^2$ ) setting *FixOD* = 1 in this example. Our GUI does not fix this value, that is, it sets *FixOD* = 0, which is the default value in the command *MCMChlogit*. We ask to replicate this example using our GUI in Exercise 4. The command *MCMChlogit* uses an adaptive algorithm to tune *v* based on an optimal acceptance rate equal to 0.44.

### R code. Doctor visits in Germany

```

1 rm(list = ls()); set.seed(12345)
2 Data <- read.csv("https://raw.githubusercontent.com/
  besmarter/BSTApp/refs/heads/master/DataApp/9VisitDoc.csv"
  , sep = ",", header = TRUE, quote = "")
3 attach(Data)
4 K1 <- 7; K2 <- 2; N <- 9197
5 b0 <- rep(0, K1); B0 <- diag(K1)
6 r0 <- 5; R0 <- diag(K2)
7 a0 <- 0.001; d0 <- 0.001
8 RegLogit <- glm(DocVis ~ Age + Male + Sport + LogInc +
  GoodHealth + BadHealth, family = binomial(link = "logit"
  ))
```

### R code. Doctor visits in Germany

```

1 SumLogit <- summary(RegLogit)
2 Beta0 <- SumLogit[["coefficients"]][,1]
3 mcmc <- 10000; burnin <- 1000; thin <- 10
4 # MCMChlogit
5 Resultshlogit <- MCMCpack::MCMChlogit(fixed = DocVis ~ Age +
   Male + Sport + LogInc + GoodHealth + BadHealth, random
   = ~ Sozh, group = "id", data = Data, burnin = burnin, mcmc
   = mcmc, thin = thin, mubeta = b0, Vbeta = B0, r = r0, R
   = R0, nu = a0, delta = d0, beta.start = Beta0, FixOD =
   1)
6 Betas <- Resultshlogit[["mcmc"]][,1:K1]
7 Sigma2RanEff <- Resultshlogit[["mcmc"]][,c(K2*N+K1+1, 2*N+K1
   +K2^2)]
8 summary(Betas); summary(Sigma2RanEff)

```

The results suggest that age, sports, income and a bad perception of health status increase the probability of visiting the physician, the posterior estimates have 95% symmetric credible intervals equal to (5.1e-03, 1.3e-02), (0.23, 0.40), (0.18, 0.37) and (1.22, 1.53), whereas men have a lower probability of visiting a physician, the 95% credible interval is (-1.19, -1.01), and individuals who have a good perception of their health status also have a lower probability of visiting the doctor, the 95% credible interval is (-1.16, -0.98). The 95% credible interval of the variances of the unobserved heterogeneity associated with the welfare program is (0.35, 1.27).

### 9.3 Poisson model

We can use same ideas as in Section 9.2 to perform inference in longitudinal/panel datasets where the dependent variable takes non-negative integers. Let's assume that  $y_{it} \sim P(\lambda_{it})$  where  $\log(\lambda_{it}) = y_{it}^*$  such that  $y_{it}^* \sim N(\mathbf{x}_{it}^\top \boldsymbol{\beta} + \mathbf{w}_{it}^\top \mathbf{b}_i, \sigma^2)$ . We can augment the model with the latent variable  $y_{it}^*$ , and again use a Metropolis-within-Gibbs algorithm to perform inference in this model.

The posterior distribution of this model is

$$\begin{aligned} \pi(\boldsymbol{\beta}, \sigma^2, \mathbf{b}_i, \mathbf{D}, \mathbf{y}^* | \mathbf{y}, \mathbf{X}, \mathbf{W}) \propto & \prod_{i=1}^N \prod_{t=1}^{T_i} \{ \lambda_{it}^{y_{it}^*} \exp\{-\lambda_{it}\} (\sigma^2)^{-1} \\ & \times \exp\{-(y_{it}^* - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i)^\top (y_{it}^* - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i)\} \end{aligned}$$

**Algorithm A26** Hierarchical longitudinal logit models

- 
- 1: Select *Hierarchical Longitudinal Model* on the top panel
  - 2: Select *Logit* model using the left radio button
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
  - 5: Write down the formula of the *fixed effects* equation in the **Main Equation: Fixed Effects** box. This formula must be written using the syntax of the *formula* command of **R** software. This equation includes intercept by default, do not include it in the equation
  - 6: Write down the formula of the *random effects* equation in the **Main Equation: Random Effects** box without writing the dependent variable, that is, starting the equation with the *tilde* ("~") symbol. This formula must be written using the syntax of the *formula* command of **R** software. This equation includes intercept by default, do not include it in the equation. If there are just random intercepts do not write anything in this box
  - 7: Write down the name of the grouping variable, that is, the variable that indicates the cross-sectional units
  - 8: Set the hyperparameters of the *fixed effects*: mean vector, covariance matrix, shape and scale parameters. This step is not necessary as by default our GUI uses non-informative priors
  - 9: Set the hyperparameters of the *random effects*: degrees of freedom and scale matrix of the inverse Wishart distribution. This step is not necessary as by default our GUI uses non-informative priors
  - 10: Click the *Go!* button
  - 11: Analyze results
  - 12: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons
- 

$$\begin{aligned}
 & \frac{1}{2\sigma^2} \left\{ \right\} \\
 & \times \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \right\} \\
 & \times \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \mathbf{b}_i^\top \mathbf{D}^{-1} \mathbf{b}_i \right\} \\
 & \times (\sigma^2)^{-\alpha_0 - 1} \exp \left\{ -\frac{\delta_0}{\sigma^2} \right\} \\
 & \times |\mathbf{D}|^{-(d_0 + K_2 + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(d_0 \mathbf{D}_0 \mathbf{D}^{-1}) \right\}.
 \end{aligned}$$

We can get samples of  $y_{it}^*$  from a normal distribution with mean equal to  $\mathbf{x}_{it}^\top \boldsymbol{\beta} + \mathbf{w}_{it}^\top \mathbf{b}_i$  and variance  $\sigma^2$ , and use these samples to get  $\lambda_{it} = \exp(y_{it}^*)$ ,  $y_{it}^{*c} = y_{it}^* + \epsilon_{it}$ , where  $\epsilon_{it} \sim \mathcal{N}(0, v^2)$ ,  $v$  is a tuning parameter to get good acceptance rates, and  $\lambda_{it}^c = \exp(y_{it}^{*c})$ . The acceptance rate of the Metropolis-Hastings algorithm is

$$\alpha = \min \left( 1, \frac{\lambda_{it}^{cy_{it}} \exp(-\lambda_{it}^c)}{\lambda_{it}^{y_{it}} \exp(-\lambda_{it})} \times \frac{\exp \left\{ -\frac{1}{2\sigma^2} (y_{it}^{*c} - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i)^\top (y_{it}^{*c} - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i) \right\}}{\exp \left\{ -\frac{1}{2\sigma^2} (y_{it}^* - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i)^\top (y_{it}^* - \mathbf{x}_{it}^\top \boldsymbol{\beta} - \mathbf{w}_{it}^\top \mathbf{b}_i) \right\}} \right).$$

In addition, we should use the posterior conditional distributions from Section 9.1 to complete the algorithm getting samples of  $\boldsymbol{\beta}$ ,  $\sigma^2$ ,  $\mathbf{b}_i$  and  $\mathbf{D}$  replacing  $y_{it}$  by  $y_{it}^*$ .

We should take into account for doing predictions that

$$\mathbb{E}[\lambda_{it}] = \exp \left\{ \mathbf{x}_{it}^\top \boldsymbol{\beta} + \mathbf{w}_{it}^\top \mathbf{b}_i + 0.5\sigma^2 \right\}.$$

See [121, pag. 137].

#### Example: Simulation exercise

Let's perform a simulation exercise to assess the performance of the hierarchical longitudinal Poisson model. The point of departure is to assume that

$$y_{it}^* = \beta_0 + \beta_1 x_{it1} + \beta_2 x_{it2} + \beta_3 x_{it3} + b_i + w_{it1} b_{i1},$$

where  $x_{itk} \sim N(0, 1)$  for  $k = 1, 2, 3$ ,  $w_{it1} \sim N(0, 1)$ ,  $b_i \sim N(0, 0.7^{1/2})$ ,  $b_{i1} \sim N(0, 0.6^{1/2})$ , and  $\boldsymbol{\beta} = [0.5 \ 0.4 \ 0.6 \ -0.6]^\top$ , with  $i = 1, 2, \dots, 50$ . Additionally,  $y_{it} \sim P(\lambda_{it})$ , where  $\lambda_{it} = \exp(y_{it}^*)$ . The sample size is 1,000 in an *unbalanced panel structure*.

We set the priors as  $\boldsymbol{\beta}_0 = \mathbf{0}_4$ ,  $\mathbf{B}_0 = \mathbf{I}_4$ ,  $\alpha_0 = \delta_0 = 0.001$ ,  $d_0 = 2$ , and  $\mathbf{D}_0 = \mathbf{I}_2$ . The number of MCMC iterations, burn-in, and thinning parameters are 15,000, 5,000, and 10, respectively.

The following code shows how to perform inference in the hierarchical longitudinal Poisson model programming the Metropolis-within-Gibbs sampler.

**R code. Simulation exercise: Hierarchical longitudinal Poisson model**

---

```
1 rm(list = ls()); set.seed(10101)
2 NT <- 1000; N <- 50
3 id <- c(1:N, sample(1:N, NT - N, replace=TRUE))
```

---

**R code. Simulation exercise: Hierarchical longitudinal Poisson model**

```

1 x1 <- rnorm(NT); x2 <- rnorm(NT); x3 <- rnorm(NT)
2 X <- cbind(1, x1, x2, x3)
3 K1 <- dim(X)[2]; w1 <- rnorm(NT)
4 W <- cbind(1, w1); K2 <- dim(W)[2]
5 B <- c(0.5, 0.4, 0.6, -0.6)
6 D <- c(0.7, 0.6); sig2 <- 0.1
7 b1 <- rnorm(N, 0, sd = D[1]^0.5)
8 b2 <- rnorm(N, 0, sd = D[2]^0.5)
9 b <- cbind(b1, b2)
10 yl <- NULL
11 for(i in 1:NT){
12   ylmeani <- X[i,] %*% B + W[i,] %*% b[id[i],]
13   yli <- rnorm(1, ylmeani, sig2^0.5)
14   yl <- c(yl, yli)
15 }
16 lambdait <- exp(yl); y <- rpois(NT, lambdait)
17 Data <- as.data.frame(cbind(y, x1, x2, x3, w1, id))
18 mcmc <- 15000; burnin <- 5000; thin <- 10; tot <- mcmc +
  burnin
19 b0 <- rep(0, K1); B0 <- diag(K1); B0i <- solve(B0)
20 r0 <- K2; R0 <- diag(K2); a0 <- 0.001; d0 <- 0.001
21 LatentMHV1 <- function(tuning, Beta, bs, sig2){
22   ylhat <- rep(0, NT)
23   accept <- NULL
24   for(i in 1:NT){
25     ids <- which(id == i)
26     yi <- y[i]
27     ylhatmeani <- X[i,] %*% Beta + W[i,] %*% bs[id[i],]
28     ylhati <- rnorm(1, ylhatmeani, sd = sig2^0.5)
29     lambdahati <- exp(ylhati)
30     ei <- rnorm(1, 0, sd = tuning)
31     ylpropi <- ylhati + ei
32     lambdapropi <- exp(ylpropi)
33     logPosthati <- sum(dpois(yi, lambdahati, log = TRUE) +
      dnorm(ylhati, ylhatmeani, sig2^0.5, log = TRUE))
34     logPostpropri <- sum(dpois(yi, lambdapropi, log = TRUE) +
      dnorm(ylpropi, ylhatmeani, sig2^0.5, log = TRUE))
35     alphai <- min(1, exp(logPostpropri - logPosthati))
36     ui <- runif(1)
37     if(ui <= alphai){
38       ylhati <- ylpropi; accepti <- 1
39     }else{
40       ylhati <- ylhati; accepti <- 0
41     }
42     ylhat[i] <- ylhati
43     accept <- c(accept, accepti)
44   }
45   res <- list(ylhat = ylhat, accept = mean(accept))
46   return(res)
47 }
```

**R code. Simulation exercise: Hierarchical longitudinal Poisson model**

```

1 PostBeta <- function(D, ylhat, sig2){
2   XVX <- matrix(0, K1, K1); XVy <- matrix(0, K1, 1)
3   for(i in 1:N){
4     ids <- which(id == i); Ti <- length(ids)
5     Wi <- W[ids, ]
6     Vi <- diag(Ti)*sig2 + Wi%*%D%*%t(Wi)
7     ViInv <- solve(Vi); Xi <- X[ids, ]
8     XVXi <- t(Xi)%*%ViInv%*%Xi
9     XVX <- XVX + XVXi
10    yi <- ylhat[ids]
11    XVy <- t(Xi)%*%ViInv%*%yi
12    XVy <- XVy + XVy
13  }
14  Bn <- solve(B0i + XVX); bn <- Bn%*%(B0i%*%b0 + XVy)
15  Beta <- MASS::mvrnorm(1, bn, Bn)
16  return(Beta)
17 }
18 Postb <- function(Beta, D, ylhat, sig2){
19   Di <- solve(D); bis <- matrix(0, N, K2)
20   for(i in 1:N){
21     ids <- which(id == i)
22     Wi <- W[ids, ]; Xi <- X[ids, ]
23     yi <- ylhat[ids]
24     Wtei <- sig2^(-1)*t(Wi)%*%(yi - Xi%*%Beta)
25     Bni <- solve(sig2^(-1)*t(Wi)%*%Wi + Di)
26     bni <- Bni%*%Wtei
27     bi <- MASS::mvrnorm(1, bni, Bni)
28     bis[i, ] <- bi
29   }
30   return(bis)
31 }
32 PostD <- function(bs){
33   rn <- r0 + N; btb <- matrix(0, K2, K2)
34   for(i in 1:N){
35     bsi <- bs[i, ]; btbi <- bsi%*%t(bsi)
36     btb <- btb + btbi
37   }
38   Rn <- d0*R0 + btb
39   Sigma <- MCMCpack::riwish(v = rn, S = Rn)
40   return(Sigma)
41 }
```

*R code. Simulation exercise: Hierarchical longitudinal Poisson model*

```

1 PostSig2 <- function(Beta, bs, ylhat){
2   an <- a0 + 0.5*NT; ete <- 0
3   for(i in 1:N){
4     ids <- which(id == i)
5     Xi <- X[ids, ]
6     yi <- ylhat[ids]
7     Wi <- W[ids, ]
8     ei <- yi - Xi%*%Beta - Wi%*%bs[i, ]
9     etei <- t(ei)%*%ei
10    ete <- ete + etei
11  }
12  dn <- d0 + 0.5*ete
13  sig2 <- MCMCpack::rinvgamma(1, shape = an, scale = dn)
14  return(sig2)
15 }
16 PostBetas <- matrix(0, tot, K1); PostDs <- matrix(0, tot, K2
  *(K2+1)/2)
17 Postbs <- array(0, c(N, K2, tot)); PostSig2s <- rep(0, tot)
18 Accepts <- rep(NULL, tot)
19 RegPois <- glm(y ~ X - 1, family = poisson(link = "log"))
20 SumPois <- summary(RegPois)
21 Beta <- SumPois[["coefficients"]][,1]
22 sig2 <- sum(SumPois[["deviance.resid"]])^2/SumPois[["df.
  residual"]]
23 D <- diag(K2); bs1 <- rnorm(N, 0, sd = D[1,1]^0.5)
24 bs2 <- rnorm(N, 0, sd = D[2,2]^0.5); bs <- cbind(bs1, bs2)
25 tuning <- 0.1; ropt <- 0.44
26 tunepariter <- seq(round(tot/10, 0), tot, round(tot/10, 0));
  1 <- 1
27 pb <- txtProgressBar(min = 0, max = tot, style = 3)
28 for(s in 1:tot){
29   LatY <- LatentMHV1(tuning = tuning, Beta = Beta, bs = bs,
    sig2 = sig2)
30   ylhat <- LatY[["ylhat"]]
31   bs <- Postb(Beta = Beta, D = D, ylhat=ylhat, sig2 = sig2)
32   D <- PostD(bs = bs)
33   Beta <- PostBeta(D = D, ylhat = ylhat, sig2 = sig2)
34   sig2 <- PostSig2(Beta = Beta, bs = bs, ylhat = ylhat)
35   PostBetas[s,] <- Beta
36   PostDs[s,] <- matrixcalc::vech(D)
37   Postbs[, , s] <- bs; PostSig2s[s] <- sig2
38   AcceptRate <- LatY[["accept"]]
39   Accepts[s] <- AcceptRate

```

**R code. Simulation exercise: Hierarchical longitudinal Poisson model**

```

1  if(AcceptRate > ropt){
2    tuning = tuning*(2-(1-AcceptRate)/(1-ropt))
3  }else{
4    tuning = tuning/(2-AcceptRate/ropt)
5  }
6  if(s == tunepariter[1]){
7    print(AcceptRate); l <- l + 1
8  }
9  setTxtProgressBar(pb, s)
10 }
11 close(pb)
12 keep <- seq((burnin+1), tot, thin)
13 Bs <- PostBetas[keep,]; Ds <- PostDs[keep,]
14 bs <- Postbs[, , keep]; sig2s <- PostSig2s[keep]
15 summary(coda::mcmc(Bs)); summary(coda::mcmc(Ds))

```

We can see that all 95% credible intervals encompass the population parameters of the *fixed effects*, the posterior medians are relatively near the population values. However, we do not get good posterior estimates of the covariance matrix of the *random effects* as the 95% credible intervals do not encompass the second element of the diagonal of this matrix. In addition, the posterior draws of this algorithm over-estimates the over-dispersion parameter.

We can perform inference for the hierarchical longitudinal Poisson model in our GUI using Algorithm A27. Our GUI is based on the *MCMChpoisson* command from the *MCMCpack* package.

## 9.4 Summary

In this chapter, we present how to perform inference in longitudinal/panel data models from a Bayesian perspective. In particular, the Bayesian approach uses a hierarchical structure, where the *random effects* have priors that depend on hyperparameters, which in turn also have priors. We cover the three most common cases: continuous, binary, and count dependent variables. The basic models presented in this chapter can be easily extended to more flexible cases, given the hierarchical structure.

**Algorithm A27** Hierarchical longitudinal Poisson models

- 1: Select *Hierarchical Longitudinal Model* on the top panel
- 2: Select *Poisson* model using the left radio button
- 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
- 4: Select MCMC iterations, burn-in and thinning parameters using the *Range sliders*
- 5: Write down the formula of the *fixed effects* equation in the **Main Equation: Fixed Effects** box. This formula must be written using the syntax of the *formula* command of **R** software. This equation includes intercept by default, do not include it in the equation
- 6: Write down the formula of the *random effects* equation in the **Main Equation: Random Effects** box without writing the dependent variable, that is, starting the equation with the *tilde* (~) symbol. This formula must be written using the syntax of the *formula* command of **R** software. This equation includes intercept by default, do not include it in the equation. If there are just random intercepts do not write anything in this box
- 7: Write down the name of the grouping variable, that is, the variable that indicates the cross-sectional units
- 8: Set the hyperparameters of the *fixed effects*: mean vector, covariance matrix, shape and scale parameters. This step is not necessary as by default our GUI uses non-informative priors
- 9: Set the hyperparameters of the *random effects*: degrees of freedom and scale matrix of the inverse Wishart distribution. This step is not necessary as by default our GUI uses non-informative priors
- 10: Click the *Go!* button
- 11: Analyze results
- 12: Download posterior chains and diagnostic plots using the *Download Posterior Chains* and *Download Posterior Graphs* buttons

\*At the time of writing this book there was an issue with the function *MCMChpoisson* from *MCMCpack*. We contact the maintainer, but users may have issues running this algorithm or running this function directly in **R**.

## 9.5 Exercises

1. Show that the posterior distribution of  $\beta \mid \sigma^2, \mathbf{D}$  is  $N(\beta_n, \mathbf{B}_n)$ , where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \sum_{i=1}^N \mathbf{X}_i^\top \mathbf{V}_i^{-1} \mathbf{X}_i)^{-1}$ ,  $\beta_n = \mathbf{B}_n(\mathbf{B}_0^{-1}\beta_0 + \sum_{i=1}^N \mathbf{X}_i^\top \mathbf{V}_i^{-1} \mathbf{y}_i)$ .
2. The relation between productivity and public investment example continues

- Perform inference of this example using our GUI.
- Program from scratch a Gibbs sampling algorithm to perform this application. Set  $\beta_0 = \mathbf{0}_5$ ,  $B_0 = \mathbf{I}_5$ ,  $\alpha_0 = \delta_0 = 0.001$ ,  $d_0 = 5$  and  $D_0 = \mathbf{I}_1$ .
- Perform inference in this example assuming that  $\mu_{it} \mid \tau_{it} \sim N(0, \sigma^2 / \tau_{it})$  and  $\tau_{it} \sim G(v/2, v/2)$  setting  $v = 5$ .

### 3. Simulation exercise of the longitudinal normal model continues

Assume that

$$y_{it} = \beta_0 + \beta_1 x_{it1} + \beta_2 x_{it2} + \beta_3 x_{it3} + \beta_4 z_{i1} + b_i + w_{it1} b_{i1} + \mu_{it},$$

where  $x_{itk} \sim N(0, 1)$ ,  $k = 1, 2, 3$ ,  $z_{i1} \sim B(0.5)$ ,  $w_{it1} \sim N(0, 1)$ ,  $b_i \sim N(0, 0.7^{1/2})$ ,  $b_{i1} \sim N(0, 0.6^{1/2})$ ,  $\mu_{it} \sim N(0, 0.1^{1/2})$ ,  $\boldsymbol{\beta} = [0.5 \ 0.4 \ 0.6 \ -0.6 \ 0.7]^\top$ ,  $i = 1, 2, \dots, 50$ , and the sample size is 2,000 in an *unbalanced panel structure*. In addition, we assume that  $\mathbf{b}_i$  depends on  $\mathbf{z}_i = [1 \ z_{i1}]^\top$  such that  $\mathbf{b}_i \sim N(\mathbf{Z}_i \boldsymbol{\gamma}, \mathbf{D})$  where  $\mathbf{Z}_i = \mathbf{I}_{K_2} \otimes \mathbf{z}_i^\top$ , and  $\boldsymbol{\gamma} = [1 \ 1 \ 1 \ 1]$ . The prior for  $\boldsymbol{\gamma}$  is  $N(\boldsymbol{\gamma}_0, \boldsymbol{\Gamma}_0)$  where we set  $\boldsymbol{\gamma}_0 = \mathbf{0}_4$  and  $\boldsymbol{\Gamma}_0 = \mathbf{I}_4$ .

- Perform inference in this model without taking into account the dependence between  $\mathbf{b}_i$  and  $z_{i1}$ , and compare the posterior estimates with the population parameters.
- Perform inference in this model taking into account the dependence between  $\mathbf{b}_i$  and  $z_{i1}$ , and compare the posterior estimates with the population parameters.

### 4. Doctor visits in Germany continues I

Replicate this example using our GUI, which by default does not fix the over-dispersion parameter ( $\sigma^2$ ), and compare the results with the results of this example in Section 9.2.

### 5. Simulation exercise of the longitudinal logit model

Perform a simulation exercise to assess the performance of the hierarchical longitudinal logit model. The point of departure is to assume that

$$y_{it}^* = \beta_0 + \beta_1 x_{it1} + \beta_2 x_{it2} + \beta_3 x_{it3} + b_i + w_{it1} b_{i1},$$

where  $x_{itk} \sim N(0, 1)$ ,  $k = 1, 2, 3$ ,  $w_{it1} \sim N(0, 1)$ ,  $b_i \sim N(0, 0.7^{1/2})$ ,  $b_{i1} \sim N(0, 0.6^{1/2})$ ,  $\boldsymbol{\beta} = [0.5 \ 0.4 \ 0.6 \ -0.6]^\top$ ,  $i = 1, 2, \dots, 50$ , and  $y_{it} \sim B(\pi_{it})$ , where  $\pi_{it} = 1/(1 + \exp(y_{it}^*))$ . The sample size is 1,000 in an *unbalanced panel structure*.

- Perform inference using the command *MCMChlogit* fixing the over-dispersion parameter, and using  $\beta_0 = \mathbf{0}_4$ ,  $B_0 = \mathbf{I}_4$ ,  $\alpha_0 = \delta_0 = 0.001$ ,  $d_0 = 2$  and  $D_0 = \mathbf{I}_2$ .
- Program from scratch a Metropolis-within-Gibbs algorithm to perform inference in this simulation.

#### 6. Doctor visits in Germany continues II

Take a sub-sample of the first 500 individuals of the dataset *9VisitDoc.csv* to perform inference in the number of visits to doctors (*DocNum*) with the same specification of the example of **Doctor visits in Germany** of Section 9.2.



# 10

---

## *Bayesian model averaging*

---

We outline in this chapter a framework for addressing model uncertainty and averaging across different models in a probabilistically consistent manner. The discussion tackles two major computational challenges in Bayesian model averaging: the vast space of possible models and the absence of analytical solutions for the marginal likelihood.

We begin by illustrating the approach within the Gaussian linear model, assuming exogeneity of the regressors, and extend the analysis to cases with endogenous regressors, and dynamic models. Additionally, we adapt the framework to generalized linear models, including the logit, gamma, and Poisson families. Lastly, we explore alternative methods for computing marginal likelihoods, especially when the Bayesian information criterion's asymptotic approximation proves inadequate.

Remember that we can run our GUI typing

***R code. How to display our graphical user interface***

---

```
1 shiny::runGitHub("besmarter/BSTApp", launch.browser = T)
```

---

in the **R** package console or any **R** code editor, and once our GUI is deployed, select *Bayesian Model Averaging*. However, users should see Chapter 5 for other options and details.

---

### 10.1 Foundation

Remember from Chapter 1 that Bayesian model averaging (BMA) is an approach which takes into account model uncertainty. In particular, we consider uncertainty in the regressors (variable selection) in a regression framework

where there are  $K$  possible explanatory variables.<sup>1</sup> This implies  $2^K$  potential models indexed by parameters  $\boldsymbol{\theta}_m$ ,  $m = 1, 2, \dots, 2^K$ .

Following [411], the posterior model probability is

$$\pi(\mathcal{M}_j|\mathbf{y}) = \frac{p(\mathbf{y}|\mathcal{M}_j)\pi(\mathcal{M}_j)}{\sum_{m=1}^{2^K} p(\mathbf{y}|\mathcal{M}_m)\pi(\mathcal{M}_m)},$$

where  $\pi(\mathcal{M}_j)$  is the prior model probability,<sup>2</sup>

$$p(\mathbf{y}|\mathcal{M}_j) = \int_{\Theta_j} p(\mathbf{y}|\boldsymbol{\theta}_j, \mathcal{M}_j)\pi(\boldsymbol{\theta}_j|\mathcal{M}_j)d\boldsymbol{\theta}_j$$

is the marginal likelihood, and  $\pi(\boldsymbol{\theta}_j|\mathcal{M}_j)$  is the prior distribution of  $\boldsymbol{\theta}_j$  conditional on model  $\mathcal{M}_j$ .

Following [349], the posterior distribution of  $\boldsymbol{\theta}$  is

$$\pi(\boldsymbol{\theta}|\mathbf{y}) = \sum_{m=1}^{2^K} \pi(\boldsymbol{\theta}_m|\mathbf{y}, \mathcal{M}_m) \cdot \pi(\mathcal{M}_m|\mathbf{y})$$

The posterior distribution of the parameter vector  $\boldsymbol{\theta}$  under model  $\mathcal{M}_m$  is denoted as  $\pi(\boldsymbol{\theta}_m|\mathbf{y}, \mathcal{M}_m)$ . The posterior mean of  $\boldsymbol{\theta}$  is given by:

$$\mathbb{E}[\boldsymbol{\theta}|\mathbf{y}] = \sum_{m=1}^{2^K} \hat{\boldsymbol{\theta}}_m \pi(\mathcal{M}_m|\mathbf{y}),$$

where  $\hat{\boldsymbol{\theta}}_m$  represents the posterior mean under model  $\mathcal{M}_m$ .

The variance of the  $k$ -th element of  $\boldsymbol{\theta}$  given the data  $\mathbf{y}$  is:

$$\text{Var}(\theta_{km}|\mathbf{y}) = \sum_{m=1}^{2^K} \pi(\mathcal{M}_m|\mathbf{y}) \widehat{\text{Var}}(\theta_{km}|\mathbf{y}, \mathcal{M}_m) + \sum_{m=1}^{2^K} \pi(\mathcal{M}_m|\mathbf{y}) (\hat{\theta}_{km} - \mathbb{E}[\theta_{km}|\mathbf{y}])^2,$$

where  $\widehat{\text{Var}}(\theta_{km}|\mathbf{y}, \mathcal{M}_m)$  denotes the posterior variance of the  $k$ -th element of  $\boldsymbol{\theta}$  under model  $\mathcal{M}_m$ .

The posterior variance highlights how BMA accounts for model uncertainty. The first term represents the weighted variance of each model, averaged across all potential models, while the second term reflects the stability of the estimates across models. The greater the variation in estimates between models, the higher the posterior variance.

<sup>1</sup>Take into account that  $K$  can increase when interaction terms and/or polynomial terms of the original control variables are included.

<sup>2</sup>We attach equal prior probabilities to each model in our GUI. However, this choice gives more prior probability to the set of models of medium size (think about the  $k$ -th row of Pascal's triangle). An interesting alternative is to use the Beta-Binomial prior proposed by [267].

The posterior predictive distribution is

$$\pi(\mathbf{y}_0|\mathbf{y}) = \sum_{m=1}^{2^K} p_m(\mathbf{y}_0|\mathbf{y}, \mathcal{M}_m), \pi(\mathcal{M}_m|\mathbf{y})$$

where  $p_m(\mathbf{y}_0|\mathbf{y}, \mathcal{M}_m) = \int_{\Theta_m} p(\mathbf{y}_0|\mathbf{y}, \boldsymbol{\theta}_m, \mathcal{M}_m) \pi(\boldsymbol{\theta}_m|\mathbf{y}, \mathcal{M}_m) d\boldsymbol{\theta}_m$  is the posterior predictive distribution under model  $\mathcal{M}_m$ .

Another important statistic in BMA is the posterior inclusion probability associated with variable  $\mathbf{x}_k$ ,  $k = 1, 2, \dots, K$ , which is

$$PIP(\mathbf{x}_k) = \sum_{m=1}^{2^K} \pi(\mathcal{M}_m|\mathbf{y}) \times \mathbb{1}_{k,m},$$

where  $\mathbb{1}_{k,m} = \begin{cases} 1 & \text{if } \mathbf{x}_k \in \mathcal{M}_m \\ 0 & \text{if } \mathbf{x}_k \notin \mathcal{M}_m \end{cases}$ .

[241] suggest that posterior inclusion probabilities (PIP) less than 0.5 are evidence against the regressor,  $0.5 \leq PIP < 0.75$  is weak evidence,  $0.75 \leq PIP < 0.95$  is positive evidence,  $0.95 \leq PIP < 0.99$  is strong evidence, and  $PIP \geq 0.99$  is very strong evidence.

There are two main computational issues in implementing BMA based on variable selection. First, the number of models in the model space is  $2^K$ , which sometimes can be enormous. For instance, three regressors imply just eight models, see Table 10.1, but 40 regressors implies approximately  $1.1e+12$  models. Take into account that models always include the intercept, and all regressors should be standardized to avoid scale issues.<sup>3</sup> The second computational issue is calculating the marginal likelihood  $p(\mathbf{y}|\mathcal{M}_j) = \int_{\Theta_j} p(\mathbf{y}|\boldsymbol{\theta}_j, \mathcal{M}_j) \pi(\boldsymbol{\theta}_j|\mathcal{M}_j) d\boldsymbol{\theta}_j$ , which most of the time does not have an analytic solution.

**TABLE 10.1**  
Space of models: Three regressors.

Regressor	Inclusion							
	$x_1$	1	1	1	1	0	0	0
$x_2$	1	1	0	0	1	1	0	0
$x_3$	1	0	1	0	1	0	1	0

Notes: “1” indicates inclusion of the regressor, and “0” indicates no inclusion. The space of models is composed by 8 models. The model always includes intercept.

The first computational issue is basically a problem of ranking models. This can be tackled using different approaches, such as Occam’s window criterion [285, 351], reversible jump Markov chain Monte Carlo computation

<sup>3</sup>Scaling variables is always an important step in variable selection.

[182], Markov chain Monte Carlo model composition [286], and multiple testing using intrinsic priors [72] or nonlocal prior densities [227]. We focus on Occam’s window and Markov chain Monte Carlo model composition in our GUI.<sup>4</sup>

In Occam’s window, a model is discarded if its predictive performance is much worse than that of the best model [285, 351]. Thus, models not belonging to  $\mathcal{M}' = \left\{ \mathcal{M}_j : \frac{\max_m \pi(\mathcal{M}_m | \mathbf{y})}{\pi(\mathcal{M}_j | \mathbf{y})} \leq c \right\}$  should be discarded, where  $c$  is chosen by the user ([285] propose  $c = 20$ ). In addition, complicated models than are less supported by the data than simpler models are also discarded, that is,  $\mathcal{M}'' = \left\{ \mathcal{M}_j : \exists \mathcal{M}_m \in \mathcal{M}', \mathcal{M}_m \subset \mathcal{M}_j, \frac{\pi(\mathcal{M}_m | \mathbf{y})}{\pi(\mathcal{M}_j | \mathbf{y})} > 1 \right\}$ . Then, the set of models used in BMA is  $\mathcal{M}^* = \mathcal{M}' \cap \mathcal{M}''^c \in \mathcal{M}$ . [351] find that the number of models in  $\mathcal{M}^*$  is normally less than 25.

However, the previous theoretical framework requires finding the model with the maximum a posteriori model probability ( $\max_m \pi(\mathcal{M}_m | \mathbf{y})$ ), which implies calculating all possible models in  $\mathcal{M}$ . This is computationally burdensome. Hence, a heuristic approach is proposed by [348] based on ideas of [285]. The search strategy is based on a series of nested comparisons of ratios of posterior model probabilities. Let  $\mathcal{M}_0$  be a model with one regressor less than model  $\mathcal{M}_1$ , then:

- If  $\log(\pi(\mathcal{M}_0 | \mathbf{y}) / \pi(\mathcal{M}_1 | \mathbf{y})) > \log(O_R)$ , then  $\mathcal{M}_1$  is rejected and  $\mathcal{M}_0$  is considered.
- If  $\log(\pi(\mathcal{M}_0 | \mathbf{y}) / \pi(\mathcal{M}_1 | \mathbf{y})) \leq -\log(O_L)$ , then  $\mathcal{M}_0$  is rejected, and  $\mathcal{M}_1$  is considered.
- If  $\log(O_L) < \log(\pi(\mathcal{M}_0 | \mathbf{y}) / \pi(\mathcal{M}_1 | \mathbf{y})) \leq \log(O_R)$ ,  $\mathcal{M}_0$  and  $\mathcal{M}_1$  are considered.

Here  $O_R$  is a number specifying the maximum ratio for excluding models in Occam’s window, and  $O_L = 1/O_R^2$  is defined by default in [348]. The search strategy can be “up”, adding one regressor, or “down”, dropping one regressor (see [285] for details about the down and up algorithms). The leaps and bounds algorithm [153] is implemented to improve the computational efficiency of this search strategy [348]. Once the set of potentially acceptable models is defined, we discard all the models that are not in  $\mathcal{M}'$ , and the models that are in  $\mathcal{M}''$  where 1 is replaced by  $\exp\{O_R\}$  due to the leaps and bounds algorithm giving an approximation to BIC, so as to ensure that no good models are discarded.

The second approach that we consider in our GUI to tackle the model space size issue is Markov chain Monte Carlo model composition (MC3) [287]. In particular, given the space of models  $\mathcal{M}_m$ , we simulate a chain of  $\mathcal{M}_s$  models,  $s = 1, 2, \dots, S << 2^K$ , where the algorithm randomly extracts a candidate

---

<sup>4</sup>Variable selection (model selection or regularization) is a topic related to model uncertainty. Approaches such as stochastic search variable selection (spike and slab) [166, 167] and Bayesian Lasso [330] are good examples of how to tackle this issue. See Chapter 12.

model  $\mathcal{M}_c$  from a neighborhood of models ( $nbd(\mathcal{M}_m)$ ) that consists of the actual model itself and the set of models with either one variable more or one variable less [351]. Therefore, there is a transition kernel in the space of models  $q(\mathcal{M}_m \rightarrow \mathcal{M}_c)$ , such that  $q(\mathcal{M}_m \rightarrow \mathcal{M}_c) = 0 \forall \mathcal{M}_c \notin nbd(\mathcal{M}_m)$  and  $q(\mathcal{M}_m \rightarrow \mathcal{M}_c) = \frac{1}{|nbd(\mathcal{M}_m)|} \forall \mathcal{M}_m \in nbd(\mathcal{M}_m)$ ,  $|nbd(\mathcal{M}_m)|$  being the number of neighbors of  $\mathcal{M}_m$ . This candidate model is accepted with probability

$$\alpha(\mathcal{M}_{s-1}, \mathcal{M}_c) = \min \left\{ \frac{|nbd(\mathcal{M}_m)| p(\mathbf{y}|\mathcal{M}_c) \pi(\mathcal{M}_c)}{|nbd(\mathcal{M}^c)| p(\mathbf{y}|\mathcal{M}_{(s-1)}) \pi(\mathcal{M}_{(s-1)})}, 1 \right\}.$$

Observe that by construction  $|nbd(\mathcal{M}_m)| = |nbd(\mathcal{M}_c)| = K$ , except in extreme cases where a model has only one regressor or has all regressors.

The Bayesian information criterion is a possible solution for the second computational issue in BMA, that is, calculating the marginal likelihood when there is no an analytic solution. Defining  $h(\boldsymbol{\theta}|\mathcal{M}_j) = -\frac{\log(p(\mathbf{y}|\boldsymbol{\theta}_j, \mathcal{M}_j) \pi(\boldsymbol{\theta}_j|\mathcal{M}_j))}{N}$ , then  $p(\mathbf{y}|\mathcal{M}_j) = \int_{\boldsymbol{\Theta}_j} \exp \{-Nh(\boldsymbol{\theta}|\mathcal{M}_j)\} d\boldsymbol{\theta}_j$ . If  $N$  is sufficiently large (technically  $N \rightarrow \infty$ ), we can make the following assumptions [208]:

- We can use the Laplace method for approximating integrals [434].
- The posterior mode is reached at the same point as the maximum likelihood estimator (MLE), denoted by  $\hat{\boldsymbol{\theta}}_{MLE}$ .

We get the following results under these assumptions:

$$p(\mathbf{y}|\mathcal{M}_j) \approx \left( \frac{2\pi}{N} \right)^{K_j/2} |\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -Nh(\hat{\boldsymbol{\theta}}_j^{MLE}|\mathcal{M}_j) \right\}, \quad N \rightarrow \infty,$$

where  $\boldsymbol{\Sigma}$  is the inverse of the Hessian matrix of  $h(\hat{\boldsymbol{\theta}}_j^{MLE}|\mathcal{M}_j)$ , and  $K_j = \dim \{\boldsymbol{\theta}_j\}$ .

This implies

$$\begin{aligned} \log(p(\mathbf{y}|\mathcal{M}_j)) &\approx \frac{K_j}{2} \log(2\pi) - \frac{K_j}{2} \log(N) - \frac{1}{2} \log(|\boldsymbol{\Sigma}|) + \log(p(\mathbf{y}|\hat{\boldsymbol{\theta}}_j^{MLE}, \mathcal{M}_j)) \\ &\quad + \log(\pi(\hat{\boldsymbol{\theta}}_j^{MLE}|\mathcal{M}_j)), \quad N \rightarrow \infty. \end{aligned}$$

Since  $\frac{K_j}{2} \log(2\pi)$  and  $\log(\pi(\hat{\boldsymbol{\theta}}_j^{MLE}|\mathcal{M}_j))$  are constants as functions of  $\mathbf{y}$ , and  $|\boldsymbol{\Sigma}|$  is bounded by a finite constant, we have

$$\log(p(\mathbf{y}|\mathcal{M}_j)) \approx -\frac{K_j}{2} \log(N) + \log(p(\mathbf{y}|\hat{\boldsymbol{\theta}}_j^{MLE}, \mathcal{M}_j)) = -\frac{BIC}{2}, \quad N \rightarrow \infty.$$

The marginal likelihood thus asymptotically converges to a linear transformation of the Bayesian Information Criterion (BIC), significantly simplifying its calculation. In addition, the BIC is consistent, that is, the probability of uncovering the population statistical model converges to one as the sample size

converges to infinity given a  $\mathcal{M}$ -closed view [41, Chap. 6], that is, one of the models in consideration is the population statistical model (data generating process) [400, 64]. In case that there is an  $\mathcal{M}$ -completed view of nature, that is, there is a true data generating process, but the space of models that we are comparing does not include it, the BIC asymptotically selects the model that minimizes the Kullback-Leiber (KL) divergence to the true (population) model [102, Chap. 4].

## 10.2 The Gaussian linear model

The Gaussian linear model specifies  $\mathbf{y} = \alpha \mathbf{i}_N + \mathbf{X}_m \boldsymbol{\beta}_m + \boldsymbol{\mu}_m$  such that  $\boldsymbol{\mu}_m \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ ,  $\mathbf{i}_N$  is a  $N$ -dimensional column of ones,  $\mathbf{X}_m$  is the design matrix without the column of ones, and  $N$  is the sample size. Following [250], the conjugate prior for the location parameters is  $\boldsymbol{\beta}_m | \sigma^2 \sim N(\boldsymbol{\beta}_{m0}, \sigma^2 \mathbf{B}_{m0})$ , and the priors for  $\sigma^2$  and  $\alpha$  can be improper, as these parameters are common to all models  $\mathcal{M}_m$ . Particularly,  $\pi(\sigma^2) \propto 1/\sigma^2$  (Jeffreys' prior for the linear Gaussian model, see [213]) and  $\pi(\alpha) \propto 1$ .

The selection of the hyperparameters of  $\boldsymbol{\beta}_m$  is more critical, as these parameters are not common to all models. A very common prior for the location parameters in the BMA literature is the Zellner's prior [466], where  $\boldsymbol{\beta}_{m0} = \mathbf{0}_m$  and  $\mathbf{B}_{m0} = (g_m \mathbf{X}_m^\top \mathbf{X}_m)^{-1}$ . Observe that this covariance matrix is similar to the covariance matrix of the ordinary least squares estimator of the location parameters. This suggests that there is compatibility between the prior information and the sample information, and the only parameter to elicit is  $g_m \geq 0$ , which facilitates the elicitation process, as eliciting covariance matrices is a very hard endeavor.

Following same steps as in Section 3.3, the posterior conditional distribution of  $\boldsymbol{\beta}_m$  has covariance matrix  $\sigma^2 \mathbf{B}_{mn}$ , where  $\mathbf{B}_{mn} = ((1 + g_m) \mathbf{X}_m^\top \mathbf{X}_m)^{-1}$  (Exercise 1), which means that  $g_m = 0$  implies a non-informative prior, whereas  $g_m = 1$  implies that prior and data information have same weights. We follow [139], who recommend

$$g_m = \begin{cases} 1/K^2, & N \leq K^2 \\ 1/N, & N > K^2 \end{cases},$$

where  $K$  is the number of regressors.

Given the likelihood function,

$$\begin{aligned} p(\boldsymbol{\beta}_m, \sigma^2 | \mathbf{y}, \mathbf{X}_m) &= (2\pi\sigma^2)^{-\frac{N}{2}} \\ &\times \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \alpha \mathbf{i}_N - \mathbf{X}_m \boldsymbol{\beta}_m)^\top (\mathbf{y} - \alpha \mathbf{i}_N - \mathbf{X}_m \boldsymbol{\beta}_m) \right\}, \end{aligned}$$

the marginal likelihood associated with model  $\mathcal{M}_m$  is proportional to (Exercise

1)

$$p(\mathbf{y}|\mathcal{M}_m) \propto \left( \frac{g_m}{1+g_m} \right)^{k_m/2} \left[ (\mathbf{y} - \bar{y}\mathbf{i}_N)^\top (\mathbf{y} - \bar{y}\mathbf{i}_N) - \frac{1}{1+g_m} (\mathbf{y}^\top \mathbf{P}_{X_m} \mathbf{y}) \right]^{-(N-1)/2},$$

where all parameters are indexed to model  $\mathcal{M}_m$ ,  $\mathbf{P}_{X_m} = \mathbf{X}_m (\mathbf{X}_m^\top \mathbf{X}_m)^{-1} \mathbf{X}_m$  is the projection matrix on the space generated by the columns of  $\mathbf{X}_m$ , and  $\bar{y}$  is the sample mean of  $\mathbf{y}$ .

We implement in our GUI four approaches to perform BMA in the Gaussian linear model: the BIC approximation using the Occam's window approach, the MC3 algorithm using the analytical expression for calculating the marginal likelihood, an instrumental variable approach based on conditional likelihoods, and dynamic variable selection.

#### Example: Simulation exercise

Let's perform a simulation exercise to assess the performance of the BIC approximation using the Occam's window, and the Markov chain Monte Carlo model composition approaches. Let's set a model where the computational burden is low and we know the data generating process (population statistical model). In particular, we set 10 regressors such that  $x_k \sim N(1, 1)$ ,  $k = 1, \dots, 6$ , and  $x_k \sim B(0.5)$ ,  $k = 7, \dots, 10$ . We set  $\beta = [1 \ 0 \ 0 \ 0.5 \ 0 \ 0, 0, 0, -0.7]^\top$  such that just  $x_1, x_5$  and  $x_{10}$  are relevant to drive  $y_i = 1 + \mathbf{x}^\top \beta + \mu_i$ ,  $\mu_i \sim N(0, 0.5^2)$ . Observe that we just have  $2^{10} = 1,024$  models in this setting, thus, we can calculate the posterior model probability for each model.

Our GUI uses the commands *bicreg* and *MC3.REG* from the package *BMA* to perform Bayesian model averaging in the linear regression model using the BIC approximation and MC3, respectively. These commands in turn are based on [347] and [351]. The following code shows how to perform the simulation and get the posterior mean and standard deviation using these commands with the default values of hyperparameters and tuning parameters.

#### R code. Simulation exercise: Bayesian model averaging, small setting

```

1 rm(list = ls()); set.seed(10101)
2 N <- 1000
3 K1 <- 6; K2 <- 4; K <- K1 + K2
4 X1 <- matrix(rnorm(N*K1, 1, 1), N, K1)
5 X2 <- matrix(rbinom(N*K2, 1, 0.5), N, K2)

```

**R code. Simulation exercise: Bayesian model averaging, small setting**

```

1 X <- cbind(X1, X2); e <- rnorm(N, 0, 0.5)
2 B <- c(1,0,0,0,0.5,0,0,0,0,-0.7)
3 y <- 1 + X%*%B + e
4 BMAGlm <- BMA::bicreg(X, y, strict = FALSE, OR = 50)
5 summary(BMAGlm)
6 BMAREG <- BMA::MC3.REG(y, X, num.its=500)
7 Models <- unique(BMAREG[["variables"]])
8 nModels <- dim(Models)[1]
9 nVistModels <- dim(BMAREG[["variables"]])[1]
10 PMP <- NULL
11 for(m in 1:nModels){
12   idModm <- NULL
13   for(j in 1:nVistModels){
14     if(sum(Models[m,] == BMAREG[["variables"]][j,]) == K){
15       idModm <- c(idModm, j)
16     }else{
17       idModm <- idModm
18     }
19   }
20   PMPm <- sum(BMAREG[["post.prob"]][idModm])
21   PMP <- c(PMP, PMPm)
22 }
23 PIP <- NULL
24 for(k in 1:K){
25   PIPk <- sum(PMP[which(Models[,k] == 1)])
26   PIP <- c(PIP, PIPk)
27 }
28 plot(PIP)
29 Means <- matrix(0, nModels, K)
30 Vars <- matrix(0, nModels, K)
31 for(m in 1:nModels){
32   idXs <- which(Models[m,] == 1)
33   if(length(idXs) == 0){
34     Regm <- lm(y ~ 1)
35   }else{
36     Xm <- X[, idXs]
37     Regm <- lm(y ~ Xm)
38     SumRegm <- summary(Regm)
39     Means[m, idXs] <- SumRegm[["coefficients"]][-1,1]
40     Vars[m, idXs] <- SumRegm[["coefficients"]][-1,2]^2
41   }
42 }
43 BMAMEANS <- colSums(Means*PMP)
44 BMASD <- (colSums(PMP*Vars) + colSums(PMP*(Means-matrix(rep
        (BMAMEANS, each = nModels), nModels, K))^2))^0.5
45 BMAMEANS; BMASD; BMAMEANS/BMASD

```

We can see from the results that the BIC approximation with the Occam's window, and the MC3 algorithm perform a good job finding the relevant regressors, and their posterior BMA means are very close to the population values. We also see that the BMA results are very similar in the two approaches.

We can perform Bayesian model averaging in our GUI for linear Gaussian models using the BIC approximation and MC3 using Algorithms A28 and A29, respectively. We ask in Exercise 2 to perform BMA using the dataset *10ExportDiversificationHHI.csv* from [225].

---

**Algorithm A28** Bayesian model averaging in linear Gaussian models using the Bayesian information criterion

---

- 1: Select *Bayesian Model Averaging* on the top panel
  - 2: Select *Normal data* model using the left radio button
  - 3: Select *BIC* using the right radio button under **Which type do you want to perform?**
  - 4: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 5: Type the *OR* number of the Occam's window in the box under **OR: Number between 5 and 50**, this is not necessary as by default there is 50
  - 6: Click the *Go!* button
  - 7: Analyze results: After a few seconds or minutes, a table appears showing, for each regressor in the dataset, the PIP (posterior inclusion probability, **p!=0**), the BMA posterior mean (**EV**), the BMA standard deviation (**SD**), and the posterior mean for models with the highest PMP. At the bottom of the table, for the models with the largest PMP, the number of variables (**nVar**), the coefficient of determination (**r2**), the BIC, and the PMP (**post prob**) are displayed
  - 8: Download posterior results using the *Download results using BIC*. There are two files, the first has the best models by row according to the PMP (last column) indicating with a 1 inclusion of the variable (0 indicates no inclusion), and the second file has the PIP, the BMA expected value and standard deviation for each variable in the dataset
- 

We show in the following code how to program a MC3 algorithm from scratch to perform BMA using the setting from the previous simulation exercise. The first part of the code is the function to calculate the log marginal likelihood. This is a small simulation setting, thus we can calculate the marginal likelihood for all 1,024 models, and then calculate the posterior model probability standardizing using the model with the largest log marginal likelihood. We see from the results that this model is the data generating process (population statistical model). We also find that the posterior inclusion probabilities

---

**Algorithm A29** Bayesian model averaging in linear Gaussian models using Markov chain Monte Carlo model composition

---

- 1: Select *Bayesian Model Averaging* on the top panel
  - 2: Select *Normal data* model using the left radio button
  - 3: Select *MC3* using the right radio button under **Which type do you want to perform?**
  - 4: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 5: Select MC3 iterations using the *Range slider* under the label **MC3 iterations:**
  - 6: Click the *Go!* button
  - 7: Analyze results: After a few seconds or minutes, a table appears showing, for each regressor in the dataset, the PIP (posterior inclusion probability, **p!=0**), the BMA posterior mean (**EV**), the BMA standard deviation (**SD**), and the posterior mean for models with the highest PMP. At the bottom of the table, for the models with the largest PMP, the number of variables (**nVar**), the coefficient of determination (**r2**), the BIC, and the PMP (**post prob**) are displayed
  - 8: Download posterior results using the *Download results using BIC*. There are two files, the first has the best models by row according to the PMP (last column) indicating with a 1 inclusion of the variable (0 indicates no inclusion), and the second file has the PIP, the BMA expected value and standard deviation for each variable in the dataset
- 

for  $x_1$ ,  $x_5$  and  $x_{10}$  are 1, whereas the PIP for the other variables are less than 0.05.

Although BMA allows incorporating model uncertainty in a regression framework, sometimes it is desirable to select just one model. Two compelling alternatives are the model with the largest posterior model probability, and the median probability model. The latter is the model which includes every predictor that has posterior inclusion probability higher than 0.5. The first model is the best alternative for prediction in the case of a 0–1 loss function [103], whereas the second is the best alternative when there is a quadratic loss function in prediction [25]. In this simulation, the two criteria indicate selection of the data generating process.

We also show how to estimate the posterior mean and standard deviation based on BMA in this code. We see that the posterior means are very close to the population parameters.

**R code. Simulation exercise: Bayesian model averaging, small setting from scratch**

```

1 rm(list = ls()); set.seed(010101)
2 N <- 1000; K1 <- 6; K2 <- 4; K <- K1 + K2
3 X1 <- matrix(rnorm(N*K1, 1, 1), N, K1)
4 X2 <- matrix(rbinom(N*K2, 1, 0.5), N, K2)
5 X <- cbind(X1, X2); e <- rnorm(N, 0, 0.5)
6 B <- c(1, 0, 0, 0, 0.5, 0, 0, 0, 0, -0.7); y <- 1 + X%*%B + e
7 LogMLfunt <- function(Model){
8   indr <- Model == 1
9   kr <- sum(indr)
10  if(kr > 0){
11    gr <- ifelse(N > kr^2, 1/N, kr^(-2))
12    Xr <- matrix(Xnew[, indr], ncol = kr)
13    PX <- Xr%*%solve(t(Xr)%*%Xr)%*%t(Xr)
14    s2pos <- c((t(y - mean(y))%*%(y - mean(y))) - t(y)%*%PX%
15      *y/(1 + gr))
16    mllMod <- (kr/2)*log(gr/(1+gr))-(N-1)/2*log(s2pos)
17  }else{
18    gr <- ifelse(N > kr^2, 1/N, kr^(-2))
19    s2pos <- c((t(y - mean(y))%*%(y - mean(y))))
20    mllMod <- (kr/2)*log(gr/(1+gr))-(N-1)/2*log(s2pos)
21  }
22  return(mllMod)
23 }
24 combs <- expand.grid(c(0,1), c(0,1), c(0,1), c(0,1), c(0,1),
25   c(0,1), c(0,1), c(0,1), c(0,1), c(0,1))
26 Xnew <- apply(X, 2, scale)
27 mll <- sapply(1:2^K, function(s){LogMLfunt(matrix(combs[s,],
28   1, K))})
29 MaxPMP <- which.max(mll); StMarLik <- exp(mll-max(mll))
30 PMP <- StMarLik/sum(StMarLik)
31 PMP[MaxPMP]; combs[MaxPMP,]
32 PIP <- NULL
33 for(k in 1:K){
34   PIPk <- sum(PMP[which(combs[,k] == 1)]); PIP <- c(PIP,
35     PIPk)
36 }
37 PIP
38 nModels <- dim(combs)[1]; Means <- matrix(0, nModels, K)

```

**R code. Simulation exercise: Bayesian model averaging, small setting from scratch**

```

1 Vars <- matrix(0, nModels, K)
2 for(m in 1:nModels){
3   idXs <- which(combs[m,] == 1)
4   if(length(idXs) == 0){
5     Regm <- lm(y ~ 1)
6   }else{
7     Xm <- X[, idXs]; Regm <- lm(y ~ Xm)
8     SumRegm <- summary(Regm)
9     Means[m, idXs] <- SumRegm[["coefficients"]][,-1]
10    Vars[m, idXs] <- SumRegm[["coefficients"]][,-1]^2
11  }
12 }
13 BMAmens <- colSums(Means*PMP)
14 BMAsd <- (colSums(PMP*Vars) + colSums(PMP*(Means-matrix(rep
15           (BMAmens, each = nModels), nModels, K))^2))^.5

```

The following part of the code demonstrates how to perform the MC3 algorithm. While this algorithm is not strictly necessary for this small-dimensional problem, it serves as a useful pedagogical exercise. The starting point is to set  $S = 100$  random models and order their log marginal likelihoods. The logic of the algorithm is to select the worst model among the  $S$  models and propose a candidate model to compete against it. We repeat this process for 1,000 iterations (as shown in the code). Note that 1,000 iterations is fewer than the number of potential models (1,024). This is the essence of the MC3 algorithm: performing fewer iterations than the number of models in the space.<sup>5</sup>

In our algorithm, we analyze all model scenarios using different conditionals and reasonably assume the same prior model probability for all models, with the same cardinality for both the actual and candidate models. The posterior model probability (PMP) can be calculated in several ways. One method is to recover the unique models from the final set of  $S$  models, calculate the log marginal likelihood for these models, and then standardize by the best model among them. Another method involves calculating the PMP using the complete set of  $S$  final models, accounting for the fact that some models may appear multiple times in the set, which requires summing the PMPs of repeated models. A third method is to calculate the PMP based on the relative frequency with which a model appears in the final set of  $S$  models. These three methods can yield different PMPs, particularly when the number

<sup>5</sup>Note that we have split the *while* loop in the code into two pages.

of MC3 iterations is small. In our example, using 1,000 MC3 iterations, the data-generating process receives the highest PMP across all three methods.

A noteworthy aspect of this algorithm is that we can obtain a single model after significantly increasing the number of iterations (for example, try using 10,000 iterations). This can be advantageous if we require only one model. However, this approach neglects model uncertainty, which could be a desirable characteristic in some cases. As a challenge, we suggest programming an algorithm that yields  $S$  different models after completing the MC3 iterations (Exercise 3).

### R code. Simulation exercise: Bayesian model averaging, small setting from scratch

```

1 M <- 100
2 Models <- matrix(rbinom(K*M, 1, p = 0.5), ncol=K, nrow = M)
3 mllnew <- sapply(1:M, function(s){LogMLfunt(matrix(Models[s,
], 1, K))})
4 oind <- order(mllnew, decreasing = TRUE)
5 mllnew <- mllnew[oind]; Models <- Models[oind, ]; iter <-
1000
6 pb <- txtProgressBar(min = 0, max = iter, style = 3); s <- 1
7 while(s <= iter){
8   ActModel <- Models[M,]; idK <- which(ActModel == 1)
9   Kact <- length(idK)
10  if(Kact < K & Kact > 1){
11    CardMol <- K; opt <- sample(1:3, 1)
12    if(opt == 1){ # Same
13      CandModel <- ActModel
14    }else{
15      if(opt == 2){ # Add
16        All <- 1:K; NewX <- sample(All[-idK], 1)
17        CandModel <- ActModel; CandModel[NewX] <- 1
18      }else{ # Subtract
19        LessX <- sample(idK, 1); CandModel <- ActModel
20        CandModel[LessX] <- 0
21      }
22    }
23  }else{
24    CardMol <- K + 1
25    if(Kact == K){
26      opt <- sample(1:2, 1)
27      if(opt == 1){ # Same
28        CandModel <- ActModel
29      }else{ # Subtract
30        LessX <- sample(1:K, 1); CandModel <- ActModel
31        CandModel[LessX] <- 0
32      }
33    }
34  }
35  if(pb < s) print(s)
36  s <- s + 1
37}

```

**R code. Simulation exercise: Bayesian model averaging, small setting from scratch**

```

1 }else{
2   if(K == 1){
3     opt <- sample(1:3, 1)
4     if(opt == 1){ # Same
5       CandModel <- ActModel
6     }else{
7       if(opt == 2){ # Add
8         All <- 1:K; NewX <- sample(All[-idK], 1)
9         CandModel <- ActModel; CandModel[NewX] <- 1
10      }else{ # Subtract
11        LessX <- sample(idK, 1); CandModel <- ActModel
12        CandModel[LessX] <- 0
13      }
14    }
15  }else{ # Add
16    NewX <- sample(1:K, 1); CandModel <- ActModel
17    CandModel[NewX] <- 1
18  }
19 }
20 }
21 LogMLact <- LogMLfunct(matrix(ActModel, 1, K))
22 LogMLcand <- LogMLfunct(matrix(CandModel, 1, K))
23 alpha <- min(1, exp(LogMLcand-LogMLact))
24 u <- runif(1)
25 if(u <= alpha){
26   mllnew[M] <- LogMLcand; Models[M, ] <- CandModel
27   oind <- order(mllnew, decreasing = TRUE)
28   mllnew <- mllnew[oind]; Models <- Models[oind, ]
29 }else{
30   mllnew <- mllnew; Models <- Models
31 }
32 s <- s + 1
33 setTxtProgressBar(pb, s)
34 }
35 close(pb)
36 ModelsUni <- unique(Models)
37 mllnewUni <- sapply(1:dim(ModelsUni)[1], function(s){
38   LogMLfunct(matrix(ModelsUni[s,], 1, K))})
39 StMarLik <- exp(mllnewUni-mllnewUni[1])
40 PMP <- StMarLik/sum(StMarLik) # PMP based on unique selected
41   models
42 nModels <- dim(ModelsUni)[1]
43 StMarLik <- exp(mllnew-mllnew[1])
44 PMPold <- StMarLik/sum(StMarLik) # PMP all selected models
45 PMPot <- NULL
46 PMPar <- NULL
47 FreqMod <- NULL

```

**R code. Simulation exercise: Bayesian model averaging, small setting from scratch**

```

1  for(m in 1:nModels){
2    idModm <- NULL
3    for(j in 1:M){
4      if(sum(ModelsUni[m,] == Models[j,]) == K){
5        idModm <- c(idModm, j)
6      }else{
7        idModm <- idModm
8      }
9    }
10   PMPm <- sum(PMPold[idModm]) # PMP unique models using sum
11   of all selected models
12   PMPot <- c(PMPot, PMPm)
13   PMPapm <- length(idModm)/M # PMP using relative frequency
14   in all selected models
15   PMPap <- c(PMPap, PMPapm)
16   FreqMod <- c(FreqMod, length(idModm))
17 }
18 PIP <- NULL
19 for(k in 1:K){
20   PIPk <- sum(PMP[which(ModelsUni[,k] == 1)])
21   PIP <- c(PIP, PIPk)
22 }
23 Means <- matrix(0, nModels, K)
24 Vars <- matrix(0, nModels, K)
25 for(m in 1:nModels){
26   idXs <- which(ModelsUni[m,] == 1)
27   if(length(idXs) == 0){
28     Regm <- lm(y ~ 1)
29   }else{
30     Xm <- X[, idXs]
31     Regm <- lm(y ~ Xm)
32     SumRegm <- summary(Regm)
33     Means[m, idXs] <- SumRegm[["coefficients"]][,-1]
34     Vars[m, idXs] <- SumRegm[["coefficients"]][,-1]^2
35   }
36 }
37 BMAmmeans <- colSums(Means*PMP)
38 BMAsd <- (colSums(PMP*Vars) + colSums(PMP*(Means-matrix(rep
39   (BMAmmeans, each = nModels), nModels, K))^2))^0.5

```

An important issue to account for regressors (model) uncertainty in the identification of causal effects, rather than finding good predictors (association relationships), is endogeneity. Thus, we also implement the instrumental variable approach of Section 7.3 to tackle this issue in BMA. We assume that  $\gamma \sim N(\mathbf{0}, \mathbf{I})$ ,  $\beta \sim N(\mathbf{0}, \mathbf{I})$ , and  $\Sigma^{-1} \sim W(3, \mathbf{I})$  [239].

[262] propose an algorithm based on conditional Bayes factors [119] that allows embedding MC3 within a Gibbs sampling algorithm. Given the candidate ( $M_c^{2nd}$ ) and actual ( $M_{s-1}^{2nd}$ ) models for the iteration  $s$  in the second stage, the conditional Bayes factor is

$$CBF^{2nd} = \frac{p(\mathbf{y}|M_c^{2nd}, \boldsymbol{\gamma}, \boldsymbol{\Sigma})}{p(\mathbf{y}|M_{s-1}^{2nd}, \boldsymbol{\gamma}, \boldsymbol{\Sigma})},$$

where

$$\begin{aligned} p(\mathbf{y}|M_c^{2nd}, \boldsymbol{\gamma}, \boldsymbol{\Sigma}) &= \int_{\mathcal{M}^{2nd}} p(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\Sigma}) \pi(\boldsymbol{\beta}|M_c^{2nd}) d\boldsymbol{\beta} \\ &\propto |\mathbf{B}_n|^{-1/2} \exp \left\{ \frac{1}{2} \boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n \right\}. \end{aligned}$$

In the first stage,

$$CBF^{1st} = \frac{p(\mathbf{y}|M_c^{1st}, \boldsymbol{\beta}, \boldsymbol{\Sigma})}{p(\mathbf{y}|M_{s-1}^{1st}, \boldsymbol{\beta}, \boldsymbol{\Sigma})},$$

where

$$\begin{aligned} p(\mathbf{y}|M_c^{1st}, \boldsymbol{\beta}, \boldsymbol{\Sigma}) &= \int_{\mathcal{M}^{1st}} p(\mathbf{y}|\boldsymbol{\gamma}, \boldsymbol{\beta}, \boldsymbol{\Sigma}) \pi(\boldsymbol{\gamma}|M_c^{1st}) d\boldsymbol{\gamma} \\ &\propto |\mathbf{G}_n|^{-1/2} \exp \left\{ \frac{1}{2} \boldsymbol{\gamma}_n^\top \mathbf{G}_n^{-1} \boldsymbol{\gamma}_n \right\}. \end{aligned}$$

These conditional Bayes factors assume  $\pi(M^{1st}, M^{2nd}) \propto 1$ . See [262] for more details of the instrumental variable BMA algorithm.<sup>6</sup>

We perform instrumental variable BMA in our GUI using the package *ivbma*. The Algorithm A30 shows how to perform this in our GUI.

Let's perform a simulation exercise to assess the performance of the instrumental variable BMA to uncover the data generating process in presence of endogeneity.

#### Example: Simulation exercise

Let's assume that  $y_i = 2 + 0.5x_{i1} - x_{i2} + x_{i3} + \mu_i$ , where  $x_{i1} = 4z_{i1} - z_{i2} + 2z_{i3} + \epsilon_{i1}$  and  $x_{i2} = -2z_{i1} + 3z_{i2} - z_{i3} + \epsilon_{i2}$ , such that  $[\epsilon_{i1} \ \epsilon_{i2} \ \mu_i]^\top \sim N(\mathbf{0}, \boldsymbol{\Sigma})$ ,

where  $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0 & 0.8 \\ 0 & 1 & 0.5 \\ 0.8 & 0.5 & 1 \end{bmatrix}$ , for  $i = 1, 2, \dots, 1000$ . The endogeneity arises due

to the correlation between  $\mu_i$  and  $x_{i1}$  and  $x_{i2}$  through the stochastic errors. In addition, there are three instruments,  $z_{il} \sim U(0, 1)$ , for  $l = 1, 2, 3$ , and another 18 regressors believed to influence  $y_i$ , which are distributed according to a standard normal distribution.

<sup>6</sup>[247] and [261] propose other frameworks for BMA taking into account endogeneity.

---

**Algorithm A30** Instrumental variable Bayesian model averaging in linear Gaussian models
 

---

- 1: Select *Bayesian Model Averaging* on the top panel
  - 2: Select *Normal data* model using the left radio button
  - 3: Select *Instrumental variable* using the right radio button under **Which type do you want to perform?**
  - 4: Upload the dataset containing the dependent variable, endogenous regressors, and exogenous regressors including the constant (see Section 5.6 for details). User should select first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 5: Upload the dataset containing the instruments (see Section 5.6 for details). User should select first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File (Instruments)** legend
  - 6: Write down the number of endogenous regressors in the box labeled **Number of Endogenous variables**
  - 7: Select MCMC iterations and burn-in using the *Range slider* under the labels **MCMC iterations** and **Burn-in Sample**
  - 8: Click the *Go!* button
  - 9: Analyze results: After a few seconds or minutes, two tables appear showing, for each regressor in the dataset, the PIP (posterior inclusion probability, **p!=0**), and the BMA posterior mean (**EV**). The top table shows the results of the second stage (main equation), and the bottom table shows the results of the first stage (auxiliary equations)
  - 10: Download posterior results using the *Download results using IV*. There are three files, the first file has the posterior inclusion probabilities of each variable, and the BMA posterior means of the coefficients in the first stage equations, the second file shows these results for the second stage (main equation), and the third file has the posteriors chains of all parameters by iteration.
- 

The following code shows how to perform IV BMA using the *ivbma* package. We see from the results that the PIP of  $x_{i1}$ ,  $x_{i2}$ , intercept and  $x_{i3}$  are equal to 1, whereas the remaining PIP are close to 0. In addition, the BMA means are also close to the population values. The PIP of the first stage equations, as well as their BMA posterior means, are very close to the populations values. The same happens with the covariance matrix.

We ask in Exercise 4 to perform BMA based on the BIC approximation and MC3 in this simulation setting. In addition, we ask in Exercise 5 to use the datasets *11ExportDiversificationHHI.csv* and *12ExportDiversificationHHIInstr.csv* to perform IV BMA assuming that the log of per capita gross domestic product is endogenous (*avgldpcap*). See [225] for details.

**R code. Simulation exercise: Instrumental variable Bayesian model averaging**

```

1 rm(list = ls())
2 set.seed(010101)
3 simIV <- function(delta1,delta2,beta0,betas1,betas2,beta2,
4 Sigma,n,z) {
5   eps <- matrix(rnorm(3*n),ncol=3) %*% chol(Sigma)
6   xs1 <- z%*%delta1 + eps[,1]
7   xs2 <- z%*%delta2 + eps[,2]
8   x2 <- rnorm(dim(z)[1])
9   y <- beta0+betas1*xs1+betas2*xs2+beta2*x2 + eps[,3]
10  X <- as.matrix(cbind(xs1,xs2,1,x2))
11  colnames(X) <- c("x1en","x2en","cte","xex")
12  y <- matrix(y,dim(z)[1],1)
13  colnames(y) <- c("y")
14  list(X=X,y=y)
15 }
16 n <- 1000 ; p <- 3
17 z <- matrix(runif(n*p),ncol=p)
18 rho31 <- 0.8; rho32 <- 0.5;
19 Sigma <- matrix(c(1,0,rho31,0,1,rho32,rho31,rho32,1),ncol=3)
20 delta1 <- c(4,-1,2); delta2 <- c(-2,3,-1); betas1 <- .5;
21       betas2 <- -1; beta2 <- 1; beta0 <- 2
22 simiv <- simIV(delta1,delta2,beta0,betas1,betas2,beta2,Sigma
23 ,n,z)
24 nW <- 18
25 W <- matrix(rnorm(nW*dim(z)[1]),dim(z)[1],nW)
26 YXW<-cbind(simiv$y, simiv$X, W)
27 y <- YXW[,1]; X <- YXW[,2:3]; W <- YXW[,-c(1:3)]
28 S <- 10000; burnin <- 1000
29 regivBMA <- ivbma::ivbma(Y = y, X = X, Z = z, W = W, s = S+
30   burnin, b = burnin, odens = S, print.every = round(S/10)
31   , run.diagnostics = FALSE)
32 PIPmain <- regivBMA[["L.bar"]] # PIP outcome
33 EVmain <- regivBMA[["rho.bar"]] # Posterior mean outcome
34 PIPaux <- regivBMA[["M.bar"]] # PIP auxiliary
35 EVAux <- regivBMA[["lambda.bar"]] # Posterior mean auxiliary
36 plot(EVAux[,1])
37 plot(EVAux[,2])
38 EVsigma <- regivBMA[["Sigma.bar"]] # Posterior mean variance
39 matrix

```

Bayesian model averaging has been also extended to state-space models. The point of departure is the univariate random walk state-space model (see Equations 8.1 and 8.2 in Chapter 8) conditional on model  $\mathcal{M}_m$ ,  $m =$

$1, 2, \dots, M$ .

$$y_t = \mathbf{x}_{mt}^\top \boldsymbol{\beta}_{mt} + \mu_{mt} \quad (10.1)$$

$$\boldsymbol{\beta}_{mt} = \boldsymbol{\beta}_{mt-1} + \mathbf{w}_{mt}, \quad (10.2)$$

where  $\mu_{mt} \sim N(0, \sigma^2)$  and  $\mathbf{w}_{mt} \sim N(\mathbf{0}, \boldsymbol{\Omega}_{mt})$ .

Given  $\boldsymbol{\beta}_{mt-1}|y_{1:t-1} \sim N(\mathbf{b}_{mt-1}, \mathbf{B}_{mt-1})$ , then, we know from Chapter 8 that  $\boldsymbol{\beta}_{mt}|y_{1:t-1} \sim N(\mathbf{b}_{mt-1}, \mathbf{R}_{mt})$ ,  $\mathbf{R}_{mt} = \mathbf{B}_{mt-1} + \boldsymbol{\Omega}_{mt}$ .

Specification of  $\boldsymbol{\Omega}_t$  can be highly demanding. Thus, a common approach is to express  $\boldsymbol{\Omega}_{mt} = \frac{1-\lambda}{\lambda} \mathbf{B}_{mt-1}$ , where  $\lambda$  is called the *forgetting parameter* or *discount factor*, because it discounts the matrix  $\mathbf{B}_{mt-1}$  that we would have with a deterministic state evolution into the matrix  $\mathbf{R}_{mt}$  [336, Chap. 4]. This parameter is typically slightly below 1, and implies that  $\mathbf{R}_{mt} = \lambda^{-1} \mathbf{B}_{mt-1}$ . ( $\lambda^{-1} > 1$ ).

[350] assume that the model changes infrequently, and its evolution is given by the transition matrix  $\mathbf{T} = [t_{ml}]$ , where  $t_{ml} = P(\mathcal{M}_t = \mathcal{M}_m | \mathcal{M}_{t-1} = \mathcal{M}_l)$ .

Then, the aim is to calculate the filtering distribution  $p(\boldsymbol{\beta}_{mt}, \mathcal{M}_t | y_t) = \sum_{m=1}^M p(\boldsymbol{\beta}_{mt} | \mathcal{M}_t = \mathcal{M}_m, y_t) p(\mathcal{M}_t = \mathcal{M}_m | y_t)$ . Thus, given the conditional distribution of the state at time  $t-1$ ,  $p(\boldsymbol{\beta}_{mt-1}, \mathcal{M}_{t-1} | y_{t-1}) = \sum_{m=1}^M p(\boldsymbol{\beta}_{mt-1} | \mathcal{M}_{t-1} = \mathcal{M}_m, y_{t-1}) p(\mathcal{M}_{t-1} = \mathcal{M}_m | y_{t-1})$ , where the conditional distribution of  $\boldsymbol{\beta}_{mt-1}$  is approximated by a Gaussian distribution,  $\boldsymbol{\beta}_{mt-1} | \mathcal{M}_{t-1} = \mathcal{M}_m, y_{t-1} \sim N(\mathbf{b}_{mt-1}, \mathbf{B}_{mt-1})$ , then the first step to get the one-step-ahead predictive distribution is getting the prediction of the model indicator,

$$\begin{aligned} p(\mathcal{M}_t = \mathcal{M}_l | y_{t-1}) &= \sum_{m=1}^M p(\mathcal{M}_{t-1} = \mathcal{M}_m | y_{t-1}) \times t_{lm} \\ &\approx \frac{p(\mathcal{M}_{t-1} = \mathcal{M}_l | y_{t-1})^\delta + c}{\sum_{m=1}^M p(\mathcal{M}_{t-1} = \mathcal{M}_m | y_{t-1})^\delta + c}, \end{aligned}$$

where the second equality is used to avoid dealing with the  $M^2$  elements of the transition matrix  $\mathbf{T}$  such that the forgetting parameter  $\delta$  is used, this parameter is slightly less than 1, and  $c = 0.001/M$  is introduced to handle a model probability being brought to computational zero by outliers (see [350]).

Then, we get the one-step-ahead predictive distribution of the state vector,  $\boldsymbol{\beta}_{mt} | \mathcal{M}_t = \mathcal{M}_m, y_{t-1} \sim N(\mathbf{b}_{mt-1}, \lambda^{-1} \mathbf{B}_{mt-1})$

Now, we consider the filtering stage, where the model filtering equation is

$$p(\mathcal{M}_t = \mathcal{M}_l | y_t) = \frac{p(\mathcal{M}_t = \mathcal{M}_l | y_{t-1}) p_l(y_t | y_{t-1})}{\sum_{m=1}^M p(\mathcal{M}_t = \mathcal{M}_m | y_{t-1}) p_m(y_t | y_{t-1})},$$

where  $p_m(y_t | y_{t-1})$  is the one-step-ahead predictive distribution of  $y_t | y_{t-1}$ , which is  $N(f_t, Q_t)$ , where  $f_t = \mathbf{x}_{mt}^\top \mathbf{b}_{t-1}$  and  $Q_t = \mathbf{x}_{mt}^\top \lambda^{-1} \mathbf{B}_{mt-1} \mathbf{x}_{mt} + \sigma^2$  (see Chapter 8).

The states filtering equation is  $\beta_{mt} | \mathcal{M}_t = \mathcal{M}_m, y_t \sim N(\mathbf{b}_{mt}, \mathbf{B}_{mt})$  where  $\mathbf{b}_{mt}$  and  $\mathbf{B}_{mt}$  are given in the Kalman filtering recursion of Chapter 8.

[350] initiate their algorithm assuming equal prior model probabilities, and  $\sigma^2$  is estimated using a recursive method of moments estimator.<sup>7</sup>

We implement dynamic Bayesian model averaging in our GUI using the function *dma* from the package *dma*. Algorithm A31 shows how to perform inference using our GUI.

---

**Algorithm A31** Dynamic Bayesian model averaging

---

- 1: Select *Bayesian Model Averaging* on the top panel
  - 2: Select *Normal data* model using the left radio button
  - 3: Select *Dynamic Bayesian model averaging* using the right radio button under **Which type do you want to perform?**
  - 4: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 5: Upload the matrix of models selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 6: Type the *forgetting parameters* in the boxes under **Lambda: Number slightly below 1** and **Delta: Number slightly below 1**. This is not necessary as by default there is 0.99 in both cases
  - 7: Click the *Go!* button
  - 8: Analyze results: After a few seconds or minutes, a figure showing the posterior model probabilities in each period appears
  - 9: Download posterior results using the *Download results for Dynamic Model Average*. There are four files: the dynamic Bayesian model averaging prediction, the dynamic Bayesian average filtering recursions for each state, and their standard deviation, and the PMP of each model.
- 

**Example: Dynamic Bayesian model averaging**

We perform a simulation exercise where there are 8 ( $2^3$ ) competing models originating from 3 regressors:  $x_{tk} \sim N(0.5, 0.8^2)$  for  $k = 2, 3, 4$ , with  $\beta_1 = 0.5$ . The sequence  $\beta_{2t}$  ranges from 1 to 2 in steps of  $1/T$ , and  $\beta_{3t}$  is given by:

$$\beta_{3t} = \begin{cases} -1, & 1 < t \leq 0.75T \\ 0, & 0.75T < t \leq T \end{cases}$$

and  $\beta_4 = 1.2$ .

Then, we have the model:

$$y_t = \beta_1 + \beta_{2t}x_{2t} + \beta_{3t}x_{3t} + \beta_{4t}x_{4t} + \mu_t,$$

where  $\mu_t \sim N(0, 1)$  for  $t = 1, 2, \dots, 500$ . This setting implies that during

---

<sup>7</sup>[356] extends this approach to Markov chain Monte Carlo model composition

the first 75% of the period, the model with all 3 regressors is the data-generating process, while after this, the model with regressors 2 and 4 is the data-generating process.

The following code shows the simulation exercise and the results of the dynamic Bayesian model averaging, setting  $\lambda = \delta = 0.99$ .

### R code. Simulation exercise: Dynamic Bayesian model averaging

```

1 rm(list = ls()); set.seed(010101)
2 T <- 500; K <- 3
3 X <- matrix(rnorm(T*K, mean = 0.5, sd = 0.8), T, K)
4 combs <- expand.grid(c(0,1), c(0,1), c(0,1))
5 B1 <- 0.5; B2t <- seq(1, 2, length.out=T)
6 a <- 0.75; B3t <- c(rep(-1, round(a*T)), rep(0, round((1-a)*T)))
7 B4 <- 1.2; sigma <- 1; mu <- rnorm(T, 0, sigma)
8 y <- B1 + X[,1]*B2t + X[,2]*B3t + X[,3]*B4 + mu
9 T0 <- 50
10 dma.test <- dma::dma(X, y, combs, lambda=.99, gamma=.99,
11   initialperiod = T0)
11 plot(dma.test[["pmp"]][,-c(1:T0),8], type = "l", col = "green",
12   main = "Posterior model probability: Model all
13   regressors vs model regressors 2 and 4", xlab = "Time",
14   ylab = "PMP")
12 lines(dma.test[["pmp"]][,-c(1:T0),6], col = "red")
13 legend(x = 0, y = 1, legend = c("Model: All regressors", "
14   Model: Regressors 2 and 4"), col = c("green", "red"),
15   lty=1:1, cex=0.8)
14 require(latex2exp)
15 plot(dma.test[["thetahat.ma"]][,-c(1:T0),1], type = "l", col
16   = "green", main = "Bayesian model average filtering
17   recursion", xlab = "Time", ylab = TeX("\beta_{1}"))
16 abline(h = B1, col = "red")
17 legend(x = 0, y = 0.4, legend = c("State filtering", "State
18   population"), col = c("green", "red"), lty=1:1, cex=0.8)
18 plot(dma.test[["thetahat.ma"]][,-c(1:T0),2], type = "l", col
19   = "green", main = "Bayesian model average filtering
20   recursion", xlab = "Time", ylab = TeX("\beta_{2t}"),
20   ylim = c(0.5,2))
19 lines(B2t[-c(1:T0)], col = "red")
20 legend(x = 0, y = 0.8, legend = c("State filtering", "State
21   population"), col = c("green", "red"), lty=1:1, cex=0.8)

```

**R code. Simulation exercise: Dynamic Bayesian model averaging**

```

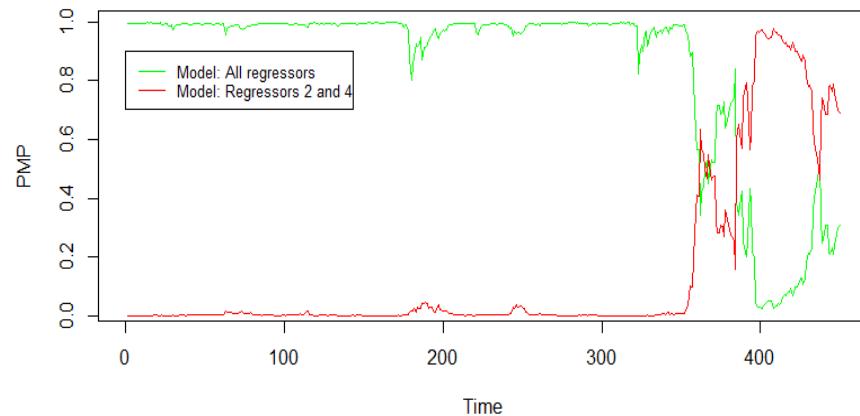
1 plot(dma.test[["thetahat.ma"]][-,3], type = "l", col
= "green", main = "Bayesian model average filtering
recursion", xlab = "Time", ylab = TeX("$\beta_{3t}$"))
2 lines(B3t[-c(1:T0)], col = "red")
3 legend(x = 0, y = -0.4, legend = c("State filtering", "State
population"), col = c("green", "red"), lty=1:1, cex
=0.8)
4 plot(dma.test[["thetahat.ma"]][-,4], type = "l", col
= "green", main = "Bayesian model average filtering
recursion", xlab = "Time", ylab = TeX("$\beta_{4t}$"))
5 abline(h = B4, col = "red")
6 legend(x = 0, y = 1.3, legend = c("State filtering", "State
population"), col = c("green", "red"), lty=1:1, cex=0.8)

```

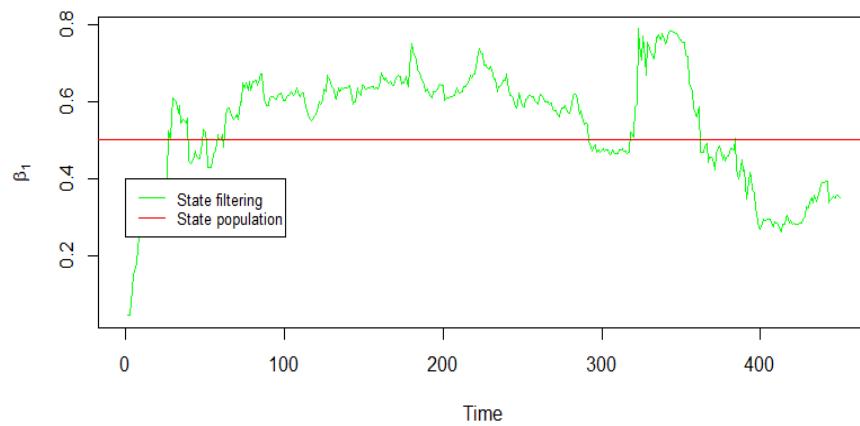
Figure 10.1 shows the posterior model probabilities for the model with all the regressors (green line) and the model with regressors 2 and 4 (red line). On one hand, we see that the model with all regressors, which is the data-generating process in the first period ( $t \leq 0.75T$ ), has a PMP close to 1, and then its PMP decreases. On the other hand, the model with regressors 2 and 4 has a PMP close to 0 in the first part of the period, and then its PMP increases to values higher than 60% on average, when this model becomes the data-generating process. These results suggest that, in this particular simulation exercise, the dynamic Bayesian model averaging method works relatively well in calculating the PMPs.

Figures 10.2, 10.3, 10.4, and 10.5 show a comparison between the Bayesian model averaging filtering recursions of the states (green lines) and their population values (red lines). We observe that the filtering recursions follow the general pattern of the population values. However, the values are not perfectly aligned. This discrepancy arises because the posterior model probabilities (PMPs) of the models that match the data-generating process are not equal to 1, which in turn affects the performance of the filtering recursions.

Dynamic Bayesian model averaging was extended to logit models by [302]. We ask in Exercise 12 to perform a simulation of this model, and perform BMA using the function *logistic.dma* from the *dma* package.

**Posterior model probability: Model all regressors vs model regressors 2 and 4****FIGURE 10.1**

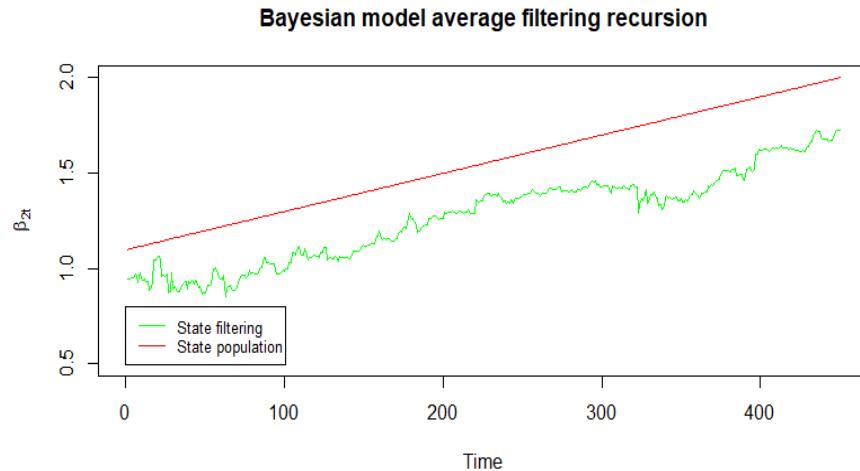
Posterior model probability: Dynamic Bayesian model averaging.

**Bayesian model average filtering recursion****FIGURE 10.2**

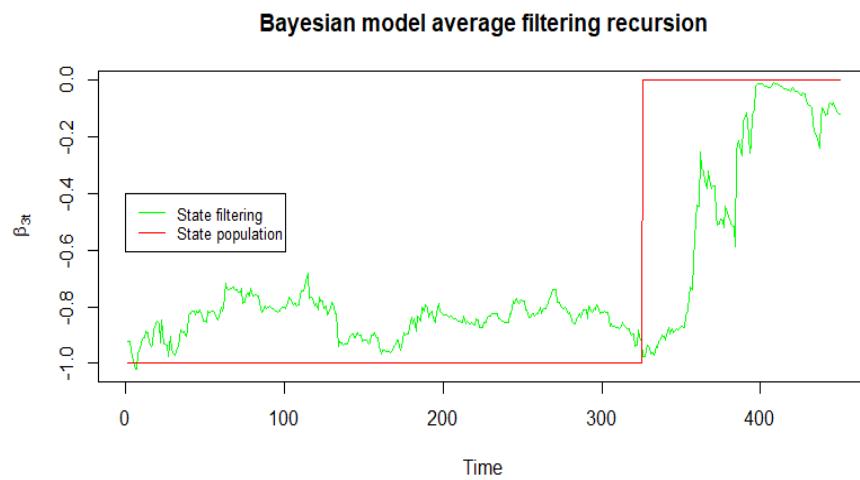
State  $\beta_1$ : Population versus dynamic Bayesian model averaging of the filtering recursion.

### 10.3 Generalized linear models

Generalized linear models (GLMs) were introduced by [320], extending the concept of linear regression to a more general setting. These models are char-

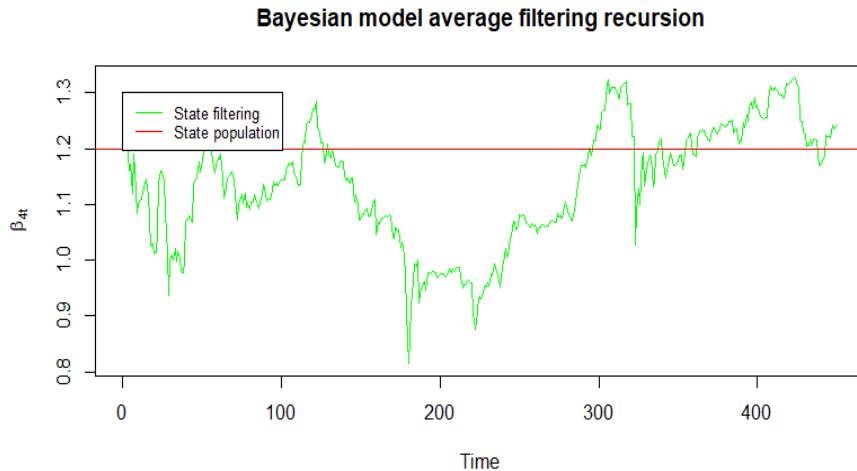
**FIGURE 10.3**

State  $\beta_{2t}$ : Population versus dynamic Bayesian model averaging of the filtering recursion.

**FIGURE 10.4**

State  $\beta_{3t}$ : Population versus dynamic Bayesian model averaging of the filtering recursion.

acterized by: i) a dependent variable  $y_i$  whose probability distribution func-

**FIGURE 10.5**

State  $\beta_4$ : Population versus dynamic Bayesian model averaging of the filtering recursion.

tion belongs to the exponential family (see Section 3.1), ii) a linear predictor  $\eta = \mathbf{x}^\top \boldsymbol{\beta}$ , and iii) a link function such that  $\mathbb{E}[y|\mathbf{x}] = g^{-1}(\mathbf{x}^\top \boldsymbol{\beta})$ , which implies that  $g(\mathbb{E}[y|\mathbf{x}]) = \mathbf{x}^\top \boldsymbol{\beta}$ . GLMs can be extended to the overdispersed exponential family [303].

As we know from Section 3.1, the Poisson distribution belongs to the exponential family, such that  $p(y|\lambda) = \frac{\exp(-\lambda) \exp(y \log(\lambda))}{y!}$ , or in the canonical form  $p(y|\eta) = \frac{\exp(\eta y - \exp(\eta))}{y!}$ , where  $\eta = \log(\lambda)$ , which means that  $\mathbf{x}^\top \boldsymbol{\beta} = \log(\lambda)$ . Consequently,  $\mathbb{E}[y|\mathbf{x}] = \nabla(\exp(\eta)) = \exp(\eta) = \lambda = \exp(\mathbf{x}^\top \boldsymbol{\beta})$ . Therefore, the link function in the Poisson case is the *log* function. In Exercise 6, we ask you to show that the link function in the Bernoulli case is the *logit* function. Other examples include the identity function in the case of the Gaussian distribution and the negative inverse in the case of the gamma distribution.

We can use the GLM framework to perform Bayesian model averaging (BMA) using the BIC approximation, following [347]. Specifically, the BIC is given by  $BIC = k_m \log(N) - 2 \log(p(\hat{\boldsymbol{\theta}}_m | \mathbf{y}))$ , where  $\hat{\boldsymbol{\theta}}_m$  is the maximum likelihood estimator. Thus, we simply need to calculate the likelihood function at the maximum likelihood estimator.

#### Example: Simulation exercises

Let's perform some simulation exercises to assess the performance of the BIC approximation using the Occam's window in GLMs. There are 27 re-

gressors, where  $x_{i1}$  and  $x_{i2}$  are just the relevant regressors in all exercises,  $i = 1, 2, \dots, 1000$ .

- Logit:  $x_k \sim N(0, 1)$ ,  $k = 1, \dots, 27$ , and  $p(y_i = 1 | \mathbf{x}_i) = \exp(0.5 + 0.8x_{i1} - 1.2x_{i2}) / (1 + \exp(0.5 + 0.8x_{i1} - 1.2x_{i2}))$ .
- Gamma:  $x_k \sim N(0, 0.5^2)$ ,  $k = 1, \dots, 27$ , and  $y_i \sim G(\alpha, \delta)$  where  $\alpha = -(0.5 + 0.2x_{i1} + 0.1x_{i2})^{-1}$  and  $\delta = 1$ .
- Poisson:  $x_k \sim N(0, 1)$ ,  $k = 1, \dots, 27$ , and  $\mathbb{E}[y_i | \mathbf{x}_i] = \lambda_i = \exp(0.5 + 1.1x_{i1} + 0.7x_{i2})$ .

Our GUI uses the command `bic.glm` from the *BMA* package to perform BMA using the BIC approximation with the Occam's window in GLMs. The Algorithm A32 shows how to do this in our GUI, and the following code shows how to perform BMA in logit models using the simulation setting.

---

**Algorithm A32** Bayesian model averaging in generalized linear models using the Bayesian information criterion

---

- 1: Select *Bayesian Model Averaging* on the top panel
  - 2: Select the generalized linear model using the left radio button. Options: *Binomial data (Logit)*, *Real positive data (Gamma)* and *Count data (Poisson)*
  - 3: Upload the dataset selecting first if there is header in the file, and the kind of separator in the *csv* file of the dataset (comma, semicolon, or tab). Then, use the *Browse* button under the **Choose File** legend
  - 4: Type the *OR* number of the Occam's window in the box under **OR: Number between 5 and 50**, this is not necessary as by default there is 50
  - 5: Type the *OL* number of the Occam's window in the box under **OL: Number between 0.0001 and 1**, this is not necessary as by default there is 0.0025
  - 6: Click the *Go!* button
  - 7: Analyze results: After a few seconds or minutes, a table appears showing, for each regressor in the dataset, the PIP (posterior inclusion probability, **p!=0**), the BMA posterior mean (**EV**), and the BMA standard deviation (**SD**). User can also see the PMP for the most relevant models
  - 8: Download posterior results using the *Download results using BIC*. There are two files, the first has the best models by row according to the PMP (last column) indicating with a 1 inclusion of the variable (0 indicates no inclusion), and the second file has the PIP, the BMA expected value and standard deviation for each variable in the dataset
- 

The results show that the PIPs of  $x_{i1}$  and  $x_{i2}$  are equal 1 in all three settings, the data generating process gets the highest PMP, and the BMA posterior means are close to the population values in each simulation setting.

The other variables get PIPs close to 0, except a few exceptions, and the BMA posterior means are also close to 0. This suggests that the BIC approximation does a good job finding the data generating process in generalized linear models.

We can take advantage of the *glm* function in **R** to perform BMA by programming an MC3 algorithm. The following code illustrates how to do this in the Poisson simulation. First, we simulate the data; second, we define a function to compute the log marginal likelihood approximation using the results from the *glm* function. Then, we initialize the models to begin the MC3 algorithm. After that, we implement the MC3 algorithm, which involves small modifications of the code used for MC3 in Gaussian linear models.<sup>8</sup> We can calculate the posterior model probabilities (PMPs), posterior inclusion probabilities (PIPs), BMA means, and standard deviations as we did previously.

The simulation setting involves  $2^{27}$  models, which corresponds to approximately 135 million models in the model space. We run our MC3 algorithm using the BIC approximation with 50,000 iterations. This takes considerably more time than the BIC approximation from the *BMA* package, but it seems to perform well in identifying the data-generating process, as the PMP of this model equals 1. The posterior inclusion probabilities (PIPs) for  $x_{i1}$  and  $x_{i2}$  are also 1, and the posterior means are 1.1 and 0.7, respectively, which are equal to the population values. The t-ratios are far greater than 2. However, running 50,000 iterations results in mass concentration in one model, in this case, the data-generating process. If we run 25,000 MC3 iterations, the highest PMP is 0.8, but it is not associated with the data-generating process. Nonetheless, the PIP is equal to 1 for  $x_{i1}$  and  $x_{i2}$ , and other regressors also have high PIPs. The BMA means for  $x_{i1}$  and  $x_{i2}$  are equal to the population values, and the BMA means for the other regressors are equal to 0. The t-ratios of the regressors in the population statistical model are much greater than 2, whereas the t-ratios of the other regressors are equal to 0. This exercise demonstrates that 25,000 iterations were not sufficient to uncover the data-generating process. However, it also emphasizes an important point: we need to analyze all the relevant results from the BMA analysis, not just the PMPs and/or PIPs.

In Exercise 10, we ask you to use this approach to perform a BMA algorithm in the logit regression, using the simulation setting for logit models from this section.

---

<sup>8</sup>Note that we have split the *while* loop in the code into two pages.

**R code. Simulation exercise: BMA for generalized linear models**

```

1 ##### Logit #####
2 rm(list = ls()); set.seed(010101)
3 n<-1000; B<-c(0.5,0.8,-1.2)
4 X<-matrix(cbind(rep(1,n),rnorm(n,0,1),rnorm(n,0,1)),n,length(B))
5 p <- exp(X%*%B)/(1+exp(X%*%B)); y <- rbinom(n, 1, p)
6 nXgar<-25; Xgar<-matrix(rnorm(nXgar*n),n,nXgar)
7 df<-as.data.frame(cbind(y,X[,-1],Xgar))
8 colnames(df) <- c("y", "x1", "x2", "x3", "x4", "x5", "x6", "x7",
9 "x8", "x9", "x10", "x11", "x12", "x13", "x14", "x15", "x16", "x17",
10 "x18", "x19", "x20", "x21", "x22", "x23", "x24", "x25", "x26", "x27")
11 BMAglmLogit <- BMA::bic.glm(y ~ x1+x2+x3+x4+x5+x6+x7+x8+x9+
12 x10+x11+x12+x13+x14+x15+x16+x17+x18+x19+x20+x21+x22+x23+
13 x24+x25+x26+x27, data = df, glm.family = binomial(link =
14 "logit"), strict = FALSE, OR = 50)
15 summary(BMAglmLogit)
16 ##### Gamma #####
17 rm(list = ls()); set.seed(010101)
18 n<-1000; B<- c(0.5, 0.2, 0.1)
19 X<-matrix(cbind(rep(1,n),rnorm(n,0,0.5),rnorm(n,0,0.5)),n,
20 length(B))
21 y1 <- (X%*%B)^(-1)
22 y <- rgamma(n,y1,scale=1)
23 nXgar<-25; Xgar<-matrix(rnorm(nXgar*n),n,nXgar)
24 df<-as.data.frame(cbind(y,X[,-1],Xgar))
25 colnames(df) <- c("y", "x1", "x2", "x3", "x4", "x5", "x6", "x7",
26 "x8", "x9", "x10", "x11", "x12", "x13", "x14", "x15",
27 "x16", "x17", "x18", "x19", "x20", "x21", "x22", "x23",
28 "x24", "x25", "x26", "x27")
29 BMAglmGamma <- BMA::bic.glm(y ~ x1+x2+x3+x4+x5+x6+x7+x8+x9+
30 x10+x11+x12+x13+x14+x15+x16+x17+x18+x19+x20+x21+x22+x23+
31 x24+x25+x26+x27, data = df, glm.family = Gamma(link =
32 "inverse"), strict = FALSE, OR = 50)
33 summary(BMAglmGamma)
34 ##### Poisson #####
35 rm(list = ls()); set.seed(010101)
36 n<-1000; B<-c(2,1.1,0.7)
37 X<-matrix(cbind(rep(1,n),rnorm(n,0,1),rnorm(n,0,1)),n,length(B))
38 y1<-exp(X%*%B); y<-rpois(n,y1)
39 nXgar<-25; Xgar<-matrix(rnorm(nXgar*n),n,nXgar)
40 df<-as.data.frame(cbind(y,X[,-1],Xgar))
41 colnames(df) <- c("y", "x1", "x2", "x3", "x4", "x5", "x6", "x7",
42 "x8", "x9", "x10", "x11", "x12", "x13", "x14", "x15",
43 "x16", "x17", "x18", "x19", "x20", "x21", "x22", "x23",
44 "x24", "x25", "x26", "x27")

```

**R code. Simulation exercise: BMA for generalized linear models**

```

1 BMAGlmPoisson <- BMA:::bic.glm(y ~ x1+x2+x3+x4+x5+x6+x7+x8+x9
+ x10+x11+x12+x13+x14+x15+x16+x17+x18+x19+x20+x21+x22+x23
+ x24+x25+x26+x27, data = df, glm.family = poisson(link =
  "log"), strict = FALSE, OR = 50)
2 summary(BMAGlmPoisson)

```

**R code. Simulation exercise: BMA for generalized linear models using MC3 from scratch**

```

1 rm(list = ls()); set.seed(010101)
2 n<-1000; B<-c(2,1.1,0.7)
3 X<-matrix(cbind(rep(1,n),rnorm(n,0,1),rnorm(n,0,1)),n,length
  (B))
4 y1<-exp(X%*%B); y<-rpois(n,y1)
5 nXgar<-25; Xgar<-matrix(rnorm(nXgar*n),n,nXgar)
6 df<-as.data.frame(cbind(y,X[,-1],Xgar))
7 colnames(df) <- c("y", "x1", "x2", "x3", "x4", "x5", "x6",
  "x7", "x8", "x9", "x10", "x11", "x12", "x13", "x14", "x15",
  "x16", "x17", "x18", "x19", "x20", "x21", "x22", "x23",
  "x24", "x25", "x26", "x27")
8 Xnew <- apply(df[,-1], 2, scale)
9 BICfunct <- function(Model){
10   indr <- Model == 1; kr <- sum(indr)
11   if(kr > 0){
12     Xr <- as.matrix(Xnew[, indr])
13     model <- glm(y ~ Xr, family = poisson(link = "log"))
14     model_bic <- BIC(model)
15     mllMod <- -model_bic/2
16   }else{
17     model <- glm(y ~ 1, family = poisson(link = "log"))
18     model_bic <- BIC(model); mllMod <- -model_bic/2
19   }
20   return(mllMod)
21 }
22 M <- 500; K <- dim(df)[2] - 1
23 Models <- matrix(rbinom(K*M, 1, p = 0.5), ncol = K, nrow = M
  )
24 mllnew <- sapply(1:M, function(s){BICfunct(matrix(Models[s,],
  1, K))})

```

**R code. Simulation exercise: BMA for generalized linear models using MC3 from scratch**

```

1 oind <- order(mllnew, decreasing = TRUE)
2 mllnew <- mllnew[oind]; Models <- Models[oind, ]
3 # Hyperparameters MC3
4 iter <- 25000
5 pb <- txtProgressBar(min = 0, max = iter, style = 3)
6 s <- 1
7 while(s <= iter){
8   ActModel <- Models[M,]
9   idK <- which(ActModel == 1)
10  Kact <- length(idK)
11  if(Kact < K & Kact > 1){
12    CardMol <- K
13    opt <- sample(1:3, 1)
14    if(opt == 1){ # Same
15      CandModel <- ActModel
16    }else{
17      if(opt == 2){ # Add
18        All <- 1:K
19        NewX <- sample(All[-idK], 1)
20        CandModel <- ActModel
21        CandModel[NewX] <- 1
22      }else{ # Subtract
23        LessX <- sample(idK, 1)
24        CandModel <- ActModel
25        CandModel[LessX] <- 0
26      }
27    }
28  }else{
29    CardMol <- K + 1
30    if(Kact == K){
31      opt <- sample(1:2, 1)
32      if(opt == 1){ # Same
33        CandModel <- ActModel
34      }else{ # Subtract
35        LessX <- sample(1:K, 1)
36        CandModel <- ActModel
37        CandModel[LessX] <- 0
38      }
39    }else{
40      if(K == 1){
41        opt <- sample(1:3, 1)
42        if(opt == 1){ # Same
43          CandModel <- ActModel

```

**R code. Simulation exercise: BMA for generalized linear models using MC3 from scratch**

---

```

1 }else{
2   if(opt == 2){ # Add
3     All <- 1:K
4     NewX <- sample(All[-idK], 1)
5     CandModel <- ActModel
6     CandModel[NewX] <- 1
7   }else{ # Subtract
8     LessX <- sample(idK, 1)
9     CandModel <- ActModel
10    CandModel[LessX] <- 0
11  }
12 }
13 }else{ # Add
14   NewX <- sample(1:K, 1)
15   CandModel <- ActModel
16   CandModel[NewX] <- 1
17 }
18 }
19 }
20 LogMLact <- BICfunct(matrix(ActModel, 1, K))
21 LogMLcand <- BICfunct(matrix(CandModel, 1, K))
22 alpha <- min(1, exp(LogMLcand-LogMLact))
23 u <- runif(1)
24 if(u <= alpha){
25   mllnew[M] <- LogMLcand
26   Models[M, ] <- CandModel
27   oind <- order(mllnew, decreasing = TRUE)
28   mllnew <- mllnew[oind]
29   Models <- Models[oind, ]
30 }else{
31   mllnew <- mllnew
32   Models <- Models
33 }
34 s <- s + 1
35 setTxtProgressBar(pb, s)
36 }
37 close(pb)
38 ModelsUni <- unique(Models)
39 mllnewUni <- sapply(1:dim(ModelsUni)[1], function(s){BICfunct
  (matrix(ModelsUni[s,], 1, K))})
40 StMarLik <- exp(mllnewUni-mllnewUni[1])
41 PMP <- StMarLik/sum(StMarLik) # PMP based on unique selected
  models

```

---

**R code. Simulation exercise: BMA for generalized linear models using MC3 from scratch**

```

1 plot(PMP)
2 ModelsUni[1,]
3 PIP <- NULL
4 for(k in 1:K){
5   PIPk <- sum(PMP[which(ModelsUni[,k] == 1)])
6   PIP <- c(PIP, PIPk)
7 }
8 plot(PIP)
9 Xnew <- df[,-1]
10 nModels <- dim(ModelsUni)[1]
11 Means <- matrix(0, nModels, K)
12 Vars <- matrix(0, nModels, K)
13 for(m in 1:nModels){
14   idXs <- which(ModelsUni[m,] == 1)
15   if(length(idXs) == 0){
16     Regm <- glm(y ~ 1, family = poisson(link = "log"))
17   }else{
18     Xm <- as.matrix(Xnew[, idXs])
19     Regm <- glm(y ~ Xm, family = poisson(link = "log"))
20     SumRegm <- summary(Regm)
21     Means[m, idXs] <- SumRegm[["coefficients"]][,-1]
22     Vars[m, idXs] <- SumRegm[["coefficients"]][,-1]^2
23   }
24 }
25 BMAmmeans <- colSums(Means*PMP)
26 BMAsd <- (colSums(PMP*Vars) + colSums(PMP*(Means-matrix(rep
27   (BMAmmeans, each = nModels), nModels, K))^2))^.5
28 plot(BMAmmeans)
29 plot(BMAsd)
30 plot(BMAmmeans/BMAsd)
```

## 10.4 Calculating the marginal likelihood

The BIC is an asymptotic approximation of the marginal likelihood, and consequently, it is used to obtain the Bayes factors. However, this method has limitations in applications with moderate and small sample sizes [157]. Therefore, other methods are available to calculate the Bayes factors when there is no analytical solution for the marginal likelihood.

Observe that calculating the Bayes factor with respect to a reference model ( $\mathcal{M}_0$ ) helps to obtain the posterior model probabilities,

$$\begin{aligned}\pi(\mathcal{M}_j|\mathbf{y}) &= \frac{p(\mathbf{y}|\mathcal{M}_j)\pi(\mathcal{M}_j)}{\sum_{m=1}^M p(\mathbf{y}|\mathcal{M}_m)\pi(\mathcal{M}_m)} \\ &= \frac{p(\mathbf{y}|\mathcal{M}_j)\pi(\mathcal{M}_j)/p(\mathbf{y}|\mathcal{M}_0)}{\sum_{m=1}^M p(\mathbf{y}|\mathcal{M}_m)\pi(\mathcal{M}_m)/p(\mathbf{y}|\mathcal{M}_0)} \\ &= \frac{BF_{j0} \times \pi(\mathcal{M}_j)}{\sum_{m=1}^M BF_{l0} \times \pi(\mathcal{M}_l)}.\end{aligned}$$

Thus,  $\pi(\mathcal{M}_j|\mathbf{y}) = \frac{BF_{j0}}{\sum_{m=1}^M BF_{l0}}$  assuming equal prior model probabilities.

In addition, it has been established in many settings that the Bayes factor is consistent. That is, the probability of identifying the true data generating process converges to 1 as the sample size increases to infinity. Alternatively, it asymptotically identifies the model that minimizes the Kullback-Leibler divergence with respect to the data generating process when this process is not part of the models under consideration [97, 445, 444].<sup>9</sup>

#### 10.4.1 Savage-Dickey density ratio

The Savage-Dickey density ratio is a way to calculate the Bayes factors when we compare nested models with particular priors [120, 441]. In particular, given the parameter space  $\boldsymbol{\theta} = (\boldsymbol{\omega}^\top, \boldsymbol{\psi}^\top)^\top \in \Theta = \Omega \times \Psi$ , where we wish to test the null hypothesis  $H_0 : \boldsymbol{\omega} = \boldsymbol{\omega}_0$  (model  $\mathcal{M}_1$ ) versus  $H_1 : \boldsymbol{\omega} \neq \boldsymbol{\omega}_0$  (model  $\mathcal{M}_2$ ), if  $\pi(\boldsymbol{\psi}|\boldsymbol{\omega}_0, \mathcal{M}_2) = \pi(\boldsymbol{\psi}|\mathcal{M}_1)$ ,<sup>10</sup> then the Bayes factor comparing  $\mathcal{M}_1$  versus  $\mathcal{M}_2$  is

$$BF_{12} = \frac{\pi(\boldsymbol{\omega} = \boldsymbol{\omega}_0|\mathbf{y}, \mathcal{M}_2)}{\pi(\boldsymbol{\omega} = \boldsymbol{\omega}_0|\mathcal{M}_2)}, \quad (10.3)$$

where  $\pi(\boldsymbol{\omega} = \boldsymbol{\omega}_0|\mathbf{y}, \mathcal{M}_2)$  and  $\pi(\boldsymbol{\omega} = \boldsymbol{\omega}_0|\mathcal{M}_2)$  are the posterior and prior densities of  $\boldsymbol{\omega}$  under  $\mathcal{M}_2$  evaluated at  $\boldsymbol{\omega}_0$  (see [441]).

Equation 10.3 is called the Savage-Dickey density ratio. A nice feature is that just requires estimation of model  $\mathcal{M}_2$ , and evaluation of the prior and posterior densities. This means no evaluation of the marginal likelihood [250, Chap. 4].

---

<sup>9</sup>[227] highlight the important distinction between pairwise consistency and model selection consistency. The latter requires the consistency of a sequence of pairwise nested comparisons.

<sup>10</sup>Note that a sufficient condition for this assumption is to assume the same prior for the parameters that are the same in each model. [441] incorporate a correction factor when this assumption is not satisfied.

### 10.4.2 Chib's methods

Another popular method to calculate the marginal likelihood is given by [89] and [96]. The former is an algorithm to calculate the marginal likelihood from the posterior draws of the Gibbs sampling algorithm, and the latter calculates the marginal likelihood from the posterior draws of the Metropolis-Hastings algorithm.

The point of departure in [89] is the identity

$$\pi(\boldsymbol{\theta}^*|\mathbf{y}, \mathcal{M}_m) = \frac{p(\mathbf{y}|\boldsymbol{\theta}^*, \mathcal{M}_m) \times \pi(\boldsymbol{\theta}^*|\mathcal{M}_m)}{p(\mathbf{y}|\mathcal{M}_m)},$$

where  $\boldsymbol{\theta}^*$  is a particular value of  $\boldsymbol{\theta}$  of high probability, for instance, the mode. This implies that

$$p(\mathbf{y}|\mathcal{M}_m) = \frac{p(\mathbf{y}|\boldsymbol{\theta}^*, \mathcal{M}_m) \times \pi(\boldsymbol{\theta}^*|\mathcal{M}_m)}{\pi(\boldsymbol{\theta}^*|\mathbf{y}, \mathcal{M}_m)}.$$

We can easily calculate the numerator of this expression. However, the critical point in this expression is to calculate the denominator as we know  $\pi(\boldsymbol{\theta}^*|\mathbf{y}, \mathcal{M}_m)$  up to a normalizing constant. We can calculate this from the posterior draws. Assume that  $\boldsymbol{\theta} = [\boldsymbol{\theta}_1^\top \ \boldsymbol{\theta}_2^\top]^\top$ , then  $\pi(\boldsymbol{\theta}^*|\mathbf{y}, \mathcal{M}_m) = \pi(\boldsymbol{\theta}_1^*|\boldsymbol{\theta}_2^*, \mathbf{y}, \mathcal{M}_m) \times \pi(\boldsymbol{\theta}_2^*|\mathbf{y}, \mathcal{M}_m)$ . We have the first term because in the Gibbs sampling algorithm the posterior conditional distributions are available. The second is

$$\begin{aligned}\pi(\boldsymbol{\theta}_2^*|\mathbf{y}, \mathcal{M}_m) &= \int_{\boldsymbol{\Theta}_1} \pi(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2^*|\mathbf{y}, \mathcal{M}_m) d\boldsymbol{\theta}_1 \\ &= \int_{\boldsymbol{\Theta}_1} \pi(\boldsymbol{\theta}_2^*|\boldsymbol{\theta}_1, \mathbf{y}, \mathcal{M}_m) \pi(\boldsymbol{\theta}_1|\mathbf{y}, \mathcal{M}_m) d\boldsymbol{\theta}_1 \\ &\approx \frac{1}{S} \sum_{s=1}^S \pi(\boldsymbol{\theta}_2^*|\boldsymbol{\theta}_1^{(s)}, \mathbf{y}, \mathcal{M}_m),\end{aligned}$$

where  $\boldsymbol{\theta}_1^{(s)}$  are the posterior draws of  $\boldsymbol{\theta}_1$  from the Gibbs sampling algorithm.

The generalization to more blocks can be seen in [89] and [184, Chap. 7]. In addition, the extension to the Metropolis-Hastings algorithm can be seen in [96], and [184, Chap. 7].

### 10.4.3 Gelfand-Dey method

We can use the Gelfand-Dey method [157] when we want to calculate the Bayes factor to compare non-nested models, models where the Savage-Dickey density ratio is hard to calculate, or the Chib's methods are difficult to implement. The Gelfand-Dey method is very general, and can be used in virtually any model [250, Chap. 5].

Given a probability density function  $q(\boldsymbol{\theta})$ , whose support is in  $\Theta$ , then

$$\mathbb{E} \left[ \frac{q(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta}|\mathcal{M}_m)p(\mathbf{y}|\boldsymbol{\theta}_m, \mathcal{M}_m)} \middle| \mathbf{y}, \mathcal{M}_m \right] = \frac{1}{p(\mathbf{y}|\mathcal{M}_m)},$$

where the expected value is with respect to the posterior distribution given the model  $\mathcal{M}_m$  (see Exercise 12).

The critical point is to select a good  $q(\boldsymbol{\theta})$ . [171] recommends to use  $q(\boldsymbol{\theta})$  equal to a truncated multivariate normal density function with mean and variance equal to the posterior mean ( $\hat{\boldsymbol{\theta}}$ ) and variance ( $\hat{\Sigma}$ ) of  $\boldsymbol{\theta}$ . The truncation region is  $\hat{\Theta} = \left\{ \boldsymbol{\theta} : (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \hat{\Sigma}^{-1} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \leq \chi^2_{1-\alpha}(K) \right\}$ , where  $\chi^2_{1-\alpha}(K)$  is the  $(1 - \alpha)$  percentile of the Chi-squared distribution with  $K$  degrees of freedom,  $K$  is the dimension of  $\boldsymbol{\theta}$ . We can pick small values of  $\alpha$ , for instance,  $\alpha = 0.01$ .

Observe that

$$\mathbb{E} \left[ \frac{q(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta}|\mathcal{M}_m)p(\mathbf{y}|\boldsymbol{\theta}_m, \mathcal{M}_m)} \middle| \mathbf{y}, \mathcal{M}_m \right] \approx \frac{1}{S} \sum_{s=1}^S \left[ \frac{q(\boldsymbol{\theta}^{(s)})}{\pi(\boldsymbol{\theta}^{(s)}|\mathcal{M}_m)p(\mathbf{y}|\boldsymbol{\theta}_m^{(s)}, \mathcal{M}_m)} \right],$$

where  $\boldsymbol{\theta}_m^{(s)}$  are draws from the posterior distribution.

Observe that we can calculate the marginal likelihoods of the models in Chapters 6, 7, 8 and 9 using the Chib's methods and the Gelfand-Dickey method.

### Example: Simulation exercise

Let's check the performance of the Savage-Dickey density ratio, Chib's method and the Gelfand-Dey method to calculate the Bayes factor in a setting where we can obtain the analytical solution for the marginal likelihood. In particular, we will consider the Gaussian linear model with a conjugate prior (see Section 3.3).

Assume that the data generating process is given by

$$y_i = 0.7 + 0.3x_{i1} + 0.7x_{i2} - 0.2x_{i3} + 0.2x_{i4} + \mu_i,$$

where  $x_{i1} \sim B(0.3)$ ,  $x_{ik} \sim N(0, 1)$ , for  $k = 2, \dots, 4$ , and  $\mu_i \sim N(0, 1)$ , for  $i = 1, 2, \dots, 500$ . Let us set  $H_0 : \beta_5 = 0$  (model  $\mathcal{M}_1$ ) versus  $H_1 : \beta_5 \neq 0$  (model  $\mathcal{M}_2$ ), that is,  $x_{i4}$  is not relevant under  $H_0$ .

We assume that  $\boldsymbol{\beta}_{m0} = \mathbf{0}_{m0}$ ,  $\mathbf{B}_{m0} = 0.5\mathbf{I}_m$ ,  $\alpha_0 = \delta_0 = 4$ . The dimensions of  $\mathbf{0}_{m0}$  and  $\mathbf{I}_m$  are 4 for model  $\mathcal{M}_1$  and 5 for model  $\mathcal{M}_2$ . In addition, we assume equal prior probabilities for both models.

We know from Section 3.3 that the marginal likelihood is

$$p(\mathbf{y}|\mathcal{M}_m) = \frac{\delta_{m0}^{\alpha_{m0}/2}}{\delta_{mn}^{\alpha_{mn}/2}} \frac{|\mathbf{B}_{mn}|^{1/2} \Gamma(\alpha_{mn}/2)}{|\mathbf{B}_{m0}|^{1/2} \Gamma(\alpha_{m0}/2)},$$

where  $\mathbf{B}_{mn} = (\mathbf{B}_{m0}^{-1} + \mathbf{X}_m^\top \mathbf{X}_m)^{-1}$ ,  $\boldsymbol{\beta}_{mn} = \mathbf{B}_{mn}(\mathbf{B}_{m0}^{-1} \boldsymbol{\beta}_{m0} + \mathbf{X}_m^\top \mathbf{X}_m \hat{\boldsymbol{\beta}}_m)$ ,

$\alpha_{mn} = \alpha_{m0} + N$ ,  $\delta_{mn} = \delta_{m0} + (\mathbf{y} - \mathbf{X}_m \hat{\boldsymbol{\beta}}_m)^\top (\mathbf{y} - \mathbf{X}_m \hat{\boldsymbol{\beta}}_m) + (\hat{\boldsymbol{\beta}}_m - \boldsymbol{\beta}_{m0})^\top ((\mathbf{X}_m^\top \mathbf{X}_m)^{-1} + \mathbf{B}_{m0})^{-1} (\hat{\boldsymbol{\beta}}_m - \boldsymbol{\beta}_{m0})$ , and  $\hat{\boldsymbol{\beta}}_m$  is the ordinary least squares estimate,  $m = 1, 2$  are the indices of the models.

The log marginal likelihoods for models  $\mathcal{M}_1$  and  $\mathcal{M}_2$  are -751.72 and -740.79, respectively. This implies a  $2 \times \log(BF_{21}) = 21.85$  which means positive evidence against model  $\mathcal{M}_1$  (see Table 1.1).

We have different ways to calculate the Bayes factor using the Savage-Dickey density ratio in this example because we know that the marginal prior and marginal posterior distributions of  $\beta_5$  have analytical solutions. In addition, we can use the posterior draws of  $\sigma^2$  to evaluate the conditional prior and conditional posterior distributions at  $\beta_5 = 0$ . We show in the following code the latter approach, as it is more general than using analytical solutions, which are not always available.

We know that the conditional posterior distribution of  $\beta_5$  is  $N(\beta_{5n}, \sigma \mathbf{B}_{55n})$ , where  $\beta_{5n}$  is the 5th element of  $\boldsymbol{\beta}_n$ , and  $\mathbf{B}_{55n}$  is the element 5,5 of  $\mathbf{B}_n$ . Then,

$$\begin{aligned}\pi(\beta_5 = 0 | \mathbf{y}, \mathcal{M}_2) &= \int_{\mathcal{R}^+} \pi(\beta_5 = 0 | \mathbf{y}, \sigma^2) \pi(\sigma^2 | \mathbf{y}) d\sigma^2 \\ &\approx \frac{1}{S} \sum_{s=1}^S \pi(\beta_5 = 0 | \mathbf{y}, \sigma^{2(s)}),\end{aligned}$$

where  $\sigma^{2(s)}$  are draws from the posterior distribution of  $\sigma^2$ .

We can follow the same logic to obtain an approximation to  $\pi(\beta_5 = 0 | \mathcal{M}_2)$  by sampling draws from the prior distribution of  $\sigma^2$ .

We obtain  $2 \times \log(BF_{21}) = 21.85$  using the Savage-Dickey density ratio, which is the same value as the analytic solution using the marginal likelihoods.

We calculate the log marginal likelihood using the Chib's method taking into account that

$$\log(p(\mathbf{y} | \mathcal{M}_m)) = \log(p(\mathbf{y} | \boldsymbol{\theta}^*, \mathcal{M}_m)) + \log(\pi(\boldsymbol{\theta}^* | \mathcal{M}_m)) - \log(\pi(\boldsymbol{\theta}^* | \mathbf{y}, \mathcal{M}_m)),$$

where  $p(\mathbf{y} | \boldsymbol{\theta}^*, \mathcal{M}_m)$  is the value of a normal density with mean  $\mathbf{X}_m \boldsymbol{\beta}_m^*$  and variance  $\sigma_m^{2*} \mathbf{I}_N$  evaluated at  $\mathbf{y}$ . In addition,  $\log(\pi(\boldsymbol{\theta}^* | \mathcal{M}_m)) = \log(\pi(\boldsymbol{\beta}_m^* | \sigma_m^{2*})) + \log(\pi(\sigma_m^{2*}))$ , where the first term is the density of a normal with mean  $\boldsymbol{\beta}_{m0}$  and variance matrix  $\sigma^{2*} \mathbf{B}_{m0}$  evaluated at  $\boldsymbol{\beta}_m^*$ , and the second term is the density of an inverse-gamma with parameters  $\alpha_{m0}/2$  and  $\delta_{m0}/2$  evaluated at  $\sigma_m^{2*}$ . Finally, the third term in the right hand of the previous expression is  $\log(\pi(\boldsymbol{\theta}^* | \mathbf{y}, \mathcal{M}_m)) = \log(\pi(\boldsymbol{\beta}_m^* | \sigma_m^{2*}, \mathbf{y})) + \log(\pi(\sigma_m^{2*} | \mathbf{y}))$ , where the first term is the density of a normal with mean  $\boldsymbol{\beta}_{mn}$  and variance matrix  $\sigma_m^{2*} \mathbf{B}_{mn}$  evaluated at  $\boldsymbol{\beta}_m^*$ , and the second term is the density of an inverse-gamma with parameters  $\alpha_{mn}/2$  and  $\delta_{mn}/2$  evaluated at  $\sigma_m^{2*}$ . We use the modes of the posterior draws of  $\boldsymbol{\beta}_m$  and  $\sigma_m^2$  as reference values.

We get the same value, up to two decimals, for the log marginal likelihood of the restricted and unrestricted models using the Chib's method and the

analytical expression. Thus,  $2 \times \log(BF_{21}) = 21.85$ , that is, positive evidence against model  $\mathcal{M}_1$  (see Table 1.1).

We calculate the log marginal likelihood using the Gelfand-Dey method taking into account that

$$\log \left[ \frac{q(\boldsymbol{\theta}^{(s)})}{\pi(\boldsymbol{\theta}^{(s)}|\mathcal{M}_m)p(\mathbf{y}|\boldsymbol{\theta}_m^{(s)}, \mathcal{M}_m)} \right] = \log(q(\boldsymbol{\theta}^{(s)})) - \log(\pi(\boldsymbol{\theta}^{(s)}|\mathcal{M}_m)) \\ - \log(p(\mathbf{y}|\boldsymbol{\theta}_m^{(s)}, \mathcal{M}_m)),$$

where  $q(\boldsymbol{\theta}^{(s)})$  is the truncated multivariate normal density of Subsection 10.4.3 evaluated at  $\boldsymbol{\theta}^{(s)} = [\boldsymbol{\beta}^{(s)\top} \ \sigma^{2(s)}]^\top$ , which is the  $s$ -th posterior draw of the Gibbs sampling algorithm, such that  $\boldsymbol{\theta}^{(s)}$  satisfies the truncation restriction  $\log(\pi(\boldsymbol{\theta}^{(s)}|\mathcal{M}_m)) = \log(\pi(\boldsymbol{\beta}_m^{(s)}|\sigma_m^{2(s)})) + \log(\pi(\sigma_m^{2(s)}))$ , where the first term is the density of a normal with mean  $\boldsymbol{\beta}_{m0}$  and variance matrix  $\sigma^{2(s)}\mathbf{B}_{m0}$  evaluated at  $\boldsymbol{\beta}_m^{(s)}$ , and the second term is the density of an inverse-gamma with parameters  $\alpha_{m0}/2$  and  $\delta_{m0}/2$  evaluated at  $\sigma_m^{2(s)}$ . The third term  $p(\mathbf{y}|\boldsymbol{\theta}^{(s)}, \mathcal{M}_m)$  is the value of a normal density with mean  $\mathbf{X}_m\boldsymbol{\beta}_m^{(s)}$  and variance  $\sigma_m^{2(s)}\mathbf{I}_N$  evaluated at  $\mathbf{y}$ .

The log marginal likelihoods of the restricted and unrestricted models using the Gelfand-Dey method are -751.79 and -740.89, respectively. This implies  $2 \times \log(BF_{21}) = 21.81$ , which is positive evidence in favor of the unrestricted model.

We see in this example that these methods give very good approximations to the true marginal likelihoods. However, the Savage-Dickey density ratio and Chib's method performed slightly better than the Gelfand-Dey method. In addition, the computational demand of the Gelfand-Dey method is by far the largest. This is because the Gelfand-Dey method requires many evaluations based on the posterior draws. However, we should keep in mind that the Gelfand-Dey method is more general.

The following code shows how to do all these calculations.

### R code. Simulation exercise: Bayes factors

```

1 rm(list = ls()); set.seed(010101)
2 N <- 500; K <- 5; K2 <- 3
3 B <- c(0.7, 0.3, 0.7, -0.2, 0.2)
4 X1 <- rbinom(N, 1, 0.3)
5 X2 <- matrix(rnorm(K2*N), N, K2)
6 X <- cbind(1, X1, X2)
7 Y <- X%*%B + rnorm(N, 0, sd = 1)

```

**R code. Simulation exercise: Bayes factors**

```

1 # Hyperparameters
2 d0 <- 4; a0 <- 4
3 b0 <- rep(0, K); cOpt <- 0.5
4 an <- N + a0; B0 <- cOpt*diag(K)
5 Bn <- solve(solve(B0)+t(X)%%X); bhat <- solve(t(X)%%X)%%t
   (X)%%Y
6 bn <- Bn%*%(solve(B0)%*%b0+t(X)%%X%%bhat)
7 dn <- as.numeric(d0 + t(Y-X%%bhat)%*%(Y-X%%bhat)+t(bhat -
   b0)%*%solve(solve(t(X)%%X)+B0)%*%(bhat - b0))
8 Hn <- as.matrix(Matrix::forceSymmetric(dn*Bn/an))
9 S <- 10000
10 LogMarLikLM <- function(X, c0){
11   K <- dim(X)[2]
12   N <- dim(X)[1]
13   # Hyperparameters
14   B0 <- c0*diag(K)
15   b0 <- rep(0, K)
16   # Posterior parameters
17   bhat <- solve(t(X)%%X)%%t(X)%%Y
18   # Force this matrix to be symmetric
19   Bn <- as.matrix(Matrix::forceSymmetric(solve(solve(B0) + t
     (X)%%X)))
20   bn <- Bn%*%(solve(B0)%*%b0 + t(X)%%X%%bhat)
21   dn <- as.numeric(d0 + t(Y)%%Y+t(b0)%*%solve(B0)%*%b0-t(bn
     )%*%solve(Bn)%*%bn)
22   an <- a0 + N
23   # Log marginal likelihood
24   logpy <- (N/2)*log(1/pi)+(a0/2)*log(d0)-(an/2)*log(dn) +
     0.5*log(det(Bn)/det(B0)) + lgamma(an/2)-lgamma(a0/2)
25   return(-logpy)
26 }
27 LogMarM2 <- -LogMarLikLM(X = X, c0 = cOpt)
28 LogMarM1 <- -LogMarLikLM(X = X[,1:4], c0 = cOpt)
29 BF12 <- exp(LogMarM1-LogMarM2)
30 BF12; 1/BF12
31 2*log(1/BF12)
32 # Savage-Dickey density ratio
33 # Posterior evaluation
34 Brest <- 0
35 sig2P <- invgamma::rinvgamma(S, shape = an/2, rate = dn/2)
36 PostRestCom <- mean(sapply(sig2P, function(x){dnorm(Brest,
   mean = bn[5], sd = (x*Bn[5,5])^0.5, log = FALSE)}))
37 # Prior evaluation
38 sig2 <- invgamma::rinvgamma(S, shape = a0/2, rate = d0/2)
39 PriorRestCom <- mean(sapply(sig2, function(x){dnorm(Brest,
   mean = 0, sd = (x*cOpt)^0.5, log = FALSE)}))
40 # Bayes factor
41 BF12SD <- PostRestCom/PriorRestCom
42 2*log(1/BF12SD)

```

### R code. Simulation exercise: Bayes factors

```

1 # Chib's method
2 sig2Post <- MCMCpack::rinvgamma(S, an/2, dn/2)
3 BetasGibbs <- sapply(1:S, function(s){MASS::mvrnorm(n = 1,
   mu = bn, Sigma = sig2Post[s]*Bn)})
4 # Mode function for continuous data
5 mode_continuous <- function(x){
6   density_est <- density(x)
7   mode_value <- density_est$x[which.max(density_est$y)]
8   return(mode_value)
9 }
10 # Unrestricted model
11 BetasMode <- apply(BetasGibbs, 1, mode_continuous)
12 Sigma2Mode <- mode_continuous(sig2Post)
13 VarModel <- Sigma2Mode*diag(N)
14 MeanModel <- X%*%BetasMode
15 LogLik <- mvtnorm::dmvnorm(c(Y), mean = MeanModel, sigma =
   VarModel, log = TRUE, checkSymmetry = TRUE)
16 LogPrior <- mvtnorm::dmvnorm(BetasMode, mean = rep(0, K),
   sigma = Sigma2Mode*cOpt*diag(K), log = TRUE,
   checkSymmetry = TRUE)+log(MCMCpack::dinvgamma(Sigma2Mode
   , a0/2, d0/2))
17 LogPost1 <- mvtnorm::dmvnorm(BetasMode, mean = bn, sigma =
   Sigma2Mode*Bn, log = TRUE, checkSymmetry = TRUE)
18 LogPost2 <- log(MCMCpack::dinvgamma(Sigma2Mode, an/2, dn/2))
19 LogMarLikChib <- LogLik + LogPrior -(LogPost1 + LogPost2)
20 # Restricted model
21 anRest <- N + a0; XRest <- X[, -5]
22 KRest <- dim(XRest)[2]; B0Rest <- cOpt*diag(KRest)
23 BnRest <- solve(solve(B0Rest)+t(XRest)%*%XRest)
24 bhatRest <- solve(t(XRest)%*%XRest)%*%t(XRest)%*%Y
25 b0Rest <- rep(0, KRest)
26 bnRest <- BnRest%*%(solve(B0Rest)%*%b0Rest+t(XRest)%*%XRest%
   *%bhatRest)
27 dnRest <- as.numeric(d0 + t(Y-XRest)%*%bhatRest)%*%(Y-XRest)%*
   %bhatRest)+t(bhatRest - b0Rest)%*%solve(solve(t(XRest)%*
   %XRest)+B0Rest)%*%(bhatRest - b0Rest)
28 sig2PostRest <- MCMCpack::rinvgamma(S, anRest/2, dnRest/2)
29 BetasGibbsRest <- sapply(1:S, function(s){MASS::mvrnorm(n =
   1, mu = bnRest, Sigma = sig2PostRest[s]*BnRest)})
30 BetasModeRest <- apply(BetasGibbsRest, 1, mode_continuous)
31 Sigma2ModeRest <- mode_continuous(sig2PostRest)
32 VarModelRest <- Sigma2ModeRest*diag(N)
33 MeanModelRest <- XRest%*%BetasModeRest
34 LogLikRest <- mvtnorm::dmvnorm(c(Y), mean = MeanModelRest,
   sigma = VarModelRest, log = TRUE, checkSymmetry = TRUE)

```

### R code. Simulation exercise: Bayes factors

```

1 LogPriorRest <- mvtnorm::dmvnorm(BetasModeRest, mean = rep
2   (0, KRest), sigma = Sigma2ModeRest*cOpt*diag(KRest), log
3   = TRUE, checkSymmetry = TRUE)+log(MCMCpack::dinvgamma(
4   Sigma2ModeRest, a0/2, d0/2))
5 LogPost1Rest <- mvtnorm::dmvnorm(BetasModeRest, mean =
6   bnRest, sigma = Sigma2ModeRest*BnRest, log = TRUE,
7   checkSymmetry = TRUE)
8 LogPost2Rest <- log(MCMCpack::dinvgamma(Sigma2ModeRest,
9   anRest/2, dnRest/2))
10 LogMarLikChibRest <- LogLikRest + LogPriorRest -(
11   LogPost1Rest + LogPost2Rest)
12 BFChibs <- exp(LogMarLikChibRest-LogMarLikChib)
13 BFChibs; 1/BFChibs; 2*log(1/BFChibs)
14 # Gelfand-Dey method
15 GDmarglik <- function(ids, X, Betas, MeanThetas, VarThetas,
16   sig2Post){
17   K <- dim(X)[2]; Thetas <- c(Betas[ids,], sig2Post[ids])
18   Lognom <- (1/(1-alpha))*mvtnorm::dmvnorm(Thetas, mean =
19     MeanThetas, sigma = VarThetas, log = TRUE, checkSymmetry
20     = TRUE)
21   Logden1 <- mvtnorm::dmvnorm(Betas[ids,], mean = rep(0, K),
22     sigma = sig2Post[ids]*cOpt*diag(K), log = TRUE,
23     checkSymmetry = TRUE) + log(MCMCpack::dinvgamma(sig2Post
24     [ids], a0/2, d0/2))
25   VarModel <- sig2Post[ids]*diag(N)
26   MeanModel <- X%*%Betas[ids,]
27   Logden2 <- mvtnorm::dmvnorm(c(Y), mean = MeanModel, sigma
28     = VarModel, log = TRUE, checkSymmetry = TRUE)
29   LogGDid <- Lognom - Logden1 - Logden2
30   return(LogGDid)
31 }
32 sig2Post <- MCMCpack::rinvgamma(S, an/2, dn/2)
33 Betas <- LaplacesDemon::rmvt(S, bn, Hn, an)
34 Thetas <- cbind(Betas, sig2Post)
35 MeanThetas <- colMeans(Thetas); VarThetas <- var(Thetas)
36 iVarThetas <- solve(VarThetas)
37 ChiSQ <- sapply(1:S, function(s){(Thetas[s,]-MeanThetas)%*%
38   iVarThetas%*%(Thetas[s,]-MeanThetas)})
39 alpha <- 0.01; criticalval <- qchisq(1-alpha, K + 1)
40 idGoodThetas <- which(ChiSQ <= criticalval)
41 pb <- txtProgressBar(min = 0, max = S, style = 3)
42 InvMargLik2 <- NULL
43 for(s in idGoodThetas){
44   LogInvs <- GDmarglik(ids = s, X = X, Betas = Betas,
45     MeanThetas = MeanThetas, VarThetas = VarThetas, sig2Post
46     = sig2Post)
47   InvMargLik2 <- c(InvMargLik2, LogInvs)
48   setTxtProgressBar(pb, s)
49 }
50 close(pb); mean(InvMargLik2)

```

### R code. Simulation exercise: Bayes factors

```

1 # Restricted model
2 anRest <- N + a0; XRest <- X[, -5]
3 KRest <- dim(XRest)[2]; B0Rest <- c0pt * diag(KRest)
4 BnRest <- solve(solve(B0Rest) + t(XRest) %*% XRest)
5 bhatRest <- solve(t(XRest) %*% XRest) %*% t(XRest) %*% Y
6 b0Rest <- rep(0, KRest)
7 bnRest <- BnRest %*% (solve(B0Rest) %*% b0Rest + t(XRest) %*% XRest %
8 *% bhatRest)
8 dnRest <- as.numeric(d0 + t(Y - XRest) %*% bhatRest) %*% (Y - XRest) %
9 %*% bhatRest + t(bhatRest - b0Rest) %*% solve(solve(t(XRest) %
10 %*% XRest) + B0Rest) %*% (bhatRest - b0Rest))
11 HnRest <- as.matrix(Matrix:::forceSymmetric(dnRest * BnRest /
12 anRest))
13 sig2PostRest <- MCMCpack::rinvgamma(S, anRest / 2, dnRest / 2)
14 BetasRest <- LaplacesDemon::rmvt(S, bnRest, HnRest, anRest)
15 ThetasRest <- cbind(BetasRest, sig2PostRest)
16 MeanThetasRest <- colMeans(ThetasRest)
17 VarThetasRest <- var(ThetasRest)
18 iVarThetasRest <- solve(VarThetasRest)
19 ChiSQRest <- sapply(1:S, function(s){(ThetasRest[s,] -
20 MeanThetasRest) %*% iVarThetasRest %*% (ThetasRest[s,] -
21 MeanThetasRest)})
22 idGoodThetasRest <- which(ChiSQRest <= criticalval)
23 pb <- txtProgressBar(min = 0, max = S, style = 3)
24 InvMargLik1 <- NULL
25 for(s in idGoodThetasRest){
26   LogInvs <- GDmarglik(ids = s, X = XRest, Betas = BetasRest
27 , MeanThetas = MeanThetasRest, VarThetas = VarThetasRest
28 , sig2Post = sig2PostRest)
29   InvMargLik1 <- c(InvMargLik1, LogInvs)
30   setTxtProgressBar(pb, s)
31 }
32 close(pb); summary(coda::mcmc(InvMargLik1))
33 mean(InvMargLik1)
34 BFFD <- exp(mean(InvMargLik2) - mean(InvMargLik1))
35 BFFD; mean(1/BFFD); 2*log(1/BFFD)

```

## 10.5 Summary

In this chapter, we introduced Bayesian model averaging (BMA) in generalized linear models. For linear Gaussian models, we perform BMA using three approaches: the Bayesian Information Criterion (BIC) approximation with

Occam’s window, the Markov Chain Monte Carlo Model Composition (MC3) algorithm, and conditional Bayes factors, which account for endogeneity. Additionally, we show how to perform dynamic Bayesian model averaging in state-space models, where forgetting parameters are used to facilitate computation. For other generalized linear models, such as logit, gamma, and Poisson, we demonstrate how to use the BIC approximation to perform BMA. Finally, we present alternative methods for calculating the marginal likelihood: the Savage-Dickey density ratio, Chib’s method, and the Gelfand-Dey method. These methods are particularly useful when the BIC approximation does not perform well due to small or moderate sample sizes.

However, a limitation of standard BMA is its implicit assumption that one of the candidate models is the true data-generating process. As Box famously noted, “Since all models are wrong, the scientist cannot obtain a ‘correct’ one by excessive elaboration. On the contrary, following William of Occam, he should seek an economical description of natural phenomena” [55].

This perspective has motivated new developments in BMA that relax the “true model” assumption and instead treat all models as misspecified. In these approaches, the Bayesian average of predictive densities is constructed using leave-one-out (LOO) predictive performance, and model weights are chosen to minimize a predictive loss function. One prominent example is *stacking of predictive distributions*, which, rather than weighting models by marginal likelihood, uses cross-validation (LOO) to assign weights that maximize out-of-sample predictive accuracy [463]. For a comprehensive review of Bayesian methods for aggregating predictive distributions, see [462].

## 10.6 Exercises

1. The Gaussian linear model specifies  $\mathbf{y} = \alpha \mathbf{i}_N + \mathbf{X}_m \boldsymbol{\beta}_m + \boldsymbol{\mu}_m$  such that  $\boldsymbol{\mu}_m \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ , and  $\mathbf{X}_m$  does not have the column of ones. Assuming that  $\pi(\sigma^2) \propto 1/\sigma^2$ ,  $\pi(\alpha) \propto 1$ , and  $\boldsymbol{\beta}_m | \sigma^2 \sim N(\mathbf{0}_{k_m}, \sigma^2 (g_m \mathbf{X}_m^\top \mathbf{X}_m)^{-1})$ .

- Show that the posterior conditional distribution of  $\boldsymbol{\beta}_m$  is  $N(\boldsymbol{\beta}_{mn}, \sigma^2 \mathbf{B}_{mn})$ , where  $\boldsymbol{\beta}_{mn} = \mathbf{B}_{mn} \mathbf{X}_m^\top \mathbf{y}$  and  $\mathbf{B}_{mn} = ((1 + g_m) \mathbf{X}_m^\top \mathbf{X}_m)^{-1}$ .
- Show that the marginal likelihood associated with model  $\mathcal{M}_m$  is proportional to

$$\begin{aligned} p(\mathbf{y} | \mathcal{M}_m) &\propto \left( \frac{g_m}{1 + g_m} \right)^{k_m/2} \\ &\times \left[ (\mathbf{y} - \bar{y} \mathbf{i}_N)^\top (\mathbf{y} - \bar{y} \mathbf{i}_N) - \frac{1}{1 + g_m} (\mathbf{y}^\top \mathbf{P}_{\mathbf{X}_m} \mathbf{y}) \right]^{-(N-1)/2}, \end{aligned}$$

where all parameter are indexed to model  $\mathcal{M}_m$ ,  $\mathbf{P}_{X_m} = \mathbf{X}_m(\mathbf{X}_m^\top \mathbf{X}_m)^{-1}\mathbf{X}_m$  is the projection matrix on the space generated by the columns of  $\mathbf{X}_m$ , and  $\bar{y}$  is the sample mean of  $\mathbf{y}$ .

Hint: Take into account that  $\mathbf{i}_N^\top \mathbf{X}_m = \mathbf{0}_{k_m}$  due to all columns being centered with respect to their means.

## 2. Determinants of export diversification I

[225] use BMA to study the determinants of export diversification. Use the dataset *10ExportDiversificationHHI.csv* to perform BMA using the BIC approximation and MC3 to check if these two approaches agree.

## 3. Simulation exercise of the Markov chain Monte Carlo model composition continues

Program an algorithm to perform MC3 where the final  $S$  models are unique. Use the simulation setting of Section 10.2 increasing the number of regressors to 40, this implies approximately  $1.1e+12$  models.

## 4. Simulation exercise of IV BMA continues

Use the simulation setting with endogeneity in Section 10.2 to perform BMA based on the BIC approximation and MC3.

## 5. Determinants of export diversification II

Use the datasets *11ExportDiversificationHHI.csv* and *12ExportDiversificationHHInstr.csv* to perform IV BMA assuming that the log of per capita gross domestic product is endogenous (*avgldpcap*). See [225] for details.

6. Show that the link function in the case of the Bernoulli distribution is  $\log\left(\frac{\theta}{1-\theta}\right)$ .
7. [356, 357] perform variable selection using the file *13InternetMed.csv*. In this dataset, the dependent variable is an indicator of Internet adoption (*internet*) for 5,000 households in Medellín (Colombia) during the period 2006–2014. This dataset contains information about 18 potential determinants, which implies 262,144 ( $2^{18}$ ) potential models, considering only variable uncertainty (see these papers for details about the dataset). Perform BMA using the logit link function with this dataset.
8. [406] use the file *14ValueFootballPlayers.csv* to analyze the market value of soccer players in the most important leagues in Europe. In particular, there are 26 potential determinants of the market value (dependent variable) of a stratified sample of 335 soccer players from the five most important leagues in Europe (see [406] for details). Use

this dataset to perform BMA using the gamma distribution, setting the default values for Occam's window.

9. Use the dataset *15Fertile2.csv* from [456, p. 547] to perform BMA using the Poisson model with the log link. This dataset contains information about 1,781 women from Botswana in 1988 (for details, see <https://rdrr.io/cran/wooldridge/man/fertil2.html>), with the caveat that we deleted some variables and omitted observations with NA values. The dependent variable is the number of children ever born (*ceb*), which is a count variable, as a function of 19 potential determinants.
10. Perform BMA in the logit model using MC3 and the BIC approximation using the simulation setting of Section 10.3.
11. Use the dataset *19ExchangeRateCOPUSD.csv* to perform dynamic BMA using four different *state-space models* to explain the annual variation of the COP to USD exchange rate:

•Interest rate parity

$$\Delta e_t = \beta_{1t}^{IRP} + \beta_{2t}^{IRP}(i_{t-1}^{Col} - i_{t-1}^{USA}) + \mu_t^{IRP}$$

•Purchasing power parity

$$\Delta e_t = \beta_{1t}^{PPP} + \beta_{2t}^{PPP}(\pi_{t-1}^{Col} - \pi_{t-1}^{USA}) + \mu_t^{PPP}$$

•Taylor rule

$$\Delta e_t = \beta_{1t}^{Taylor} + \beta_{2t}^{Taylor}(\pi_{t-1}^{Col} - \pi_{t-1}^{USA}) + \beta_{2t}^{Taylor}(g_{t-1}^{Col} - g_{t-1}^{USA}) + \mu_t^{Taylor}$$

•Money supply

$$\Delta e_t = \beta_{1t}^{Money} + \beta_{2t}^{Money}(g_{t-1}^{Col} - g_{t-1}^{USA}) + \beta_{2t}^{Money}(m_{t-1}^{Col} - m_{t-1}^{USA}) + \mu_t^{Money}$$

where *varTRM* ( $\Delta e_t$ ) represents the annual variation rate of the exchange rate from COP to USD, *TES\_COL10* ( $i_t^{Col}$ ) and *TES\_USA10* ( $i_t^{USA}$ ) denote the annual return rates of Colombian and U.S. public debts over 10 years, *inflation\_COL* ( $\pi_t^{Col}$ ) and *inflation\_USA* ( $\pi_t^{USA}$ ) are the annual inflation rates for Colombia and the U.S., *varISE\_COL* ( $g_t^{Col}$ ) and *varISE\_USA* ( $g_t^{USA}$ ) represent the annual variations of economic activity indices, and *varCOL\_M3* ( $m_t^{Col}$ ) and *varUSA\_M3* ( $m_t^{USA}$ ) are the annual variations of the money supply. In addition,  $\mu_t$  is the stochastic error. The dataset includes monthly variations from January 2006 to November 2023.

Perform Bayesian model averaging using these four models to explain the annual variation of the exchange rate, calculate the posterior model probabilities, and plot the posterior mean and credible interval of  $\beta_{2t}^{Money}$ .

12. Perform a simulation of the dynamic logistic model, where there are  $7$  ( $2^3 - 1$ , excluding the model without regressors) competing

models originating from 3 regressors:  $x_{tk} \sim N(0.5, 0.8^2)$ ,  $k = 2, 3, 4$ , and  $\beta_1 = 0.5$ ,  $\beta_{2t}$  is a sequence from 1 to 2 in steps given by  $1/T$ , and  $\beta_{3t} = \begin{cases} -1, & 1 < t \leq 0.5T \\ 0, & 0.5T < t \leq T \end{cases}$ , with  $\beta_4 = 1.2$ . Then,  $\mathbf{x}_t^\top \boldsymbol{\beta}_t = \beta_1 + \beta_{2t}x_{2t} + \beta_{3t}x_{3t} + \beta_4x_{4t}$ , where

$$P[Y_t = 1 | \mathbf{x}_t, \boldsymbol{\beta}_t] = \frac{\exp(\mathbf{x}_t^\top \boldsymbol{\beta}_t)}{1 + \exp(\mathbf{x}_t^\top \boldsymbol{\beta}_t)}, \quad t = 1, 2, \dots, 1100.$$

Use the function *logistic.dma* from the *dma* package to obtain the posterior model probabilities, first setting the forgetting parameter of the models to 0.99, and then to 0.95. Compare the results.

13. Show that

$$\mathbb{E} \left[ \frac{q(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta} | \mathcal{M}_m) p(\mathbf{y} | \boldsymbol{\theta}_m, \mathcal{M}_m)} \middle| \mathbf{y}, \mathcal{M}_m \right] = \frac{1}{p(\mathbf{y} | \mathcal{M}_m)},$$

where the expected value is with respect to the posterior distribution given the model  $\mathcal{M}_m$ , and  $q(\boldsymbol{\theta})$  is the proposal distribution whose support is  $\Theta$ .



## Part III

# *Advanced* methods: A brief introduction



# 11

---

## Semi-parametric and non-parametric models

---

Non-parametric models are characterized by making minimal assumptions about the data-generating process. Unlike parametric models, which have a finite-dimensional parameter space, non-parametric models often involve infinite-dimensional parameter spaces. A major challenge in non-parametric modeling is the *curse of dimensionality*, as these models require dense data coverage, necessitating large datasets to achieve reliable estimates.

Semi-parametric methods, on the other hand, combine parametric assumptions for part of the model with non-parametric assumptions for the rest. This approach offers a balance between flexibility, tractability and applicability.

In this chapter, we introduce finite Gaussian mixture models (GMM) and Dirichlet mixture processes (DMP), the latter representing an infinite mixture. Both can be used to specify an entire statistical model (nonparametric specification) or to model stochastic error distributions in a semiparametric framework.

Additionally, we present spline models, where the outcome depends linearly on smooth nonparametric functions. To address the curse of dimensionality, we introduce partially linear models, which mitigate this issue while remaining interpretable and flexible for practical applications.

We let other useful Bayesian non-parametric approaches like Bayesian additive random trees (BART) and Gaussian process (GP) for Chapter 12.

---

### 11.1 Mixture models

Mixture models naturally arise in situations where a sample consists of draws from different *subpopulations* (*clusters*) that cannot be easily distinguished based on observable characteristics. However, performing inference on specific identified subpopulations can be misleading if the assumed distribution for each cluster is misspecified.

Even when distinct subpopulations do not exist, finite and infinite mixture models provide a useful framework for semi-parametric inference. They

effectively approximate distributions with skewness, excess kurtosis, and multimodality, making them useful for modeling stochastic errors.

In addition, mixture models help capture unobserved heterogeneity. That is, as data modelers, we may observe individuals with identical sets of observable variables but entirely different response variables. These differences cannot be explained solely by sampling variability; rather, they suggest the presence of an unobserved underlying process, independent of the observable features, that accounts for this pattern.

### 11.1.1 Finite Gaussian mixtures

A finite Gaussian mixture model for regression with  $H$  known components assumes that a sample  $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^\top$  consists of observations  $y_i$ , for  $i = 1, 2, \dots, N$ , where each  $y_i$  is generated from one of the  $H$  components,  $h = 1, 2, \dots, H$ , conditional on the regressors  $\mathbf{x}_i$ . Specifically, we assume

$$y_i \mid \mathbf{x}_i \sim N(\mathbf{x}_i^\top \boldsymbol{\beta}_h, \sigma_h^2).$$

Thus, the sampling distribution of  $y_i$  is given by

$$p(y_i \mid \{\lambda_h, \boldsymbol{\beta}_h, \sigma_h^2\}_{h=1}^H, \mathbf{x}_i) = \sum_{h=1}^H \lambda_h \phi(y_i \mid \mathbf{x}_i^\top \boldsymbol{\beta}_h, \sigma_h^2),$$

where  $\phi(y_i \mid \mathbf{x}_i^\top \boldsymbol{\beta}_h, \sigma_h^2)$  is the Gaussian density with mean  $\mathbf{x}_i^\top \boldsymbol{\beta}_h$  and variance  $\sigma_h^2$ ,  $0 < \lambda_h < 1$  represents the proportion of the population belonging to subpopulation  $h$ , and the weights satisfy  $\sum_{h=1}^H \lambda_h = 1$ .

Then, we allow cross-sectional units to differ according to unobserved clusters (subpopulations) that exhibit homogeneous behavior within each cluster.

To model a finite Gaussian mixture, we introduce an individual cluster indicator or latent class  $\psi_{ih}$  such that

$$\psi_{ih} = \begin{cases} 1, & \text{if the } i\text{-th unit is drawn from the } h\text{-th cluster,} \\ 0, & \text{otherwise.} \end{cases}$$

Thus,  $P(\psi_{ih} = 1) = \lambda_h$  for all clusters  $h = 1, 2, \dots, H$  and units  $i = 1, 2, \dots, N$ . Note that a high probability of individuals belonging to the same cluster suggests that these clusters capture similar sources of unobserved heterogeneity.

This setting implies that

$$\boldsymbol{\psi}_i = [\psi_{i1} \ \psi_{i2} \ \dots \ \psi_{iH}]^\top \sim \text{Categorical}(\boldsymbol{\lambda}),$$

where  $\boldsymbol{\lambda} = [\lambda_1 \ \lambda_2 \ \dots \ \lambda_H]^\top$  represents the event probabilities.

We know from Subsection 3.2.1 that the Dirichlet prior distribution is

conjugate to the multinomial distribution, where the categorical distribution is a special case in which the number of trials is one. Thus, we assume that

$$\pi(\boldsymbol{\lambda}) \sim \text{Dir}(\boldsymbol{\alpha}_0),$$

where  $\boldsymbol{\alpha}_0 = [\alpha_{10} \ \alpha_{20} \ \dots \ \alpha_{H0}]^\top$ ,  $\alpha_{h0} = 1/H$  is recommended by [160, p. 535].

Observe that we are using a hierarchical structure, as we specify a prior on  $\boldsymbol{\lambda}$ , which serves as the hyperparameter for the cluster indicators. In addition, we can assume conjugate families for the location and scale parameters to facilitate computation, that is,  $\boldsymbol{\beta}_h \sim N(\boldsymbol{\beta}_{h0}, \mathbf{B}_{h0})$  and  $\sigma_h^2 \sim IG(\alpha_{h0}/2, \delta_{h0}/2)$ .

This setting allows to obtain standard conditional posterior distributions:

$$\boldsymbol{\beta}_h \sim N(\boldsymbol{\beta}_{hn}, \mathbf{B}_{hn}),$$

where  $\mathbf{B}_{hn} = (\mathbf{B}_{h0}^{-1} + \sigma_h^{-2} \sum_{\{i:\psi_{ih}=1\}} \mathbf{x}_i \mathbf{x}_i^\top)^{-1}$  and  $\boldsymbol{\beta}_{hn} = \mathbf{B}_{hn}(\mathbf{B}_{h0}^{-1} \boldsymbol{\beta}_{h0} + \sigma_h^{-2} \sum_{\{i:\psi_{ih}=1\}} \mathbf{x}_i y_i)$ .

$$\sigma_h^2 \sim IG(\alpha_{hn}/2, \delta_{hn}/2),$$

where  $\alpha_{hn} = \alpha_{h0} + N_h$ ,  $\delta_{hn} = \delta_{h0} + \sum_{\{i:\psi_{ih}=1\}} (y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_h)^2$ , and  $N_h$  is the number of units in cluster  $h$ .

$$\boldsymbol{\lambda} \sim \text{Dir}(\boldsymbol{\alpha}_n),$$

where  $\boldsymbol{\alpha}_n = [\alpha_{1n} \ \alpha_{2n} \ \dots \ \alpha_{Hn}]^\top$ , and  $\alpha_{hn} = \alpha_{h0} + N_h$ .

$$\psi_{in} \sim \text{Categorical}(\boldsymbol{\lambda}_n),$$

where  $P(\psi_{ih} = 1) = \frac{\lambda_h \phi(y_i | \mathbf{x}_i^\top \boldsymbol{\beta}_h, \sigma_h^2)}{\sum_{j=1}^H \lambda_j \phi(y_i | \mathbf{x}_i^\top \boldsymbol{\beta}_j, \sigma_j^2)}$ .

In general, it is always safer to perform inference in mixture models using informative priors, as non-informative priors may have unintended consequences on posterior inference. One way to facilitate prior elicitation is to work with standardized data, as this removes dependence on measurement units. Another useful approach is to specify the model in log-log form so that the coefficients can be interpreted as elasticities or semi-elasticities. As always in Bayesian inference, it makes sense to perform a sensitivity analysis with respect to the hyperparameters.

Mixture models have the label-switching identification problem, meaning they are nonidentifiable because the distribution remains unchanged if the group labels are permuted [440]. For instance, a mixture model with two components can be characterized by  $\{\lambda_1, \boldsymbol{\beta}_1, \sigma_1^2\}$  for the first cluster and  $\{1 - \lambda_1, \boldsymbol{\beta}_2, \sigma_2^2\}$  for the second. However, an alternative characterization is  $\{1 - \lambda_1, \boldsymbol{\beta}_2, \sigma_2^2\}$  for cluster 1 and  $\{\lambda_1, \boldsymbol{\beta}_1, \sigma_1^2\}$  for cluster 2. This parametrization yields exactly the same likelihood as the first one, meaning any permutation of the cluster labels leaves the likelihood unchanged. Consequently, the posterior draws of each component-specific parameter target the same distribution.

Label switching may pose challenges when performing inference on specific mixture components, such as in the regression analysis presented here. However, it is not an issue when inference on specific components is unnecessary, as in cases where mixtures are used to model stochastic errors in semi-parametric settings. In the former case, post-processing strategies can mitigate the issue, such as *random permutation of latent classes* (see the simulation exercise below, [160, p. 534] and Algorithm 3.5 in [150, p. 82]).

A semi-parametric regression imposes a specific structure in part of the model and uses flexible assumptions in another part, for instance

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \mu_i,$$

$$p(\mu_i \mid \{\lambda_h, \mu_h, \sigma_h^2\}_{h=1}^H) = \sum_{h=1}^H \lambda_h \phi(\mu_i \mid \mu_h, \sigma_h^2).$$

Thus, the distribution of the stochastic error is a finite Gaussian mixture. Note that the mean of the stochastic error is not equal to zero; consequently, the intercept in the regression should be removed, as these two parameters are not separately identifiable [440]. Additionally, this approach allows for multiple modes and asymmetric distributions of the stochastic errors, providing greater flexibility.

We can use a Gibbs sampling algorithm in this semi-parametric specification if we assume conjugate families. The difference from the previous setting is that we have the same slope parameters; thus,  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$ . Additionally, we must specify the prior distribution for the means of the stochastic errors, given by  $\mu_h \sim N(\mu_{h0}, \sigma_{h0}^2)$ . Then, the posterior distributions are:

$$\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_n, \mathbf{B}_n),$$

where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \sum_{h=1}^H \sum_{\{i:\psi_{ih}=1\}} \sigma_h^{-2} \mathbf{x}_i \mathbf{x}_i^\top)^{-1}$  and  $\boldsymbol{\beta}_n = \mathbf{B}_n (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{h=1}^H \sum_{\{i:\psi_{ih}=1\}} \sigma_h^{-2} \mathbf{x}_i (y_i - \mu_h))$ .

$$\sigma_h^2 \sim IG(\alpha_{hn}/2, \delta_{hn}/2),$$

where  $\alpha_{hn} = \alpha_{h0} + N_h$ ,  $\delta_{hn} = \delta_{h0} + \sum_{\{i:\psi_{ih}=1\}} (y_i - \mu_h - \mathbf{x}_i^\top \boldsymbol{\beta}_h)^2$ , and  $N_h$  is the number of units in cluster  $h$ .

$$\mu_h \sim N(\mu_{hn}, \sigma_{hn}^2),$$

where  $\sigma_{hn}^2 = \left( \frac{1}{\sigma_{h0}^2} + \frac{N_h}{\sigma_h^2} \right)^{-1}$  and  $\mu_{hn} = \sigma_{hn}^2 \left( \frac{\mu_{h0}}{\sigma_{h0}^2} + \frac{\sum_{\{i:\psi_{ih}=1\}} (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})}{\sigma_h^2} \right)$ .

$$\boldsymbol{\lambda} \sim \text{Dir}(\boldsymbol{\alpha}_n),$$

where  $\boldsymbol{\alpha}_n = [\alpha_{1n} \ \alpha_{2n} \ \dots \ \alpha_{Hn}]^\top$ , and  $\alpha_{hn} = \alpha_{h0} + N_h$ .

$$\boldsymbol{\psi}_{in} \sim \text{Categorical}(\boldsymbol{\lambda}_n),$$

where  $P(\psi_{ih} = 1) = \frac{\lambda_h \phi(y_i | \mu_h + \mathbf{x}_i^\top \boldsymbol{\beta}, \sigma_h^2)}{\sum_{j=1}^H \lambda_j \phi(y_i | \mu_j + \mathbf{x}_i^\top \boldsymbol{\beta}, \sigma_j^2)}$ .

A potential limitation of finite mixture models is the need to specify the number of components in advance. One approach is to estimate the model for different values of  $H$  and then compute the marginal likelihood to select the model best supported by the data. However, this procedure can be tedious. A simpler strategy is to set  $H$  large enough (e.g., 10 components), assign  $\alpha_{h0} = 1/H$ , and perform an initial run of the algorithm. If we are not interested in the specific composition of clusters, this approach is sufficient. Otherwise, the posterior distribution of  $H$  can be obtained by tracking the number of nonempty clusters in each iteration. In a second run,  $H$  can then be fixed at the mode of this posterior distribution.

However, the previous approaches ultimately fix the number of components. Consequently, finite mixtures cannot be considered a non-parametric method [378], as they lack an automatic mechanism to increase  $H$  as the sample size grows. An alternative is to avoid pre-specifying the number of components altogether by using a Dirichlet process mixture (DPM). This is the topic of the next section.

#### **Example: Simulation exercises**

First, let's illustrate the label-switching issue using a simple model without regressors, assuming the same known variance. Consider the following distribution:

$$p(y_i) = 0.75\phi(y_i | \beta_{01}, 1^2) + 0.25\phi(y_i | \beta_{02}, 1^2), \quad i = 1, 2, \dots, 500.$$

Initially, we set  $\beta_{01} = 0.5$  and  $\beta_{02} = 2.5$ . We perform 1,000 MCMC iterations, with a burn-in period of 500 and a thinning factor of 2. The following code demonstrates how to implement the Gibbs sampler using a prior normal distribution with mean 0 and variance 10, with the hyperparameters of the Dirichlet distribution set to 1/2.

#### **R code. Simulation exercise: Label switching issue**

```

1 ### Simulation exercise: Label switching issue ###
2 rm(list = ls()); set.seed(10101); library(ggplot2)
3 # Simulate data from a 2-component mixture model
4 n <- 500; z <- rbinom(n, 1, 0.75) # Latent class indicator
5 y <- ifelse(z == 1, rnorm(n, 0.5, 1), rnorm(n, 2.5, 1))
6 data <- data.frame(y)

```

### R code. Simulation exercise: Label switching issue

```

1 # Plot
2 ggplot(data, aes(x = y)) +
3   geom_density(fill = "blue", alpha = 0.3) + labs(title = "
4     Density Plot", x = "y", y = "Density") + theme_minimal()
5 # Hyperparameters
6 mu0 <- 0; sig2mu0 <- 10; H <- 2; a0h <- rep(1/H, H)
7 # MCMC parameters
8 mcmc <- 1000; burnin <- 500; tot <- mcmc + burnin; thin <- 2
9 # Gibbs sampling functions
10 Postmu <- function(yh){
11   Nh <- length(yh)
12   sig2mu <- (1/sig2mu0 + Nh)^(-1)
13   mun <- sig2mu*(mu0/sig2mu0 + sum(yh))
14   mu <- rnorm(1, mun, sig2mu^0.5)
15   return(mu)
16 }
17 PostPsi <- matrix(NA, tot, n); PostMu <- matrix(NA, tot, H)
18 PostLambda <- rep(NA, tot)
19 Id1 <- which(y <= 1) # 1 is from inspection of the density
20           plot of y
21 Id2 <- which(y > 1)
22 N1 <- length(Id1); N2 <- length(Id2)
23 Lambda <- c(N1/n, N2/n)
24 MU <- c(mean(y[Id1]), mean(y[Id2])); Psi <- rep(NA, n)
25 pb <- txtProgressBar(min = 0, max = tot, style = 3)
26 for(s in 1:tot){
27   for(i in 1:n){
28     lambdai <- NULL
29     for(h in 1:H){
30       lambdaih <- Lambda[h]*dnorm(y[i], MU[h], 1)
31       lambdai <- c(lambdai, lambdaih)
32     }
33     Psi[i] <- sample(1:H, 1, prob = lambdai)
34   }
35   PostPsi[s, ] <- Psi
36   for(h in 1:H){
37     idh <- which(Psi == h);      MU[h] <- Postmu(yh = y[idh])
38   }
39   PostMu[s, ] <- MU;
40   Lambda <- sort(MCMCpack::rdirichlet(1, a0h + table(Psi)),
41                 decreasing = TRUE)
42   PostLambda[s] <- Lambda[1]
43   setTxtProgressBar(pb, s)
44 }
45 close(pb)

```

### R code. Simulation exercise: Label switching issue

```

1 keep <- seq(burnin, tot, thin)
2 PosteriorMUs <- coda::mcmc(PostMu[keep,])
3 summary(PosteriorMUs); plot(PosteriorMUs)
4 dfMU <- data.frame(mu1 = PostMu[keep,1], mu2 = PostMu[keep
   ,2])
5 # Plot
6 require(latex2exp)
7 ggplot(dfMU) + geom_density(aes(x = mu1, color = "mu1"),
   linewidth = 1) + geom_density(aes(x = mu2, color = "mu2"
   ), linewidth = 1) + labs(title = "Density Plot", x = TeX
   ("$\\mu$"), y = "Density", color = "Variable") + theme_
   minimal() + scale_color_manual(values = c("mu1" = "blue"
   , "mu2" = "red"))
8

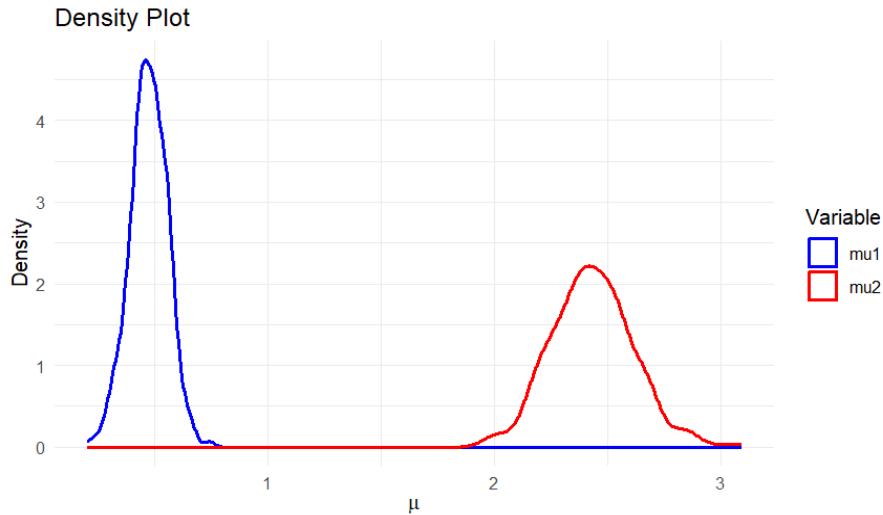
```

Figure 11.1 shows the posterior densities of the location parameters. The posterior means are 0.42 and 2.50, with 95% credible intervals of (0.07, 0.71) and (2.34, 2.65), respectively. The posterior mean of the probability is 0.27, with a 95% credible interval of (0.19, 0.35). Note that in this simple simulation exercise, we did not observe unintended consequences from using non-informative priors and not standardizing the data. However, real-world applications should take these aspects into account.

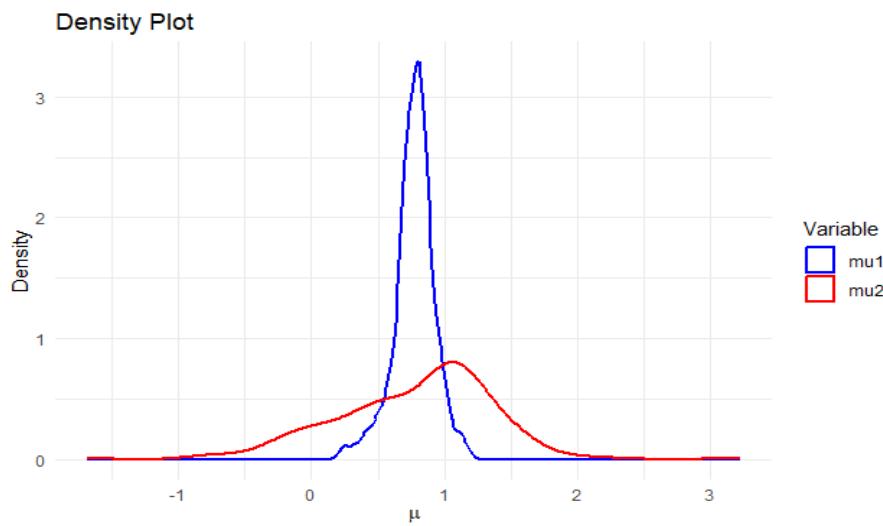
We perform the same exercise assuming  $\beta_{01} = 0.5$  and  $\beta_{02} = 1$ . Figure 11.2 shows the posterior densities, where we observe significant overlap. The posterior means are 0.77 in both cases, with 95% credible intervals of (0.40, 1.05) and (-0.44, 1.71). The posterior mean of the probability is 0.84.

In the second setting, the posterior draws of the Gibbs sampler can switch between the two means because they are relatively close. This situation contrasts with the first example, where there is a relatively large separation between the means, resulting in a region of the parameter space with zero probability (see the flat region between the two posterior distributions in Figure 11.1). The key point is that, given a sufficiently large number of Gibbs sampler iterations, the algorithm should eventually explore the entire parameter space and encounter the label-switching issue. This occurs because both posterior chains should exhibit similar behavior, as they are targeting the same distribution.

We can implement *random permutation of latent classes* to address this issue. This involves sampling a random permutation of the labels at each iteration of the MCMC algorithm. For example, with three clusters, there are  $3! = 6$  possible label permutations. Let the permutations be labeled as  $p_k = \{p_k(1), p_k(2), \dots, p_k(H)\}, k = 1, 2, \dots, H!$ . At the end of each iteration

**FIGURE 11.1**

Posterior distributions: Mean parameters, population values  $\beta_{01} = 0.5$  and  $\beta_{02} = 2.5$ .

**FIGURE 11.2**

Posterior distributions: Mean parameters, population values  $\beta_{01} = 0.5$  and  $\beta_{02} = 1$ .

in the MCMC algorithm, we randomly select one of the permutations  $\boldsymbol{p}_k$  and replace the cluster probabilities  $\lambda_1^{(s)}, \dots, \lambda_H^{(s)}$  with  $\lambda_{p_k(1)}^{(s)}, \dots, \lambda_{p_k(H)}^{(s)}$ . We apply the same permutation to  $\boldsymbol{\beta}^{(s)}$ ,  $\sigma^2(s)$ , and  $\boldsymbol{\psi}_i^{(s)}$ , for  $i = 1, 2, \dots, n$ . The following algorithm illustrates how to implement this in our simple example.

### R code. Simulation exercise: Random permutation of latent classes

```

1 ##### Permutations #####
2 rm(list = ls()); set.seed(010101); library(ggplot2)
3 # Simulate data from a 2-component mixture model
4 n <- 500
5 z <- rbinom(n, 1, 0.75) # Latent class indicator
6 y <- ifelse(z == 1, rnorm(n, 0.5, 1), rnorm(n, 2.5, 1))
7 # Hyperparameters
8 mu0 <- 0; sig2mu0 <- 10; H <- 2; a0h <- rep(1/H, H)
9 # MCMC parameters
10 mcmc <- 2000; burnin <- 500
11 tot <- mcmc + burnin; thin <- 2
12 # Gibbs sampling functions
13 Postmu <- function(yh){
14   Nh <- length(yh)
15   sig2mu <- (1/sig2mu0 + Nh)^(-1)
16   mun <- sig2mu*(mu0/sig2mu0 + sum(yh))
17   mu <- rnorm(1, mun, sig2mu^0.5)
18   return(mu)
19 }
20 PostPsi <- matrix(NA, tot, n); PostMu <- matrix(NA, tot, H)
21 PostLambda <- rep(NA, tot)
22 Id1 <- which(y <= 1); Id2 <- which(y > 1)
23 N1 <- length(Id1); N2 <- length(Id2)
24 Lambda <- c(N1/n, N2/n); MU <- c(mean(y[Id1]), mean(y[Id2]))
25 Psi <- rep(NA, n); per1 <- c(1,2); per2 <- c(2,1)
26 pb <- txtProgressBar(min = 0, max = tot, style = 3)
27 for(s in 1:tot){
28   for(i in 1:n){
29     lambdai <- NULL
30     for(h in 1:H){
31       lambdaih <- Lambda[h]*dnorm(y[i], MU[h], 1)
32       lambdai <- c(lambdai, lambdaih)
33     }
34     Psi[i] <- sample(1:H, 1, prob = lambdai)
35   }
36   for(h in 1:H){
37     idh <- which(Psi == h)
38     MU[h] <- Postmu(yh = y[idh])
39   }
40   Lambda <- MCMCpack::rdirichlet(1, a0h + table(Psi))

```

**R code. Simulation exercise: Random permutation of latent classes**

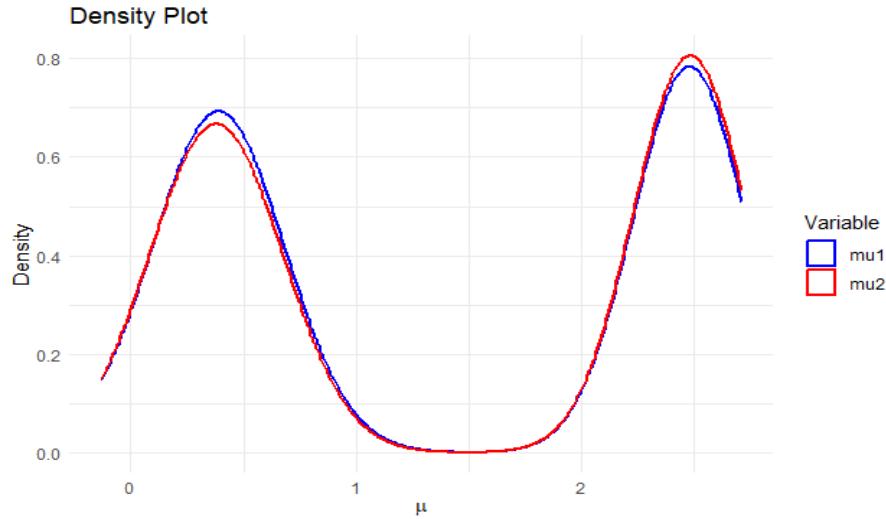
```

1  # Permutations
2  labels <- sample(1:2, 1, prob = c(0.5, 0.5))
3  if(labels == 2){
4      Lambda <- Lambda[per2]
5      MU <- MU[per2]
6      for(i in 1:n){
7          if(Psi[i] == 1){Psi[i] <- 2
8          }else{Psi[i] <- 1}
9      }
10 }
11 PostPsi[s, ] <- Psi; PostMu[s, ] <- MU
12 PostLambda[s] <- Lambda[1]
13 setTxtProgressBar(pb, s)
14 }
15 close(pb)
16 keep <- seq(burnin, tot, thin)
17 PosteriorMUs <- coda::mcmc(PostMu[keep,])
18 summary(PosteriorMUs)
19 plot(PosteriorMUs)
20 dfMU <- data.frame(mu1 = PostMu[keep,1], mu2 = PostMu[keep
21 ,2])
22 # Plot
23 require(latex2exp)
24 ggplot(dfMU) +
25 geom_density(aes(x = mu1, color = "mu1"), linewidth = 1) +
# First density plot
26 geom_density(aes(x = mu2, color = "mu2"), linewidth = 1) +
# Second density plot
27 labs(title = "Density Plot", x = TeX("$\mu$"), y = "Density
", color = "Variable") + theme_minimal() + scale_color_
manual(values = c("mu1" = "blue", "mu2" = "red")) #
Custom colors
28 PosteriorLAMBDA <- coda::mcmc(PostLambda[keep])
29 summary(PosteriorLAMBDA)
30 plot(PosteriorLAMBDA)
31

```

Figure 11.3 shows the posterior distributions from the random permutation of latent classes in the first simulation setting. We observe that both posterior distributions look similar, as both are targeting a bimodal distribution given by two clusters.

In the following setting we simulate a simple regression mixture with two components such that  $\psi_{i1} \sim \text{Ber}(0.5)$ , consequently,  $\psi_{i2} = 1 - \psi_{i1}$ , and assume

**FIGURE 11.3**

Posterior distributions: Mean parameters, population values  $\beta_{01} = 0.5$  and  $\beta_{02} = 2.5$ .

one regressor,  $x_i \sim N(0, 1)$ ,  $i = 1, 2, \dots, 1,000$ . Then,

$$p(y_i | \mathbf{x}_i) = 0.5\phi(y_i | 2 + 1.5x_i, 1^2) + 0.5\phi(y_i | -1 + 0.5x_i, 0.8^2).$$

The following code shows how to perform inference in this model, assuming  $N(0, 5)$  and  $N(0, 2)$  priors for the intercepts and slopes, respectively. Additionally, we use a  $Cauchy(0, 2)$  prior truncated at 0 for the standard deviations, and a  $Dirichlet(1, 1)$  prior for the probabilities. We use the *brms* package in **R**, which in turn uses *Stan*, setting number of MCMC iterations 2,000, a burn-in (warm-up) equal to 1,000, and 4 chains. Remember that *Stan* software uses Hamiltonian Monte Carlo.

The following code performs inference in this simulation from scratch using Gibbs sampling. We do not implement the random permutation of latent classes algorithm for facilitating exposition and comparability with the results from the package *brms*. We use non-informative priors, setting  $\alpha_{h0} = \delta_{h0} = 0.01$ ,  $\beta_{h0} = \mathbf{0}_2$ ,  $\mathbf{B}_{h0} = \mathbf{I}_2$ , and  $\boldsymbol{\alpha}_0 = [1/2 \ 1/2]^\top$ . The number of MCMC iterations is 5,000, the burn-in is 1,000, and the thinning parameter is 2. In general, the Gibbs sampler appears to yield good posterior results as all 95% intervals encompass the population parameters.

**R code. Simulation exercise: Gaussian mixture with 2 components using brms package**

```

1 ##### Simulation exercise: Gaussian mixture: 2 components #####
2 rm(list = ls()); set.seed(010101)
3 library(brms); library(ggplot2)
4 # Simulate data from a 2-component mixture model
5 n <- 1000
6 x <- rnorm(n)
7 z <- rbinom(n, 1, 0.5) # Latent class indicator
8 y <- ifelse(z == 1, rnorm(n, 2 + 1.5*x, 1), rnorm(n, -1 +
  0.5*x, 0.8))
9 data <- data.frame(y, x)
10 # Plot
11 ggplot(data, aes(x = y)) +
12   geom_density(fill = "blue", alpha = 0.3) + # Density plot
  with fill color
13   labs(title = "Density Plot", x = "y", y = "Density") + theme
  _minimal()
14 # Define priors
15 priors <- c(
16   set_prior("normal(0, 5)", class = "Intercept", dpar = "mu1"),
  , # First component intercept
17   set_prior("normal(0, 5)", class = "Intercept", dpar = "mu2"),
  , # Second component intercept
18   set_prior("normal(0, 2)", class = "b", dpar = "mu1"), #
  First component slope
19   set_prior("normal(0, 2)", class = "b", dpar = "mu2"), #
  Second component slope
20   set_prior("cauchy(0, 2)", class = "sigma1", lb = 0), #
  First component sigma
21   set_prior("cauchy(0, 2)", class = "sigma2", lb = 0), #
  Second component sigma
22   set_prior("dirichlet(1, 1)", class = "theta") # Mixing
  proportions
23 )
24 # Fit a 2-component Gaussian mixture regression model
25 fit <- brm(
26   bf(y ~ 1 + x, family = mixture(gaussian, gaussian)), # Two
  normal distributions
27   data = data,
28   prior = priors,
29   chains = 4, iter = 2000, warmup = 1000, cores = 4
30 )
31 prior_summary(fit) # Summary of priors
32 summary(fit) # Summary of posterior draws
33 plot(fit) # Plots of posterior draws

```

**R code. Simulation exercise: Gaussian mixture with 2 components using Gibbs sampling**

```

1 ##### Perform inference from scratch #####
2 rm(list = ls()); set.seed(010101)
3 library(brms); library(ggplot2)
4 # Simulate data from a 2-component mixture model
5 n <- 1000
6 x <- rnorm(n)
7 z <- rbinom(n, 1, 0.5) # Latent class indicator
8 y <- ifelse(z == 1, rnorm(n, 2 + 1.5*x, 1), rnorm(n, -1 +
  0.5*x, 0.8))
9 # Hyperparameters
10 d0 <- 0.001; a0 <- 0.001
11 b0 <- rep(0, 2); B0 <- diag(2); B0i <- solve(B0)
12 a01 <- 1/2; a02 <- 1/2
13 # MCMC parameters
14 mcmc <- 5000; burnin <- 1000
15 tot <- mcmc + burnin; thin <- 2
16 # Gibbs sampling functions
17 PostSig2 <- function(Betah, Xh, yh){
18   Nh <- length(yh); an <- a0 + Nh
19   dn <- d0 + t(yh - Xh%*%Betah)%*%(yh - Xh%*%Betah)
20   sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
21   return(sig2)
22 }
23 PostBeta <- function(sig2h, Xh, yh){
24   Bn <- solve(B0i + sig2h^(-1)*t(Xh)%*%Xh)
25   bn <- Bn%*%(B0i%*%b0 + sig2h^(-1)*t(Xh)%*%yh)
26   Beta <- MASS::mvrnorm(1, bn, Bn)
27   return(Beta)
28 }
29 PostBetas1 <- matrix(0, mcmc+burnin, 2)
30 PostBetas2 <- matrix(0, mcmc+burnin, 2)
31 PostSigma21 <- rep(0, mcmc+burnin)
32 PostSigma22 <- rep(0, mcmc+burnin)
33 PostPsi <- matrix(0, mcmc+burnin, n)
34 PostLambda <- rep(0, mcmc+burnin)
35 Id1 <- which(y<1) # 1 is from inspection of the density plot
  of y
36 N1 <- length(Id1); Lambda1 <- N1/n
37 Id2 <- which(y>=1)
38 N2 <- length(Id2); Lambda2 <- N2/n
39 Reg1 <- lm(y ~ x, subset = Id1)
40 SumReg1 <- summary(Reg1); Beta1 <- Reg1$coefficients
41 sig21 <- SumReg1$sigma^2
42 Reg2 <- lm(y ~ x, subset = Id2); SumReg2 <- summary(Reg2)
43 Beta2 <- Reg2$coefficients
44 sig22 <- SumReg2$sigma^2
45 X <- cbind(1, x); Psi <- rep(NA, n)

```

**R code. Simulation exercise: Gaussian mixture with 2 components using Gibbs sampling**

```

1 pb <- txtProgressBar(min = 0, max = tot, style = 3)
2 for(s in 1:tot){
3   for(i in 1:n){
4     lambdai1 <- Lambda1*dnorm(y[i], X[i,] %*% Beta1, sig21
5       ^0.5)
5     lambdai2 <- Lambda2*dnorm(y[i], X[i,] %*% Beta2, sig22
6       ^0.5)
6     Psi[i] <- sample(c(1,2), 1, prob = c(lambdai1, lambdai2)
7       )
7   }
8   PostPsi[s, ] <- Psi
9   Id1 <- which(Psi == 1); Id2 <- which(Psi == 2)
10  N1 <- length(Id1); N2 <- length(Id2)
11  sig21 <- PostSig2(Betah = Beta1, Xh = X[Id1, ], yh = y[Id1
12    ])
12  sig22 <- PostSig2(Betah = Beta2, Xh = X[Id2, ], yh = y[Id2
13    ])
13  PostSigma21[s] <- sig21; PostSigma22[s] <- sig22
14  Beta1 <- PostBeta(sig2h = sig21, Xh = X[Id1, ], yh = y[Id1
15    ])
15  Beta2 <- PostBeta(sig2h = sig22, Xh = X[Id2, ], yh = y[Id2
16    ])
16  PostBetas1[s,] <- Beta1; PostBetas2[s,] <- Beta2
17  Lambda <- sort(MCMCpack::rdirichlet(1, c(a01 + N1, a02 +
18    N2)), decreasing = TRUE)
18  Lambda1 <- Lambda[1]; Lambda2 <- Lambda[2]
19  PostLambda[s] <- Lambda1
20  setTxtProgressBar(pb, s)
21 }
22 close(pb)
23 keep <- seq((burnin+1), tot, thin)
24 PosteriorBetas1 <- coda::mcmc(PostBetas1[keep,])
25 summary(PosteriorBetas1)
26 plot(PosteriorBetas1)
27 PosteriorBetas2 <- coda::mcmc(PostBetas2[keep,])
28 summary(PosteriorBetas2)
29 plot(PosteriorBetas2)
30 PosteriorSigma21 <- coda::mcmc(PostSigma21[keep])
31 summary(PosteriorSigma21)
32 plot(PosteriorSigma21)
33 PosteriorSigma22 <- coda::mcmc(PostSigma22[keep])
34 summary(PosteriorSigma22)
35 plot(PosteriorSigma22)

```

Let's perform another simulation exercise in which we conduct a semi-parametric analysis where the stochastic error follows a Student's t-

distribution with 3 degrees of freedom. Specifically,

$$y_i = 1 - 0.5x_{i1} + 1.5x_{i2} + \mu_i, \quad i = 1, 2, \dots, 500.$$

The variables  $x_{i1}$  and  $x_{i2}$  are standard normally distributed. Let's set  $H = 5$ , and use non-informative priors setting  $\alpha_{h0} = \delta_{h0} = 0.01$ ,  $\beta_0 = \mathbf{0}_2$ ,  $B_0 = I_2$ ,  $\mu_{h0} = 0$ ,  $\sigma_{h0}^2 = 10$  and  $\alpha_0 = [1/H \dots 1/H]^\top$ . Use 6,000 MCMC iterations, burn-in equal to 4,000, and thinning parameter equal to 2. In this exercise, there is no need to address the label-switching issue, as we are not specifically interested in the individual components of the posterior distributions of the clusters. Exercise 1 asks how to get the posterior density of the stochastic errors in semi-parametric specifications.

We can see from the posterior estimates that three components disappear after the burn-in iterations. The 95% credible intervals encompass the population values of the slope parameters. The 95% credible intervals for the probabilities are (0.70, 0.89) and (0.11, 0.30), and the 95% credible interval for the weighted average of the intercepts encompasses the population parameter.

### **R code. Simulation exercise: Semi-parametric model using Gibbs sampling**

```

1 rm(list = ls()); set.seed(010101)
2 library(ggplot2)
3 n <- 500
4 x1 <- rnorm(n); x2 <- rnorm(n)
5 X <- cbind(x1, x2); B <- c(-0.5, 1.5)
6 u <- rt(n, 3); y <- 1 + X %*% B + u
7 Reg <- lm(y ~ X)
8 Res <- Reg$residuals
9 data <- data.frame(Res)
10 # Plot
11 ggplot(data, aes(x = Res)) +
12   geom_density(fill = "blue", alpha = 0.3) + # Density plot
13   labs(title = "Density Plot", x = "Residuals", y = "Density")
14   + theme_minimal()
15 # Hyperparameters
16 d0 <- 0.001; a0 <- 0.001; b0 <- rep(0, 2)
17 B0 <- diag(2); B0i <- solve(B0)
18 mu0 <- 0; sig2mu0 <- 10; H <- 5; a0h <- rep(1/H, H)
19 mcmc <- 2000; burnin <- 4000
20 tot <- mcmc + burnin; thin <- 2

```

**R code. Simulation exercise: Semi-parametric model using Gibbs sampling**

```

1 # Gibbs sampling functions
2 PostSig2 <- function(Beta, muh, Xh, yh){
3   Nh <- length(yh); an <- a0 + Nh
4   dn <- d0 + t(yh - muh - Xh%*%Beta)%*%(yh - muh - Xh%*%Beta
      )
5   sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
6   return(sig2)
7 }
8 PostBeta <- function(sig2, mu, X, y, Psi){
9   XtX <- matrix(0, 2, 2); Xty <- matrix(0, 2, 1)
10  Hs <- length(mu)
11  for(h in 1:Hs){
12    idh <- which(Psi == h)
13    if(length(idh) == 1){
14      Xh <- matrix(X[idh,], 1, 2)
15      XtXh <- sig2[h]^( -1) * t(Xh) %*% Xh
16      yh <- y[idh]
17      Xtyh <- sig2[h]^( -1) * t(Xh) %*% (yh - mu[h])
18    }else{
19      Xh <- X[idh,]
20      XtXh <- sig2[h]^( -1) * t(Xh) %*% Xh
21      yh <- y[idh]
22      Xtyh <- sig2[h]^( -1) * t(Xh) %*% (yh - mu[h])
23    }
24    XtX <- XtX + XtXh; Xty <- Xty + Xtyh
25  }
26  Bn <- solve(B0i + XtX); bn <- Bn%*%(B0i%*%b0 + Xty)
27  Beta <- MASS::mvrnorm(1, bn, Bn)
28  return(Beta)
29 }
30 Postmu <- function(sig2h, Beta, Xh, yh){
31   Nh <- length(yh)
32   sig2mu <- (1/sig2mu0 + Nh/sig2h)^(-1)
33   mun <- sig2mu*(mu0/sig2mu0 + sum((yh - Xh%*%Beta))/sig2h)
34   mu <- rnorm(1, mun, sig2mu^0.5)
35   return(mu)
36 }
37 PostBetas <- matrix(0, mcmc+burnin, 2)
38 PostPsi <- matrix(0, mcmc+burnin, n)
39 PostSigma2 <- list(); PostMu <- list()
40 PostLambda <- list()
41 Resq <- quantile(Res, c(0.2, 0.4, 0.6, 0.8))
42 Id1 <- which(Res <= Resq[1])
43 Id2 <- which(Res > Resq[1] & Res <= Resq[2])
44 Id3 <- which(Res > Resq[2] & Res <= Resq[3])
45 Id4 <- which(Res > Resq[3] & Res <= Resq[4])
46 Id5 <- which(Res > Resq[4])

```

**R code. Simulation exercise: Semi-parametric model using Gibbs sampling**

```

1 Nh <- rep(n/H, H); Lambda <- rep(1/H, H)
2 MU <- c(mean(Res[Id1]), mean(Res[Id2]), mean(Res[Id3]), mean
   (Res[Id4]), mean(Res[Id5]))
3 Sig2 <- c(var(Res[Id1]), var(Res[Id2]), var(Res[Id3]), var(
   Res[Id4]), var(Res[Id5]))
4 Beta <- Reg$coefficients[2:3]
5 Psi <- rep(NA, n); Hs <- length(MU)
6 pb <- txtProgressBar(min = 0, max = tot, style = 3)
7 for(s in 1:tot){
8   for(i in 1:n){
9     lambdai <- NULL
10    for(h in 1:Hs){
11      lambdaih <- Lambda[h]*dnorm(y[i] - X[i,]*Beta, MU[h
12      ], Sig2[h]^0.5)
13      lambdai <- c(lambdai, lambdaih)
14    }
15    Psi[i] <- sample(1:Hs, 1, prob = lambdai)
16  }
17  PostPsi[s, ] <- Psi
18  Hs <- length(table(Psi))
19  for(h in 1:Hs){
20    idh <- which(Psi == h)
21    Sig2[h] <- PostSig2(Beta = Beta, muh = MU[h], Xh = X[idh
22    ,], yh = y[idh])
23    MU[h] <- Postmu(sig2h = Sig2[h], Beta = Beta, Xh = X[idh
24    ,], yh = y[idh])
25  }
26  PostSigma2[[s]] <- Sig2
27  PostMu[[s]] <- MU
28  Beta <- PostBeta(sig2 = Sig2, mu = MU, X = X, y = y, Psi =
   Psi)
29  PostBetas[s,] <- Beta
30  Lambda <- sort(MCMCpack::rdirichlet(1, a0h[1:Hs] + table(
   Psi)), decreasing = TRUE)
31  PostLambda[[s]] <- Lambda
32  setTxtProgressBar(pb, s)
33 }
34 close(pb)
35 keep <- seq(burnin, tot, thin)
36 PosteriorBetas <- coda::mcmc(PostBetas[keep,])
37 summary(PosteriorBetas)
38 plot(PosteriorBetas)
39 PosteriorPsi <- PostPsi[keep,]
40 Clusters <- sapply(1:length(keep), function(i){length(table(
   PosteriorPsi[i,]))})
41 NClus <- 2
42 PosteriorSIGMA <- matrix(NA, length(keep), NClus)
43 PosteriorMU <- matrix(NA, length(keep), NClus)
44 PosteriorLAMBDA <- matrix(NA, length(keep), NClus)

```

**R code. Simulation exercise: Semi-parametric model using Gibbs sampling**

```

1 l <- 1
2 for (s in keep){
3   PosteriorSIGMA[1,] <- PostSigma2[[s]][1:NCLUS]
4   PosteriorMU[1,] <- PostMu[[s]][1:NCLUS]
5   PosteriorLAMBDA[1,] <- PostLambda[[s]][1:NCLUS]
6   l <- l + 1
7 }
8 summary(coda::mcmc(PosteriorSIGMA))
9 summary(coda::mcmc(PosteriorMU))
10 summary(coda::mcmc(PosteriorLAMBDA))

```

### 11.1.2 Dirichlet processes

A Dirichlet process (DP) is a probability distribution over probability distributions, and a generalization of the Dirichlet distribution. It was introduced by [138], and it is commonly used as a prior for unknown distributions, making it particularly useful in non-parametric settings. Unlike finite Gaussian mixture models, a Dirichlet process does not require pre-specifying the number of components, allowing for greater flexibility in modeling complex data structures. In this sense, DPs can be viewed as the limiting case of finite mixtures when the number of components approaches infinity.

A Dirichlet process  $G \sim DP(\alpha G_0)$  is defined by a precision parameter  $\alpha > 0$ , and a base probability measure  $G_0$  on a probability space  $\Omega$ .<sup>1</sup> The DP assigns probability  $G(B)$  to any (measurable) set  $B$  in  $\Omega$  such that for any finite (measurable) partition  $\{B_1, B_2, \dots, B_k\}$  of  $\Omega$ ,

$$G(B_1), G(B_2), \dots, G(B_k) \sim \text{Dirichlet}(\alpha G_0(B_1), \alpha G_0(B_2), \dots, \alpha G_0(B_k)).$$

In particular,  $\mathbb{E}[G(B)] = G_0(B)$  and  $\text{Var}[G(B)] = G_0(B)[1 - G_0(B)]/(1+\alpha)$  [314, p. 8]. Observe that as  $\alpha \rightarrow \infty$ ,  $G$  concentrates at  $G_0$ , which is why  $\alpha$  is called the precision parameter.

A key property of the DP is its discrete nature, which allows it to be expressed as

$$G(\cdot) = \sum_{h=1}^{\infty} \lambda_h \delta_{\theta_h}(\cdot),$$

---

<sup>1</sup>Measurability ensures that probabilities and expectations are well-defined and internally consistent. If a set is not measurable, we cannot assign a meaningful probability to it because doing so would violate the axioms of probability, typically through contradictions with countable additivity or complementation.

where  $\lambda_h$  is the probability mass at  $\boldsymbol{\theta}_h$ , and  $\delta_{\boldsymbol{\theta}_h}(\cdot)$  denotes the Dirac measure that assigns mass one to the atom  $\boldsymbol{\theta}_h$ .<sup>2</sup> Given this property, a particularly useful construction of the DP is the *stick-breaking* representation [407], which is given by

$$\begin{aligned} G(\cdot) &= \sum_{h=1}^{\infty} \lambda_h \delta_{\boldsymbol{\theta}_h}(\cdot), \\ \lambda_h &= V_h \prod_{m < h} (1 - V_m), \quad V_h \sim \text{Beta}(1, \alpha), \\ \boldsymbol{\theta}_h &\stackrel{iid}{\sim} G_0. \end{aligned}$$

The intuition behind this representation is straightforward. We begin with a stick of length 1 and break off a random proportion  $V_1$  drawn from a  $\text{Beta}(1, \alpha)$  distribution. This assigns a probability mass of  $\lambda_1 = V_1$  to  $\boldsymbol{\theta}_1$ , which is sampled from  $G_0$ . Next, from the remaining stick of length  $(1 - V_1)$ , we break off a fraction proportional to  $V_2 \sim \text{Beta}(1, \alpha)$ , assigning a probability mass of  $\lambda_2 = V_2(1 - V_1)$  to a new draw  $\boldsymbol{\theta}_2$  from  $G_0$ . This process continues indefinitely.

Since  $\mathbb{E}(V_h) = \frac{1}{1+\alpha}$ , as  $\alpha \rightarrow \infty$ , we have  $\mathbb{E}(V_h) \rightarrow 0$ . Consequently, the DP places mass on a large number of atoms, leading to convergence to the base distribution  $G_0$ .

Note that DPs give as realizations discrete distributions, this poses challenges when working with continuous distributions. One way to overcome this limitation is to use the DP as a mixing distribution for simple parametric distributions, such as the normal distribution [137]. This leads to Dirichlet process mixtures (DPM), which are defined as

$$f_G(y) = \int f_{\boldsymbol{\theta}}(y) dG(\boldsymbol{\theta}).$$

Observe that the mixing measure  $G$  is discrete when a DP is the prior, with mass concentrated at an infinite number of atoms  $\boldsymbol{\theta}_h$ . Consequently, if  $f_{\boldsymbol{\theta}}(y)$  follows a Gaussian distribution, the resulting mixture resembles a finite Gaussian mixture. However, unlike finite mixtures, the DP-based approach eliminates the need to pre-specify the number of components, as it provides an automatic mechanism for determining them.

Thus, if we assume  $y \sim N(\mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2)$ , then the DPM is

$$p(y_i \mid \{\lambda_h, \boldsymbol{\beta}_h, \sigma_h^2\}_{h=1}^{\infty}, \mathbf{x}_i) = \sum_{h=1}^{\infty} \lambda_h \phi(y_i \mid \mathbf{x}_i^\top \boldsymbol{\beta}_h, \sigma_h^2),$$

where  $\lambda_h$  are drawn from the stick-breaking representation of the DP, and  $\boldsymbol{\theta}_h = [\boldsymbol{\beta}_h^\top \ \sigma_h^2]^\top$ .

---

<sup>2</sup>An atom is a set that has positive probability but cannot be further decomposed into smaller sets with strictly smaller positive probability.

This mixture model can be expressed in a hierarchical structure:

$$\begin{aligned} y_i \mid \boldsymbol{\theta}_i &\stackrel{ind}{\sim} f_{\boldsymbol{\theta}_i} \\ \boldsymbol{\theta}_i \mid G &\stackrel{iid}{\sim} G \\ G \mid \alpha, G_0 &\sim DP(\alpha G_0). \end{aligned}$$

Note that the hierarchical representation induces specific unit parameters, leading to a probabilistic clustering model [21], similar to a finite Gaussian mixture. However, the DPM is not consistent in estimating the number of clusters, it tends to overestimate the number of clusters; although there is posterior asymptotic concentration in the other model components [309].

The hierarchical representation implies that there are latent assignment variables  $s_i = h$ , such that when  $\boldsymbol{\theta}_i$  is equal to the  $h$ -th unique  $\boldsymbol{\theta}_h^*$  — that is,  $\boldsymbol{\theta}_i = \boldsymbol{\theta}_h^*$  — then  $s_i = h$ . Then,

$$\begin{aligned} s_i &\sim \sum_{h=0}^{\infty} \lambda_h \delta_h, \\ y_i \mid s_i, \boldsymbol{\theta}_{s_i} &\sim N(\mathbf{x}_i^\top \boldsymbol{\beta}_{s_i}, \sigma_{s_i}^2), \end{aligned}$$

where  $\lambda_h = P(\boldsymbol{\theta}_i = \boldsymbol{\theta}_h^*)$ .

This latent assignment structure, the *Pólya urn* representation of the DP [48], which is obtained when  $G$  is marginalized out to avoid infinite atoms, and the use of conjugate priors allow for convenient computational inference in DPM.

Specifically, we assume  $\boldsymbol{\beta}_i \mid \sigma_i^2 \sim N(\boldsymbol{\beta}_0, \sigma_i^2 \mathbf{B}_0)$  and  $\sigma_i^2 \sim IG(\alpha_0/2, \delta_0/2)$ . In addition, we can add a layer in the hierarchical representation,  $\alpha \sim G(a, b)$  such that introducing the latent variable  $\xi \mid \alpha, N \sim Be(\alpha + 1, N)$ , allows to easily sample the posterior draws of  $\alpha \mid \xi, H, \pi_\xi \sim \pi_\xi G(a + H, b - \log(\xi)) + (1 - \pi_\xi)G(a + H - 1, b - \log(\xi))$ , where  $\frac{\pi_\xi}{1 - \pi_\xi} = \frac{a + H - 1}{N(b - \log(\xi))}$ ,  $H$  is the number of atoms (mixture components) [137].

The conditional posterior distribution of  $\boldsymbol{\theta}_i$  is

$$\begin{aligned} \boldsymbol{\theta}_i \mid \{\boldsymbol{\theta}_{i'}, \mathbf{s}_{i'} : i' \neq i\}, y_i, \alpha &\sim \sum_{i' \neq i} \frac{N_h^{(i)}}{\alpha + N - 1} f_N(y_i \mid \mathbf{x}_i^\top \boldsymbol{\beta}_h, \sigma_h^2) \\ &\quad + \frac{\alpha}{\alpha + N - 1} \int_{\mathcal{R}^K} \int_0^\infty f_N(y_i \mid \mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2) \\ &\quad \times f_N(\boldsymbol{\beta} \mid \boldsymbol{\beta}_0, \sigma^2 \mathbf{B}_0) f_{IG}(\sigma^2 \mid \alpha_0, \delta_0) d\sigma^2 d\boldsymbol{\beta}, \end{aligned}$$

where  $N_h^{(i)}$  is the number of observations such that  $s_{i'} = h, i' \neq i$ .

Observe that the probability of belonging to a particular cluster has a reinforcement property, as it increases with the cluster size; therefore, a DPM exhibits a self-reinforcing property, the more often a given value has been sampled in the past, the more likely it is to be sampled again.

Observe that the integral in the previous equation has exactly the same form as in the marginal likelihood presented in Section 3.3. Thus,

$$\begin{aligned} p(y_i) &= \int_{\mathcal{R}^K} \int_0^\infty f_N(y_i | \mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2) f_N(\boldsymbol{\beta} | \boldsymbol{\beta}_0, \sigma^2 \mathbf{B}_0) f_{IG}(\sigma^2 | \alpha_0, \delta_0) d\sigma^2 d\boldsymbol{\beta} \\ &= \frac{1}{\pi^{1/2}} \frac{\delta_0^{\alpha_0/2}}{\delta_n^{\alpha_n/2}} \frac{|\mathbf{B}_n|^{1/2} \Gamma(\alpha_n/2)}{|\mathbf{B}_0|^{1/2} \Gamma(\alpha_0/2)}, \end{aligned}$$

where  $\alpha_n = 1 + \alpha_0$ ,  $\delta_n = \delta_0 + y_i^2 + \boldsymbol{\beta}_0^\top \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 - \boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n$ ,  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \mathbf{x}_i \mathbf{x}_i^\top)^{-1}$  and  $\boldsymbol{\beta}_n = \mathbf{B}_n (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{x}_i y_i)$ .

Therefore, we sample  $s_i$  as follows,

$$s_i | \{\boldsymbol{\beta}_{i'}, \sigma_{i'}^2, s_{i'} : i' \neq i\}, y_i, \alpha \sim \left\{ \begin{array}{l} P(s_i = 0 | \cdot) = q_0^* \\ P(s_i = h | \cdot) = q_h^*, h = 1, 2, \dots, H^{(i)} \end{array} \right\},$$

where  $H^{(i)}$  is the number of clusters excluding  $i$ , which may have its own cluster (singleton cluster),  $q_c^* = \frac{q_c}{q_0 + \sum_h q_h}$ ,  $q_c = \{q_0, q_h\}$ ,  $q_h = \frac{N_h^{(i)}}{\alpha + N - 1} f_N(y_i | \mathbf{x}_i^\top \boldsymbol{\beta}_h, \sigma_h^2)$  and  $q_0 = \frac{\alpha}{\alpha + N - 1} p(y_i)$ .

If  $s_i = 0$  is sampled, then  $s_i = H + 1$ , and a new  $\sigma_h^2$  is sampled from  $IG(\alpha_n/2, \delta_n/2)$ , a new  $\boldsymbol{\beta}_h$  is sample from  $N(\boldsymbol{\beta}_n, \sigma_h^2 \mathbf{B}_n)$ .

Discarding  $\boldsymbol{\theta}_h$ 's from last step, we use  $\mathbf{s}$  and the total number of components to sample  $\sigma_h^2$  from

$$IG\left(\frac{\alpha_0 + N_m}{2}, \frac{\delta_0 + \mathbf{y}_h^\top \mathbf{y}_h + \boldsymbol{\beta}_0^\top \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 - \boldsymbol{\beta}_{hn}^\top \mathbf{B}_{hn}^{-1} \boldsymbol{\beta}_{hn}}{2}\right),$$

where  $\mathbf{B}_{hn} = (\mathbf{B}_0^{-1} + \mathbf{X}_h^\top \mathbf{X}_h)^{-1}$  and  $\boldsymbol{\beta}_{hn} = \mathbf{B}_{hn} (\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{X}_h^\top \mathbf{y}_h)$ ,  $\mathbf{X}_h$  and  $\mathbf{y}_h$  have the  $\mathbf{x}_i$  and  $y_i$  of individuals in component  $h$ . We sample  $\boldsymbol{\beta}_h$  from

$$N(\boldsymbol{\beta}_{hn}, \sigma_h^2 \mathbf{B}_{hn}),$$

$$h = 1, 2, \dots, H.$$

We can also use DPMs in a semi-parametric setting, as we did in the finite Gaussian mixtures case. In Exercise 3, we ask to get the posterior sampler in this setting, that is,

$$\begin{aligned} y_i &= \mathbf{x}_i^\top \boldsymbol{\beta} + e_i \\ e_i | \mu_i, \sigma_i^2 &\stackrel{iid}{\sim} N(\mu_i, \sigma_i^2). \end{aligned}$$

We should not include the intercept in  $\boldsymbol{\beta}$ , such that  $\mu_i$  allows to get flexibility in the distribution of the stochastic errors.

Note that mixture models can be easily extended to non-linear models, where posterior inference is based on data augmentation, such as the probit and tobit models in Chapter 6. The basic idea is to incorporate the mixture (finite or infinite) into the specification of the latent variable implicit in the

data-generating process. For instance, [29] perform inference in the probit model using DPMs. Furthermore, mixtures can also be used in the multivariate models presented in Chapter 7 and the hierarchical models presented in Chapter 9. [360] perform semi-parametric inference in the exact affine demand system [266] using DPMs, while [29] apply them in hierarchical models.

To sum up, a Dirichlet process mixture is a flexible way to model non-parametrically a distribution using an infinite weighted sum of discrete distributions, where each individual weight increases with the number of observations that belongs to it.

### Example: Simulation exercises

Let's simulate again the simple regression mixture with two components such that  $\psi_{i1} \sim \text{Ber}(0.5)$ , consequently,  $\psi_{i2} = 1 - \psi_{i1}$ , and assume one regressor,  $x_i \sim N(0, 1)$ ,  $i = 1, 2, \dots, 1,000$ . Then,

$$p(y_i | \mathbf{x}_i) = 0.5\phi(y_i | 2 + 1.5x_i, 1^2) + 0.5\phi(y_i | -1 + 0.5x_i, 0.8^2).$$

We use non-informative priors again, setting  $\alpha_0 = \delta_0 = 0.01$ ,  $\beta_0 = \mathbf{0}_2$ ,  $B_0 = I_2$ , and  $a = b = 0.1$ . The number of MCMC iterations is 5,000, the burn-in is 1,000, and the thinning parameter is 2. However, we do not have to fix in advance the number of clusters using the DPM.

The following code demonstrates how to perform inference using the stated Gibbs sampler for the DPM in the **R** package.

We see from the results that the posterior plots illustrate the label-switching problem. Eventually, the posterior chains of clusters reach different posterior modes, causing a high variability of the posterior draws. In Exercise 5, we ask to fix this issue using *random permutation of latent classes*.

### R code. Simulation exercise: Dirichlet process mixture using Gibbs sampling

```

1 rm(list = ls()); set.seed(010101)
2 # Simulate data from a 2-component mixture model
3 N <- 1000; x <- rnorm(N); z <- rbinom(N, 1, 0.5)
4 y <- ifelse(z == 1, rnorm(N, 2 + 1.5*x, 1), rnorm(N, -1 +
  0.5*x, 0.8))
5 X <- cbind(1, x); k <- 2
6 data <- data.frame(y, x); Reg <- lm(y ~ x); SumReg <-
  summary(Reg)

```

**R code. Simulation exercise: Dirichlet process mixture using Gibbs sampling**

```

1 # Hyperparameters
2 a0 <- 0.001; d0 <- 0.001
3 b0 <- rep(0, k); B0 <- diag(k)
4 B0i <- solve(B0)
5 a <- 0.1; b <- 0.1
6 # MCMC parameters
7 mcmc <- 5000; burnin <- 1000
8 tot <- mcmc + burnin; thin <- 2
9 # Gibbs sampling functions
10 PostSig2 <- function(Xh, yh){
11   Nh <- length(yh)
12   yh <- matrix(yh, Nh, 1)
13   if(Nh == 1){
14     Xh <- matrix(Xh, k, 1)
15     Bn <- solve(Xh%*%t(Xh) + B0i)
16     bn <- Bn%*%(B0i%*%b0 + Xh%*%yh)
17   }else{
18     Xh <- matrix(Xh, Nh, k)
19     Bn <- solve(t(Xh)%*%Xh + B0i)
20     bn <- Bn%*%(B0i%*%b0 + t(Xh)%*%yh)
21   }
22   Bni <- solve(Bn)
23   an <- a0 + Nh
24   dn <- d0 + t(yh)%*%yh + t(b0)%*%B0i%*%b0 - t(bn)%*%Bni%*%
25   bn
26   sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
27 }
28 PostBeta <- function(sig2h, Xh, yh){
29   Nh <- length(yh)
30   yh <- matrix(yh, Nh, 1)
31   if(Nh == 1){
32     Xh <- matrix(Xh, k, 1)
33     Bn <- solve(Xh%*%t(Xh) + B0i)
34     bn <- Bn%*%(B0i%*%b0 + Xh%*%yh)
35   }else{
36     Xh <- matrix(Xh, Nh, k)
37     Bn <- solve(t(Xh)%*%Xh + B0i)
38     bn <- Bn%*%(B0i%*%b0 + t(Xh)%*%yh)
39   }
40   Beta <- MASS::mvrnorm(1, bn, sig2h*Bn)
41   return(Beta)
42 }
```

**R code. Simulation exercise: Dirichlet process mixture using Gibbs sampling**

```

1 PostAlpha <- function(s, alpha){
2   H <- length(unique(s))
3   psi <- rbeta(1, alpha + 1, N)
4   pi.ratio <- (a + H - 1) / (N * (b - log(psi)))
5   pi <- pi.ratio / (1 + pi.ratio)
6   components <- sample(1:2, prob = c(pi, (1 - pi)), size =
7     1)
8   cs <- c(a + H, a + H - 1)
9   ds <- b - log(psi)
10  alpha <- rgamma(1, cs[components], ds)
11  return(alpha)
12 }
13 LogMarLikLM <- function(xh, yh){
14   xh <- matrix(xh, k, 1)
15   Bn <- solve(xh%*%t(xh) + B0i)
16   Bni <- solve(Bn)
17   bn <- Bn%*%(B0i%*%b0 + xh%*%yh)
18   an <- a0 + 1
19   dn <- d0 + yh^2 + t(b0)%*%B0i%*%b0 - t(bn)%*%Bni%*%bn
20   # Log marginal likelihood
21   logpy <- (1/2)*log(1/pi)+(a0/2)*log(d0)-(an/2)*log(dn) +
22   0.5*log(det(Bn)/det(B0)) + lgamma(an/2)-lgamma(a0/2)
23   return(logpy)
24 }
25 PostS <- function(BETA, SIGMA, Alpha, s, i){
26   Nl <- table(s[-i]); H <- length(Nl)
27   qh <- sapply(1:H, function(h){(Nl[h]/(N+Alpha-1))*dnorm(y[
28     i], mean = t(X[i,])%*%BETA[,h], sd = SIGMA[h])})
29   q0 <- (Alpha/(N+Alpha-1))*exp(LogMarLikLM(xh = X[i,], yh =
30     y[i]))
31   qh <- c(q0, qh)
32   Clust <- as.numeric(names(Nl))
33   si <- sample(c(0, Clust), 1, prob = qh)
34   if(si == 0){
35     si <- Clust[H] + 1
36     Sig2New <- PostSig2(Xh = X[i,], yh = y[i])
37     SIGMA <- c(SIGMA, Sig2New^0.5)
38     BetaNew <- PostBeta(sig2h = Sig2New, Xh = X[i,], yh =
39     y[i])
40     BETA <- cbind(BETA, BetaNew)
41   }else{si == si}
42   return(list(si = si, BETA = BETA, SIGMA = SIGMA))
43 }
44 PostBetas <- list(); PostSigma <- list()
45 Posts <- matrix(0, tot, N); PostAlphas <- rep(0, tot)
46 S <- sample(1:3, N, replace = T, prob = c(0.5, 0.3, 0.2))

```

**R code. Simulation exercise: Dirichlet process mixture using Gibbs sampling**

```

1 BETA <- cbind(Reg$coefficients, Reg$coefficients, Reg$  

2 coefficients)
3 SIGMA <- rep(SumReg$sigma, 3); Alpha <- rgamma(1, a, b)
4 pb <- txtProgressBar(min = 0, max = tot, style = 3)
5 for(s in 1:tot){
6   for(i in 1:N){
7     Rests <- PostS(BETA = BETA, SIGMA = SIGMA, Alpha = Alpha  

8       , s = S, i = i)
9     S[i] <- Rests$si
10    BETA <- Rests$BETA; SIGMA <- Rests$SIGMA
11  }
12  sFreq <- table(S)
13  lt <- 1
14  for(li in as.numeric(names(sFreq))){
15    Index <- which(S == li)
16    if(li == lt){S[Index] <- li
17    } else {S[Index] <- lt
18    }
19    lt <- lt + 1
20  }
21  Alpha <- PostAlpha(s = S, alpha = Alpha)
22  Nl <- table(S); H <- length(Nl)
23  SIGMA <- rep(NA, H)
24  BETA <- matrix(NA, k, H)
25  l <- 1
26  for(h in unique(S)){
27    Idh <- which(S == h)
28    SIGMA[1] <- (PostSig2(Xh = X[Idh, ], yh = y[Idh]))^0.5
29    BETA[,1] <- PostBeta(sig2h = SIGMA[1]^2, Xh = X[Idh, ],  

30      yh = y[Idh])
31    l <- l + 1
32  }
33  PostBetas[[s]] <- BETA
34  PostSigma[[s]] <- SIGMA
35  Posts[s, ] <- S
36  PostAlphas[s] <- Alpha
37  setTxtProgressBar(pb, s)
38 }
39 close(pb)
40 keep <- seq((burnin+1), tot, thin)
41 PosteriorS<- Posts[keep,]
42 Clusters <- sapply(1:length(keep), function(i){length(table(  

43   PosteriorS[i,]))})
44 table(Clusters)
45 PosteriorBeta1 <- matrix(NA, length(keep), k); j <- 1
46 for(s in keep){
47   PosteriorBeta1[j,] <- PostBetas[[s]][,1]; j <- j + 1
48 }
49 print(summary(coda::mcmc(PosteriorBeta1)))
50 plot(coda::mcmc(PosteriorBeta1))

```

### R code. Simulation exercise: Dirichlet process mixture using Gibbs sampling

```

1 PosteriorBeta2 <- matrix(NA, length(keep), k); j <- 1
2 for(s in keep){
3   PosteriorBeta2[j,] <- PostBetas[[s]][,2]; j <- j + 1
4 }
5 print(summary(coda::mcmc(PosteriorBeta2)))
6 plot(coda::mcmc(PosteriorBeta2))
7 # # If there are 3 clusters
8 # PosteriorBeta3 <- matrix(NA, length(keep), k); j <- 1
9 # for(s in keep){
10  # PosteriorBeta3[j,] <- PostBetas[[s]][,3]; j <- j + 1
11  # }
12 # print(summary(coda::mcmc(PosteriorBeta3)))
13 # plot(coda::mcmc(PosteriorBeta3))
14

```

#### Example: Consumption of marijuana in Colombia

Let's use the dataset *MarijuanaColombia.csv* from our GitHub repository to perform inference on the demand for marijuana in Colombia. This dataset contains information on the (log) monthly demand in 2019 from the National Survey of the Consumption of Psychoactive Substances. It includes variables such as the presence of a drug dealer in the neighborhood (*Dealer*), gender (*Female*), indicators of good physical and mental health (*PhysicalHealthGood* and *MentalHealthGood*), age (*Age* and *Age2*), years of schooling (*YearsEducation*), and (log) prices of marijuana, cocaine, and crack by individual (*LogPriceMarijuana*, *LogPriceCocaine*, and *LogPriceCrack*). The sample size is 1,156.

We are interested in the own-price and cross-price elasticities of marijuana demand. However, there is a potential endogeneity issue between price and demand due to strategic provider search (omission of a relevant regressor) or measurement errors in self-reported prices. Endogeneity should be properly addressed for rigorous scientific analysis. This example is purely illustrative.

In Exercise 2 we ask to perform inference using a finite Gaussian mixture regression starting with five potential clusters. Here we perform inference using a Dirichlet process mixture.

The following code shows how to perform this analysis using non-informative priors, setting  $\alpha_0 = \delta_0 = 0.01$ ,  $\beta_0 = \mathbf{0}_K$ ,  $B_0 = I_K$ , and  $a = b = 0.1$ ,  $K$  is the number of regressors, 11 including the intercept. The number of MCMC iterations is 5,000, the burn-in is 1,000, and the thinning parameter is 2.

***R code. Demand of marijuana in Colombia:  
Dirichlet process mixture***

---

```

1 rm(list = ls()); set.seed(010101)
2 Data <- read.csv("https://raw.githubusercontent.com/
  BEsmarter-consultancy/BSTApp/reviews/heads/master/DataApp/
  MarijuanaColombia.csv")
3 attach(Data)
4 y <- LogMarijuana; X <- as.matrix(cbind(1, Data[,-1]))
5 Reg <- lm(y ~ X - 1); SumReg <- summary(Reg)
6 k <- dim(X)[2]
7 N <- dim(X)[1]
8 library(ggplot2)
9 ggplot(Data, aes(x = LogMarijuana)) + geom_density(fill =
  "blue", alpha = 0.3) + labs(title = "Density Plot:
  Marijuana (log) monthly consumption in Colombia", x = "y
  ", y = "Density") + theme_minimal()
10 a0 <- 0.001; d0 <- 0.001
11 b0 <- rep(0, k); B0 <- diag(k)
12 B0i <- solve(B0)
13 a <- 0.1; b <- 0.1
14 # MCMC parameters
15 mcmc <- 5000; burnin <- 1000
16 tot <- mcmc + burnin; thin <- 2
17 # Gibbs sampling functions
18 PostSig2 <- function(Xh, yh){
19   Nh <- length(yh)
20   yh <- matrix(yh, Nh, 1)
21   if(Nh == 1){
22     Xh <- matrix(Xh, k, 1)
23     Bn <- solve(Xh%*%t(Xh) + B0i)
24     bn <- Bn%*%(B0i%*%b0 + Xh%*%yh)
25   }else{
26     Xh <- matrix(Xh, Nh, k)
27     Bn <- solve(t(Xh)%*%Xh + B0i)
28     bn <- Bn%*%(B0i%*%b0 + t(Xh)%*%yh)
29   }
30   Bni <- solve(Bn)
31   an <- a0 + Nh
32   dn <- d0 + t(yh)%*%yh + t(b0)%*%B0i%*%b0 - t(bn)%*%Bni%*%
33   bn
34   sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
35   return(sig2)
36 }
37 PostBeta <- function(sig2h, Xh, yh){
38   Nh <- length(yh)
39   yh <- matrix(yh, Nh, 1)
40   if(Nh == 1){
41     Xh <- matrix(Xh, k, 1)
42     Bn <- solve(Xh%*%t(Xh) + B0i)

```

---

**R code. Demand of marijuana in Colombia:  
Dirichlet process mixture**

```

1 }else{
2   Xh <- matrix(Xh, Nh, k)
3   Bn <- solve(t(Xh) %*% Xh + B0i)
4   bn <- Bn %*% (B0i %*% b0 + t(Xh) %*% yh)
5 }
6 Beta <- MASS::mvrnorm(1, bn, sig2h*Bn)
7 return(Beta)
8 }
9 PostAlpha <- function(s, alpha){
10   H <- length(unique(s))
11   psi <- rbeta(1, alpha + 1, N)
12   pi.ratio <- (a + H - 1) / (N * (b - log(psi)))
13   pi <- pi.ratio / (1 + pi.ratio)
14   components <- sample(1:2, prob = c(pi, (1 - pi)), size =
15     1)
16   cs <- c(a + H, a + H - 1)
17   ds <- b - log(psi)
18   alpha <- rgamma(1, cs[components], ds)
19   return(alpha)
20 }
21 LogMarLikLM <- function(xh, yh){
22   xh <- matrix(xh, k, 1)
23   Bn <- solve(xh %*% t(xh) + B0i)
24   Bni <- solve(Bn)
25   bn <- Bn %*% (B0i %*% b0 + xh %*% yh)
26   an <- a0 + 1
27   dn <- d0 + yh^2 + t(b0) %*% B0i %*% b0 - t(bn) %*% Bni %*% bn
28   logpy <- (1/2)*log(1/pi)+(a0/2)*log(d0)-(an/2)*log(dn) +
29   0.5*log(det(Bn)/det(B0)) + lgamma(an/2)-lgamma(a0/2)
30   return(logpy)
31 }
32 Posts <- function(BETA, SIGMA, Alpha, s, i){
33   Nl <- table(s[-i]); H <- length(Nl)
34   qh <- sapply(1:H, function(h){(Nl[h]/(N+Alpha-1))*dnorm(y[
35     i], mean = t(X[i,]) %*% BETA[,h], sd = SIGMA[h])})
36   q0 <- (Alpha/(N+Alpha-1))*exp(LogMarLikLM(xh = X[i,], yh =
37     y[i]))
38   qh <- c(q0, qh)
39   Clust <- as.numeric(names(Nl))
40   si <- sample(c(0, Clust), 1, prob = qh)
41   if(si == 0){
42     si <- Clust[H] + 1
43     Sig2New <- PostSig2(Xh = X[i,], yh = y[i])
44     SIGMA <- c(SIGMA, Sig2New^0.5)
45     BetaNew <- PostBeta(sig2h = Sig2New, Xh = X[i,], yh = y[
46       i])
47     BETA <- cbind(BETA, BetaNew)
48   }else{si == si}
49 }
50 return(list(si = si, BETA = BETA, SIGMA = SIGMA))
51 }
```

***R code. Demand of marijuana in Colombia:  
Dirichlet process mixture***

---

```

1 PostBetas <- list(); PostSigma <- list()
2 Posts <- matrix(0, tot, N); PostAlphas <- rep(0, tot)
3 S <- sample(1:3, N, replace = T, prob = c(0.5, 0.3, 0.2))
4 BETA <- cbind(Reg$coefficients, Reg$coefficients, Reg$coefficients)
5 SIGMA <- rep(SumReg$sigma, 3)
6 Alpha <- rgamma(1, a, b)
7 pb <- txtProgressBar(min = 0, max = tot, style = 3)
8 for(s in 1:tot){
9   for(i in 1:N){
10     Rests <- PostS(BETA = BETA, SIGMA = SIGMA, Alpha = Alpha,
11       , s = S, i = i)
12     S[i] <- Rests$si
13     BETA <- Rests$BETA; SIGMA <- Rests$SIGMA
14   }
15   sFreq <- table(S)
16   lt <- 1
17   for(li in as.numeric(names(sFreq))){
18     Index <- which(S == li)
19     if(li == lt){S[Index] <- li
20     } else {S[Index] <- lt
21     }
22     lt <- lt + 1
23   }
24   Alpha <- PostAlpha(s = S, alpha = Alpha)
25   N1 <- table(S); H <- length(N1)
26   SIGMA <- rep(NA, H)
27   BETA <- matrix(NA, k, H)
28   l <- 1
29   for(h in unique(S)){
30     Idh <- which(S == h)
31     SIGMA[1] <- (PostSig2(Xh = X[Idh, ], yh = y[Idh]))^0.5
32     BETA[,1] <- PostBeta(sig2h = SIGMA[1]^2, Xh = X[Idh, ],
33       yh = y[Idh])
34     l <- l + 1
35   }
36   PostBetas[[s]] <- BETA
37   PostSigma[[s]] <- SIGMA
38   Posts[s, ] <- S
39   PostAlphas[s] <- Alpha
40   setTxtProgressBar(pb, s)
41 }
42 close(pb)
43 keep <- seq((burnin+1), tot, thin)
44 PosteriorS<- Posts[keep,]
45 Clusters <- sapply(1:length(keep), function(i){length(table(PosteriorS[i,]))})
46 table(Clusters)

```

---

***R code. Demand of marijuana in Colombia:  
Dirichlet process mixture***

---

```

1 NClus <- 3
2 sapply(1:length(keep), function(i){print(table(PosteriorS[i
    ,]))})
3 PosteriorSIGMA <- matrix(NA, length(keep), NClus)
4 l <- 1
5 for (s in keep){
6   PosteriorSIGMA[1,] <- PostSigma[[s]][1:NClus]
7   l <- l + 1
8 }
9 summary(coda::mcmc(PosteriorSIGMA))
10 plot(coda::mcmc(PosteriorSIGMA))
11 l <- 1
12 for (s in keep){
13   PosteriorSIGMA[1,] <- PostSigma[[s]][1:NClus]
14   l <- l + 1
15 }
16 summary(coda::mcmc(PosteriorSIGMA))
17 plot(coda::mcmc(PosteriorSIGMA))
18 PosteriorBeta1 <- matrix(NA, length(keep), k)
19 j <- 1
20 for(s in keep){
21   PosteriorBeta1[j,] <- PostBetas[[s]][,1]
22   j <- j + 1
23 }
24 colnames(PosteriorBeta1) <- c("Ct", names(Data[, -1]))
25 HDI <- HDInterval::hdi(PosteriorBeta1, credMass = 0.95)
26 print(summary(coda::mcmc(PosteriorBeta1)))
27 plot(coda::mcmc(PosteriorBeta1))
28 PosteriorBeta2 <- matrix(NA, length(keep), k)
29 j <- 1
30 for(s in keep){
31   PosteriorBeta2[j,] <- PostBetas[[s]][,2]
32   j <- j + 1
33 }
34 colnames(PosteriorBeta2) <- c("Ct", names(Data[, -1]))
35 HDI <- HDInterval::hdi(PosteriorBeta2, credMass = 0.95)
36 print(summary(coda::mcmc(PosteriorBeta2)))
37 plot(coda::mcmc(PosteriorBeta2))
38 PosteriorBeta3 <- matrix(NA, length(keep), k)
39 j <- 1
40 for(s in keep){
41   PosteriorBeta3[j,] <- PostBetas[[s]][,3]
42   j <- j + 1
43 }
44 colnames(PosteriorBeta3) <- c("Ct", names(Data[, -1]))
45 HDI <- HDInterval::hdi(PosteriorBeta3, credMass = 0.95)
46 print(summary(coda::mcmc(PosteriorBeta3)))
47 plot(coda::mcmc(PosteriorBeta3))

```

---

***R code. Demand of marijuana in Colombia:  
Dirichlet process mixture***

```

1 DataElast <- data.frame(Marijuana = PosteriorBeta1[,2],
  Cocaine = PosteriorBeta1[,3],
  Crack = PosteriorBeta1[,4])
2 dens1 <- ggplot(DataElast, aes(x = Marijuana)) + geom_
  density(fill = "blue", alpha = 0.3) + labs(x = "
  Marijuana", y = "Density") + theme_minimal()
3 dens2 <- ggplot(DataElast, aes(x = Cocaine)) + geom_density(
  fill = "blue", alpha = 0.3) + labs(x = "Cocaine", y = "
  Density") + theme_minimal()
4 dens3 <- ggplot(DataElast, aes(x = Crack)) + geom_density(
  fill = "blue", alpha = 0.3) +
5 labs(x = "Crack", y = "Density") + theme_minimal()
6 library(ggpubr)
7 ggarrange(dens1, dens2, dens3, labels = c("A", "B", "C"),
  ncol = 3, nrow = 1, legend = "bottom", common.legend =
  TRUE)
8
9

```

In this application, there are three clusters after the burn-in period, with individuals in the sample relatively evenly allocated between them. The posterior density plots of the elasticities associated with the three clusters appear very similar, which implies that the posterior draws associated with each cluster are targeting the same posterior distribution. Consequently, Figure 11.4 presents the posterior density estimates of the demand price elasticities of marijuana in Colombia for the first cluster.

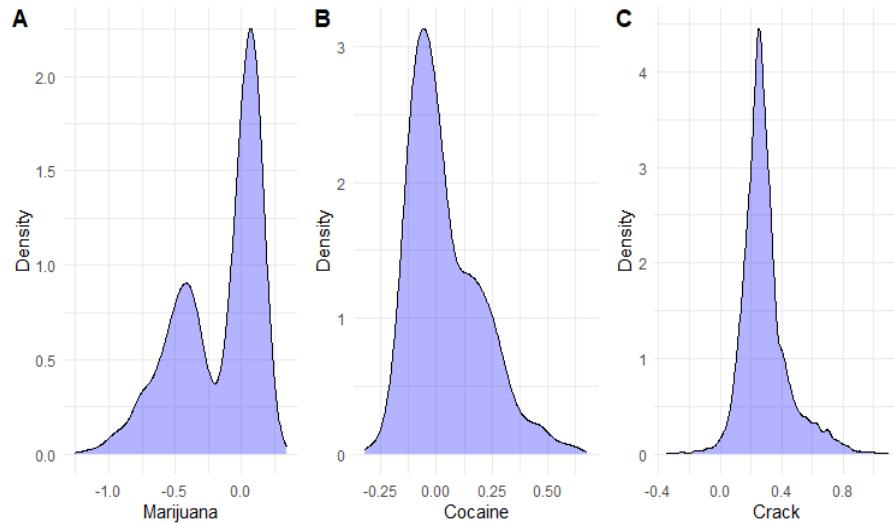
Panel A of Figure 11.4 shows that the own-price elasticity of marijuana demand is bimodal. The first mode suggests that some individuals have an elasticity near -0.5, meaning that a 10% increase in the price of marijuana leads to a 5% decrease in their demand. The second mode is near 0, indicating that these individuals do not adjust their marijuana consumption in response to price changes.

Panel B presents the cross-price elasticity of marijuana demand with respect to the price of cocaine. This posterior distribution is also bimodal. The first mode suggests that some individuals do not change their marijuana demand in response to cocaine price variations. The second mode indicates that marijuana and cocaine are substitutes: a 10% increase in the price of cocaine leads to an approximately 2.5% increase in marijuana demand.

Panel C shows that marijuana and crack are also substitutes. This posterior distribution is unimodal, indicating that most individuals increase their marijuana consumption when the price of crack rises. On average, a 10% increase

in the price of crack results in an approximately 3% increase in marijuana demand.

These results seem plausible. However, they should be interpreted with caution, as potential endogeneity issues may affect the estimates.



**FIGURE 11.4**

Posterior distributions: Elasticities of marijuana consumption in Colombia.

## 11.2 Splines

Another common approach to performing non-parametric or semi-parametric inference is using splines. These are piecewise polynomial functions that allow for approximating complex relationships in data.

To clarify ideas, let's begin with a simple univariate case. The starting point is to assume that

$$y_i = f(x_i) + \mu_i,$$

where  $\mu_i \sim N(0, \sigma^2)$  is independent of the regressor  $x_i$ , and  $f(x_i)$  is a smooth function such that nearby points in the support of  $x_i$  should have similar values. Observe that this means we must always sort the data in increasing order according to  $x_i$ . This has no effect, as we assume a random sample. The index  $i = 1, 2, \dots, N$  represents the observations.

The spline approach defines

$$f(x_i) = \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \cdots + \beta_H b_H(x_i) = \mathbf{b}(x_i)^\top \boldsymbol{\beta}, \quad (11.1)$$

where  $\beta = [\beta_1 \ \beta_2 \ \dots \ \beta_H]^\top$  and  $\mathbf{b}(x_i) = [b_1(x_i) \ b_2(x_i) \ \dots \ b_H(x_i)]^\top$  are the vector of parameters and the vector of spline basis (B-splines) function values at  $x_i$ , respectively.

Note that in this setting, there is only one regressor; however, there are  $H$  location parameters. Given  $N$  observations, the design matrix is:

$$\mathbf{W} = \begin{bmatrix} b_1(x_1) & b_2(x_1) & \dots & b_H(x_1) \\ b_1(x_2) & b_2(x_2) & \dots & b_H(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ b_1(x_N) & b_2(x_N) & \dots & b_H(x_N) \end{bmatrix}.$$

Thus, we have a standard linear model where we can proceed as in Chapter 3.

Actually, Equation 11.1 represents the basis function approach, where  $\mathbf{b}(x_i)$  is a specified set of basis functions, such as B-splines. However, other functions derived from the Fourier transform, Gaussian radial basis functions, wavelets, etc., can also be used. Here, we focus on the most popular case: cubic B-splines. However, the same ideas apply to other basis functions. In particular, cubic B-splines provide the lowest degree necessary to generate sufficiently smooth curves; the first and second derivatives are nonzero, which is not the case with constant and linear splines [299].

Thus, from Equation 11.1, we see that a spline is a linear combination of  $H$  B-splines. The latter are polynomials of a certain degree, meaning that splines are piecewise polynomial functions. Any spline function of order  $n$  can be expressed as a linear combination of B-splines of order  $n$ , and B-splines serve as the basis functions for the spline function space. The cubic B-spline is

$$b_{h,3}(x_i) = \begin{cases} \frac{u^3}{6} & x_h \leq x_i < x_{h+1}, \quad u = (x_i - x_h)/\delta \\ \frac{1+3u+3u^2-3u^3}{6} & x_{h+1} \leq x_i < x_{h+2}, \quad u = (x_i - x_{h+1})/\delta \\ \frac{4-6u^2+3u^3}{6} & x_{h+2} \leq x_i < x_{h+3}, \quad u = (x_i - x_{h+2})/\delta \\ \frac{1-3u+3u^2-3u^3}{6} & x_{h+3} \leq x_i < x_{h+4}, \quad u = (x_i - x_{h+3})/\delta \\ 0 & \text{otherwise} \end{cases},$$

where  $x_{h+k}$  are called the knots, for  $k = 0, 1, \dots, 4$ . These are the points in the domain of the function where the pieces of the spline join, satisfying  $x_{h+k} \leq x_{h+k+1}$ . The knots control the shape and flexibility of the spline curve; more knots imply greater flexibility but less smoothness. Note that  $\delta$  determines the distance between the knots, and that this expression clearly shows that the cubic B-spline is a piecewise polynomial function.

The following code implements this function in **R** using a delta of 1.5 and 5 knots  $\{2, 3.5, 5, 6.5, 8\}$ . The knots  $\{2, 8\}$  are called boundary knots, which can be given by the range of  $x$ , and the knots  $\{3.5, 5, 6.5\}$  are called internal knots. B-splines with evenly distributed knots are called uniform B-splines. However, this is not always the case.

Figure 11.5 compares the results of our own function (black circles) with those obtained using the *bs* function from the *splines* package (red line).

### R code. B-spline

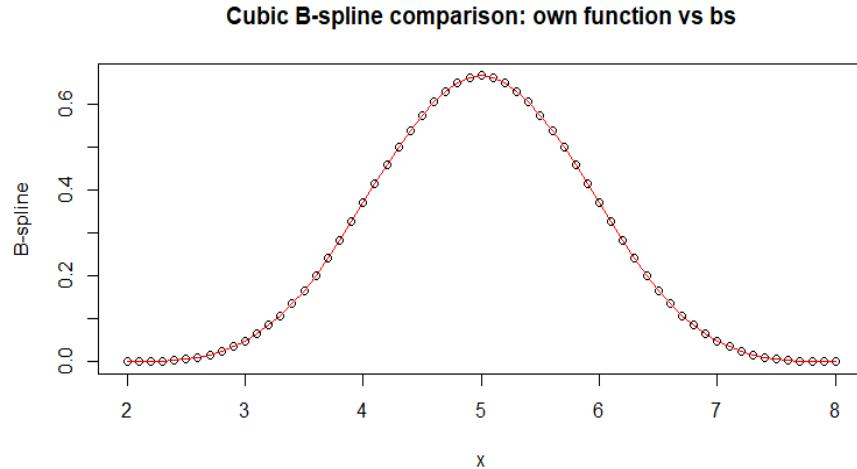
```

1 SplineOwn <- function(x, knots, delta){
2   if(knots[1] <= x & x < knots[2]){
3     u <- (x - knots[1])/delta
4     b <- u^3/6
5   }else{
6     if(knots[2] <= x & x < knots[3]){
7       u <- (x - knots[2])/delta
8       b <- (1/6)*(1 + 3*u + 3*u^2 - 3*u^3)
9     }else{
10      if(knots[3] <= x & x < knots[4]){
11        u <- (x - knots[3])/delta
12        b <- (1/6)*(4 - 6*u^2 + 3*u^3)
13      }else{
14        if(knots[4] <= x & x < knots[5]){
15          u <- (x - knots[4])/delta
16          b <- (1/6)*(1 - 3*u + 3*u^2 - u^3)
17        }else{
18          b <- 0
19        }
20      }
21    }
22  }
23  return(b)
24 }
25 delta <- 1.5
26 knotsA <- seq(2, 8, delta)
27 xA <- seq(2, 8, 0.1)
28 Ens <- sapply(xA, function(xi) {SplineOwn(xi, knots = knotsA,
29 , delta = delta)})
30 plot(xA, Ens, xlab = "x", ylab = "B-spline", main = "Cubic B
31 -spline comparison: own function vs bs")
32 require(splines)
33 BSfunc <- bs(xA, knots = knotsA, degree = 3)
34 lines(xA, BSfunc[,4], col = "red")

```

Figure 11.5 shows a particular cubic B-spline, but  $f(x_i)$  in Equation 11.1 involves  $H$  splines over the range of potential values of  $x$ , as we need to evaluate the splines centered at different knots. Here, we should clarify that to avoid issues with the use of splines at boundary regions, it is essential to artificially increase the number of boundary knots by repeating these knots as many times as the degree of the polynomial. For instance, in our example, the final set of knots is  $\{2, 2, 2, 2, 3.5, 5, 6.5, 8, 8, 8, 8\}$ .

The set of B-splines can be constructed using the Cox-de Boor recursion



**FIGURE 11.5**  
Cubic B-spline comparison.

formula. The starting point is the B-spline of degree 0,

$$b_{h,0}(x_i) = \begin{cases} 1 & x_h \leq x_i < x_{h+1} \\ 0 & \text{otherwise} \end{cases},$$

and then, the other B-splines are defined by the recursion

$$b_{h,p}(x_i) = \frac{x_i - x_h}{x_{h+p} - x_h} b_{h,p-1}(x_i) + \frac{x_{h+p+1} - x_i}{x_{h+p+1} - x_{h+1}} b_{h+1,p-1}(x_i). \quad (11.2)$$

Note that B-splines are defined by their degree and their knots.

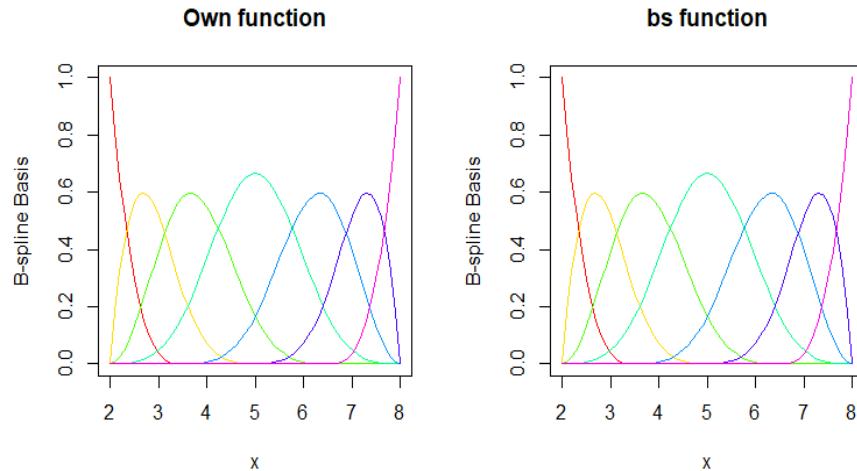
The following code demonstrates how to perform this recursion using the same settings as in the previous code and compares the results with the *bs* function from the *splines* package. As shown in Figure 11.6, we obtain the same results. We will continue using the *bs* function from the *splines* package for convenience in the process, but it seems that we have gained some understanding of the underlying process.

### R code. Splines

```

1 cubic_bspline <- function(x, knots, degree = 3) {
2   extended_knots <- c(rep(knots[1], degree), knots, rep(
3     knots[length(knots)], degree))
4   num_basis <- length(knots) + degree - 1
5   basis_matrix <- matrix(0, nrow = length(x), ncol = num_
6     basis)
7   # Function to compute B-spline basis recursively
8   b_spline_basis <- function(x, degree, i, knots) {
9     if (degree == 0) {
10       return(ifelse(x >= knots[i] & x < knots[i + 1], 1, 0))
11     } else {
12       left_num <- (x - knots[i])
13       left_den <- (knots[i + degree] - knots[i])
14       left <- ifelse(left_den != 0, (left_num / left_den) *
15         b_spline_basis(x, degree - 1, i, knots), 0)
16
17       right_num <- (knots[i + degree + 1] - x)
18       right_den <- (knots[i + degree + 1] - knots[i + 1])
19       right <- ifelse(right_den != 0, (right_num / right_den) *
20         b_spline_basis(x, degree - 1, i + 1, knots), 0)
21
22       return(left + right)
23     }
24   }
25   for (i in 1:num_basis) {
26     basis_matrix[, i] <- sapply(x, function(xi) b_spline_
27       basis(xi, degree, i, extended_knots))
28   }
29   if(x[length(x)] == knots[length(knots)]){
30     basis_matrix[length(x), num_basis] <- 1
31   }
32   return(basis_matrix)
33 }
34 delta <- 1.5
35 knotsA <- seq(2, 8, delta)
36 xA <- seq(2, 8, 0.1)
37 basis_matrix <- cubic_bspline(xA, knots = knotsA, degree =
38   3)
39 library(splines)
40 bs_matrix <- bs(xA, knots = knotsA[-c(1, length(knotsA))],
41   degree = 3, intercept = TRUE, Boundary.knots = range(
42     knotsA))
43 bs_matrix_matrix <- as.matrix(bs_matrix)
44 par(mfrow = c(1,2))
45 matplot(xA, basis_matrix, type = "l", lty = 1, col = rainbow
46   (ncol(basis_matrix)), ylab = "B-spline Basis", xlab = "x"
47   , main = "Own function")
48 matplot(xA, bs_matrix_matrix, type = "l", lty = 1, col =
49   rainbow(ncol(basis_matrix)), ylab = "B-spline Basis",
50   xlab = "x", main = "bs function")

```

**FIGURE 11.6**

Sequence of cubic B-splines comparison.

The idea in splines is to calculate the linear combination of the  $H$  B-splines, as those shown in Figure 11.6, that bets fit the dependent variable.

Splines can also be extended to non-linear models based on data augmentation, such as the probit and tobit models. Again, the basic idea is to incorporate splines into the implicit latent variables. [250, Chap. 10] presents these extensions.

Let's perform a simulation exercise to fix ideas.

#### Simulation exercise: Splines

Let's assume that the data generating process is

$$y_i = 0.4 + 0.25 \sin(8x_i - 5) + 0.4 \exp(-16(4x_i - 2.5)^2) + \mu_i,$$

where  $\mu_i \sim N(0, 0.15^2)$ , and  $x_i$  is a sequence in  $(0, 1)$ , with 100 random draws of the pairs  $(x_i, y_i)$ . In addition, we choose the knots  $\{0, 0.25, 0.5, 0.75, 1\}$ . We then calculate the cubic B-splines; however, we should take into account that the last two B-spline columns generate multicollinearity issues. Therefore, we exclude them when generating random realizations of the splines by drawing the coefficients from  $N(0, 0.35^2)$ . The intercept is fixed to maintain the scale in Figure 11.7, where we observe the data points,  $f(x)$ , and the different splines associated with various random realizations of the coefficients. The following code demonstrates how to perform this simulation.

### R code. Realizations of splines

```

1 ##### Simulation: B-splines #####
2 rm(list = ls())
3 library(ggplot2); library(splines)
4 set.seed(010101)
5 x <- seq(0, 1, 0.001)
6 ysignal <- 0.4 + 0.25*sin(8*x - 5) + 0.4*exp(-16*(4*x - 2.5)
    ^2)
7 sig <- 0.15
8 e <- rnorm(length(ysignal), 0, sd = sig)
9 y <- ysignal + e
10 N <- 100
11 ids <- sort(sample(1:length(ysignal), N))
12 xobs <- x[ids]
13 yobs <- y[ids]
14 knots <- seq(0, 1, 0.25)
15 BS <- bs(xobs, knots = knots, degree = 3, Boundary.knots =
    range(x), intercept = FALSE)
16 # Splines
17 Spline1 <- 0.56 + BS[,-c(7:8)] %*% rnorm(6, 0, 0.35)
18 Spline2 <- 0.56 + BS[,-c(7:8)] %*% rnorm(6, 0, 0.35)
19 Spline3 <- 0.56 + BS[,-c(7:8)] %*% rnorm(6, 0, 0.35)
20 Spline4 <- 0.56 + BS[,-c(7:8)] %*% rnorm(6, 0, 0.35)
21 # Create data frames for the true signal, observed data, and
    Splines
22 data_true_signal <- data.frame(x = x, y = ysignal, Type = "
    True Signal")
23 data_obs <- data.frame(x = xobs, y = yobs, Type = "Observed
    Data")
24 # Create separate data frames for each Spline
25 data_Spline1 <- data.frame(x = xobs, y = Spline1, Type = "
    Spline 1")
26 data_Spline2 <- data.frame(x = xobs, y = Spline2, Type = "
    Spline 2")
27 data_Spline3 <- data.frame(x = xobs, y = Spline3, Type = "
    Spline 3")
28 data_Spline4 <- data.frame(x = xobs, y = Spline4, Type = "
    Spline 4")
29 data <- rbind(data_true_signal, data_obs, data_Spline1, data
    _Spline2, data_Spline3, data_Spline4)

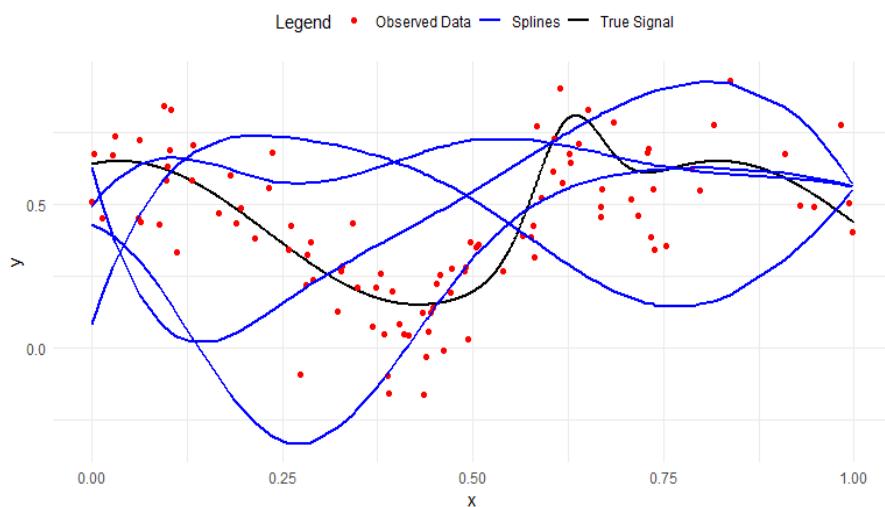
```

### *R code. Realizations of splines*

```

1 ggplot(data, aes(x = x, y = y)) + geom_line(data = subset(
  data, Type == "True Signal"), aes(color = "True Signal")
, linewidth = 1) + geom_point(data = subset(data, Type
== "Observed Data"), aes(color = "Observed Data"), shape
= 16) + geom_line(data = subset(data, Type == "Spline 1"
), aes(color = "Splines"), linewidth = 1, linetype =
"solid") + geom_line(data = subset(data, Type == "Spline
2"), aes(color = "Splines"), linewidth = 1, linetype =
"solid") + geom_line(data = subset(data, Type == "Spline
3"), aes(color = "Splines"), linewidth = 1, linetype =
"solid") + geom_line(data = subset(data, Type == "Spline
4"), aes(color = "Splines"), linewidth = 1, linetype =
"solid") + scale_color_manual(values = c("True Signal" =
"black", "Observed Data" = "red", "Splines" = "blue")) +
  labs(y = "y", color = "Legend") + theme_minimal() +
  theme(legend.position = "top")
2

```



**FIGURE 11.7**  
Different realizations of splines.

We can also use this simulation to examine the effect of changing the knots.

The following code performs the fit using least squares, which is equivalent to using a non-informative prior in a Bayesian linear regression framework.

Figure 11.8 illustrates the fit of different splines based on varying numbers of knots. We observe that increasing the number of knots enhances the flexibility of the spline, providing better local control. However, too many knots can result in a wiggly, overly complex fit that captures noise rather than the true underlying trend and may also increase multicollinearity. Therefore, selecting an appropriate number of knots is a crucial aspect of spline modeling. To address this issue, we can apply the strategies discussed in Chapter 10, and the strategies that will be discussed in Section 12.2 in the next chapter.

### **R code. Demand of marijuana in Colombia: Splines in age**

```

1 rm(list = ls())
2 library(ggplot2); library(splines)
3 # Data generation
4 set.seed(10101)
5 x <- seq(0, 1, 0.001)
6 ysignal <- 0.4 + 0.25*sin(8*x - 5) + 0.4*exp(-16*(4*x - 2.5)^2)
7 sig <- 0.15; e <- rnorm(length(ysignal), 0, sd = sig)
8 y <- ysignal + e; N <- 100
9 ids <- sort(sample(1:length(ysignal), N))
10 xobs <- x[ids]
11 yobs <- y[ids]
12 # Generate Fits with different knot placements
13 knots_list <- list(seq(0, 1, 0.33), seq(0, 1, 0.25), seq(0,
1, 0.2), seq(0, 1, 0.1))
14 Fits <- list()
15 for (i in 1:4) {
16   BS <- bs(xobs, knots = knots_list[[i]], degree = 3,
Boundary.knots = range(x), intercept = FALSE)
17   fm <- lm(yobs ~ BS)
18   Fits[[i]] <- predict(fm)
19 }
20 # Create data frames
21 data_true_signal <- data.frame(x = x, y = ysignal, Type =
"True Signal")
22 data_obs <- data.frame(x = xobs, y = yobs, Type = "Observed
Data")
23 data_preds <- data.frame(
24 x = rep(xobs, 4),
25 y = c(Fits[[1]], Fits[[2]], Fits[[3]], Fits[[4]]),
26 Type = rep(c("Fit 1", "Fit 2", "Fit 3", "Fit 4"), each =
length(xobs)))
27 )

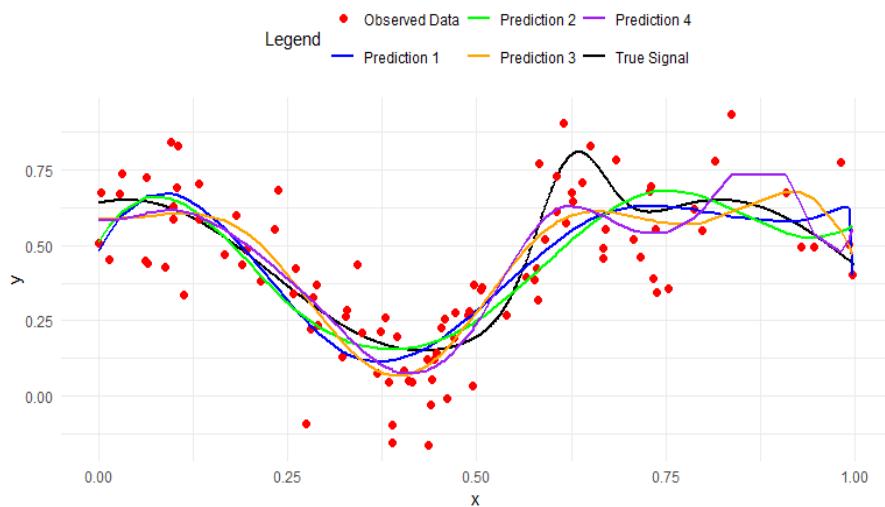
```

**R code. Demand of marijuana in Colombia:  
Splines in age**

```

1 data <- rbind(data_true_signal, data_obs, data_preds)
2 # Create ggplot
3 ggplot(data, aes(x = x, y = y, color = Type)) + geom_line(
  data = subset(data, Type == "True Signal"), linewidth =
  1) + geom_point(data = subset(data, Type == "Observed
  Data"), shape = 16, size = 2) + geom_line(data = subset(
  data, grepl("Fit", Type)), linewidth = 1, linetype =
  "solid") + scale_color_manual(values = c("True Signal" =
  "black", "Observed Data" = "red", "Fit 1" = "blue", "Fit
  2" = "green", "Fit 3" = "orange", "Fit 4" = "purple"))
+ labs(y = "y", color = "Legend") + theme_minimal() +
  theme(legend.position = "top")

```



**FIGURE 11.8**  
Different fits of splines.

Note that splines use local information to perform the fit. This implies the *curse of dimensionality* issue, where increasing the number of variables leads to more dispersed regions in the regressor space. In other words, the observations become increasingly further apart as we have more regressors. As a result, splines become less reliable in high-dimensional spaces.

A strategy to overcome the curse of dimensionality is to specify the *partial linear model*:

$$y_i = \mathbf{z}_i^\top \boldsymbol{\gamma} + f(x_i) + \mu_i. \quad (11.3)$$

In this specification, some regressors enter in a linear way, while potentially the most relevant regressor — the one under primary investigation — enters in a non-parametric way. Note that we only need to increase the dimension of the matrix  $\mathbf{W}$  by adding the columns from  $\mathbf{z}$ , and proceed as we did previously.

#### **Example: Consumption of marijuana in Colombia continues**

Let's continue using the dataset *MarijuanaColombia.csv* to perform a partial linear model as in Equation 11.3, where  $y_i$  is the (log) marijuana monthly consumption,  $\mathbf{z}_i$  represents the presence of a drug dealer in the neighborhood (*Dealer*), gender (*Female*), indicators of good physical and mental health (*PhysicalHealthGood* and *MentalHealthGood*), years of education (*YearsEducation*), and the (log) prices of marijuana, cocaine, and crack by individual.

We set the knots as the percentiles  $\{0, 0.05, \dots, 0.95, 1\}$  of age, use cubic B-splines, non-informative conjugate priors in the linear regression model, 5,000 MCMC iterations, and 5,000 burn-in iterations.

#### **R code. Fit using splines: Different knots**

```

1 rm(list = ls()); set.seed(010101)
2 library(splines); library(ggplot2)
3 Data <- read.csv("https://raw.githubusercontent.com/
   BEsmarter-consultancy/BSTApp/reviews/heads/master/DataApp/
   MarijuanaColombia.csv")
4 attach(Data)
5 IdOrd <- order(Age)
6 y <- LogMarijuana[IdOrd]
7 Z <- as.matrix(cbind(Data[IdOrd,-c(1, 6, 7)]))
8 x <- Age[IdOrd]
9 knots <- quantile(x, seq(0, 1, 0.05))
10 BS <- bs(x, knots = knots, degree = 3, Boundary.knots =
   range(x), intercept = FALSE)
11 matplot(x, BS, type = "l", lty = 1, col = rainbow(ncol(BS)))
12 X <- cbind(1, BS[,1:22], Z) # Get B-splines without
   multicollinearity issues
13 k <- dim(X)[2]; N <- dim(X)[1]
14 # Hyperparameters
15 d0 <- 0.001; a0 <- 0.001
16 b0 <- rep(0, k); c0 <- 1000; B0 <- c0*diag(k); B0i <- solve(
   B0)

```

### R code. Fit using splines: Different knots

```

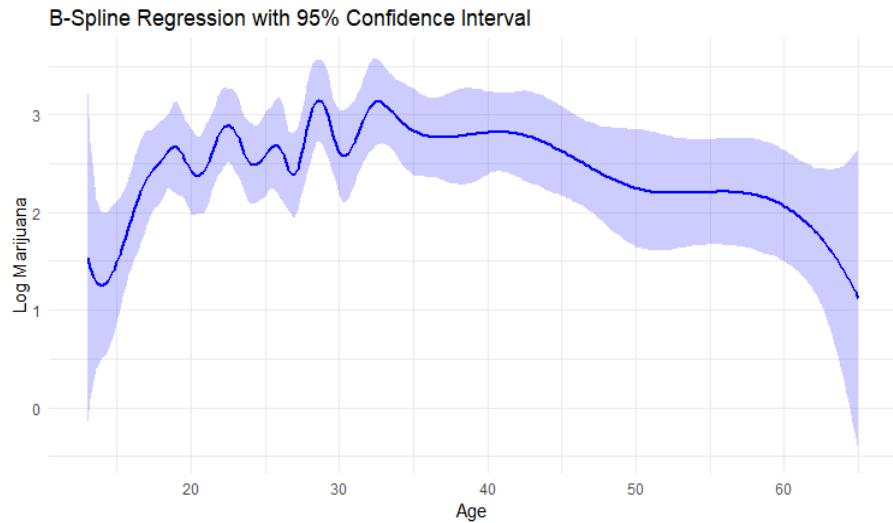
1 # MCMC parameters
2 mcmc <- 5000; burnin <- 5000
3 tot <- mcmc + burnin; thin <- 1
4 posterior <- MCMCpack::MCMCregress(y~X-1, b0=b0, B0 = B0i,
   c0 = a0, d0 = d0, burnin = burnin, mcmc = mcmc, thin =
   thin)
5 summary(coda::mcmc(posterior))
6 # Predict values with 95% credible intervals
7 xfit <- seq(min(x), max(x), 0.2)
8 H <- length(xfit)
9 i <- sample(1:N, 1)
10 idfit <- sample(1:N, H)
11 BSfit <- bs(xfit, knots = knots, degree = 3, Boundary.knots
   = range(x), intercept = FALSE)
12 Xfit <- cbind(1, BSfit[,1:22], Z[rep(i, H),]) # Relevant
   regressors, PIP > 0.5
13 Fit <- matrix(NA, mcmc, H)
14 # posterior[posterior > 0] <- 0
15 for(s in 1:mcmc){
16   Fit[s,] <- Xfit%*%posterior[s,1:31]
17 }
18 # Create a data frame for ggplot
19 plot_data <- data.frame(x = xfit, fit = colMeans(Fit),
   liminf = apply(Fit, 2, quantile, 0.025), limsup = apply(
   Fit, 2, quantile, 0.975))
20 ggplot() + geom_line(data = plot_data, aes(x, fit), color =
   "blue", linewidth = 1) + geom_ribbon(data = plot_data,
   aes(x, ymin = liminf, ymax = limsup), fill = "blue",
   alpha = 0.2) + labs(title = "B-Spline Regression with
   95% Confidence Interval", x = "Age", y = "Log Marijuana"
   ) + theme_minimal()

```

The previous code illustrates the procedure. The posterior 95% credible intervals indicate that marijuana is an inelastic good, there is substitution between marijuana and crack, more educated female individuals consume less, and the presence of a drug dealer in the neighborhood increases consumption.

Figure 11.9 shows the fitted spline, specifically the mean (blue line) and the 95% credible intervals (shaded blue). Notice that there is a lot of wigginess due to the use of many knots, suggesting that this relationship could be expressed in a more parsimonious way. In Exercise 7, we ask to use variable selection methods, such as the BIC approximation presented in Chapter 10, to perform regressor selection, particularly for the splines.

In general, the selection of  $H$  is a problem of variable selection (regressor uncertainty) and can be handled using other approaches, such as those discussed in the regularization section 12.2 of the following chapter.

**FIGURE 11.9**

Spline: Marijuana monthly consumption vs age.

### 11.3 Summary

In this chapter, we introduce the most common Bayesian methods for inference in non-parametric and semi-parametric models. Finite Gaussian and Dirichlet process mixtures, as well as splines, are highly flexible methods; however, they also have limitations, such as the label-switching issue in mixtures and potential multicollinearity in splines. Additionally, it is wise to elicit informative priors and conduct sensitivity analyses to assess the robustness of the results in real-world applications.

### 11.4 Exercises

1. Simulate a semi-parametric regression where

$$\begin{aligned} y_i &= 0.5x_{i1} - 1.2x_{i2} + \mu_i, \\ p(\mu_i) &= 0.7\phi(\mu_i \mid -0.5, 0.5^2) + 0.3\phi(\mu_i \mid 1, 0.8^2). \end{aligned}$$

Assume that  $x_{i1}$  and  $x_{i2}$  follow a standard normal distribution and that the sample size is 1,000. Perform inference in this model assuming that the number of components is unknown. Start with

$H = 5$  and use non-informative priors, setting  $\alpha_{h0} = \delta_{h0} = 0.01$ ,  $\beta_0 = \mathbf{0}_2$ ,  $B_0 = I_2$ ,  $\mu_{h0} = 0$ ,  $\sigma_{h0}^2 = 10$ , and  $\alpha_0 = [1/H \dots 1/H]^\top$ . Use 6,000 MCMC iterations, a burn-in period of 4,000, and a thinning parameter of 2. Compare the population parameters with the posterior estimates and plot the population density along with the posterior density estimate of  $\mu$  (the mean, and the 95% credible interval).

2. **Example: Consumption of marijuana in Colombia continues I**

Use the dataset *MarijuanaColombia.csv* from our GitHub repository to perform inference on the demand for marijuana in Colombia. This dataset contains information on the (log) monthly demand in 2019 from the National Survey of the Consumption of Psychoactive Substances. It includes variables such as the presence of a drug dealer in the neighborhood (*Dealer*), gender (*Female*), indicators of good physical and mental health (*PhysicalHealthGood* and *MentalHealthGood*), age (*Age* and *Age2*), years of schooling (*YearsEducation*), and (log) prices of marijuana, cocaine, and crack by individual (*LogPriceMarijuana*, *LogPriceCocaine*, and *LogPriceCrack*). The sample size is 1,156.

Estimate a finite Gaussian mixture regression using non-informative priors, that is,  $\alpha_0 = \delta_0 = 0.01$ ,  $\beta_0 = \mathbf{0}_K$ ,  $B_0 = I_K$ , and  $\alpha_0 = [1/H \dots 1/H]^\top$ ,  $K$  is the number of regressors, 11 including the intercept. The number of MCMC iterations is 5,000, the burn-in is 1,000, and the thinning parameter is 2. Start with five potential clusters. Obtain the posterior distribution of the own-price elasticity of marijuana and the cross-price elasticities of marijuana demand with respect to the prices of cocaine and crack.

3. Get the posterior sampler in the semi-parametric setting using a Dirichlet process mixture:

$$\begin{aligned} y_i &= \mathbf{x}_i^\top \boldsymbol{\beta} + e_i \\ e_i \mid \mu_i, \sigma_i^2 &\stackrel{iid}{\sim} N(\mu_i, \sigma_i^2), \end{aligned}$$

Do not include the intercept in  $\boldsymbol{\beta}$  to get flexibility in the distribution of the stochastic errors.

Let's assume  $\boldsymbol{\beta} \sim N(\beta_0, B_0)$ ,  $\sigma_i^2 \sim IG(\alpha_0/2, \delta_0/2)$ ,  $\mu_i \sim N(\mu_0, \sigma_i^2/\beta_0)$ ,  $\alpha \sim G(a, b)$  such that introducing the latent variable  $\xi \mid \alpha, N \sim Be(\alpha + 1, N)$ , allows to easily sample the posterior draws of  $\alpha \mid \xi, H, \pi_\xi \sim \pi_\xi G(a+H, b-\log(\xi)) + (1-\pi_\xi)G(a+H-1, b-\log(\xi))$ , where  $\frac{\pi_\xi}{1-\pi_\xi} = \frac{a+H-1}{N(b-\log(\xi))}$ ,  $H$  is the number of atoms (mixture components).

4. **Example: Exercise 1 and 3 continue**

Perform inference in the simulation of the semi-parametric model of Exercise 1 using the sampler of Exercise 3. Use non-informative priors, setting  $\alpha_0 = \delta_0 = 0.01$ ,  $\beta_0 = \mathbf{0}_2$ ,  $B_0 = I_2$ , and  $a = b = 0.1$ . The number of MCMC iterations is 5,000, the burn-in is 1,000, and the thinning parameter is 2.

**5. Example: Simulation exercise continues I**

Fix the label-switching problem of the simulation exercise of the DPM using *random permutation of latent classes*.

**6. Example: Simulation exercise continues II**

Obtain the density estimate of  $y$  in the simulation exercise of the DPM, evaluated at  $x = 0$ .

**7. Example: Consumption of marijuana in Colombia continues II**

Perform the application of marijuana consumption with the following specification:

$$y_i = \mathbf{z}_i^\top \boldsymbol{\gamma} + f(\text{Age}_i) + \mu_i,$$

where  $y_i$  is the (log) marijuana monthly consumption,  $\mathbf{z}_i$  represents the presence of a drug dealer in the neighborhood (*Dealer*), gender (*Female*), indicators of good physical and mental health (*PhysicalHealthGood* and *MentalHealthGood*), years of education (*YearsEducation*), and the (log) prices of marijuana, cocaine, and crack by individual.

Initially, set the knots as the percentiles  $\{0, 0.05, \dots, 0.95, 1\}$  of age and use cubic B-splines. Then, apply the BIC approximation to perform variable selection in this model with non-informative conjugate priors, 5,000 MCMC iterations, and 5,000 burn-in iterations.

Do you think that using a linear regression with a second-degree polynomial in age provides a good approximation to the relationship found using splines in this application?

# 12

---

## Bayesian machine learning

---

In this chapter, we focus on Bayesian approaches to *supervised Machine Learning* (ML) problems, where the outcome variable is observed and used to guide prediction or inference. In contrast, *unsupervised ML* refers to settings in which the outcome variable is not observed, such as in clustering or dimensionality reduction.

Machine learning methods are often characterized by high-dimensional parameter spaces, particularly in the context of nonparametric inference. It is important to note that nonparametric inference does not imply the absence of parameters, but rather models with potentially infinitely many parameters. This setting is often referred to as the *wide* problem, where the number of input variables, and consequently parameters, can exceed the sample size. Another common challenge in ML is the *tall* problem, which occurs when the sample size is extremely large, necessitating scalable algorithms.

Specifically, we introduce Bayesian ML tools for regression, including regularization techniques, regression trees, and Gaussian processes. Extensions of these methods for binary classification are explored in some of the exercises. The section begins with a discussion on the relationship between cross-validation and Bayes factors, and concludes with Bayesian approaches for addressing large-scale data challenges.

---

### 12.1 Cross-validation and Bayes factors

Prediction is central to machine learning, particularly in supervised learning. The starting point is a set of raw regressors or features, denoted by  $\mathbf{x}$ , which are used to construct a set of input variables fed into the model:  $\mathbf{w} = T(\mathbf{x})$ , where  $T(\cdot)$  represents a *dictionary of transformations* such as polynomials, interactions between variables, or the application of functions like exponentials, and so on. These inputs are then used to predict  $y$  through a model  $f(\mathbf{w})$ :

$$y_i = f(\mathbf{w}_i) + \mu_i,$$

where  $\mu_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$ .

We predict  $y$  using  $\hat{f}(\mathbf{w})$ , a model trained on the data. The expected squared error (ESE) at a fixed input  $\mathbf{w}_i$  is given by:

$$\begin{aligned}
\text{ESE} &= \mathbb{E}_{\mathcal{D},y} [(y_i - \hat{f}(\mathbf{w}_i))^2] \\
&= \mathbb{E}_{\mathcal{D},\mu} \left[ (f(\mathbf{w}_i) + \mu_i - \hat{f}(\mathbf{w}_i))^2 \right] \\
&= \mathbb{E}_{\mathcal{D},\mu} \left[ (f(\mathbf{w}_i) - \hat{f}(\mathbf{w}_i))^2 + 2\mu_i(f(\mathbf{w}_i) - \hat{f}(\mathbf{w}_i)) + \mu_i^2 \right] \\
&= \mathbb{E}_{\mathcal{D}} \left[ (f(\mathbf{w}_i) - \hat{f}(\mathbf{w}_i))^2 \right] + \mathbb{E}_{\mu} [\mu_i^2] \\
&= \mathbb{E}_{\mathcal{D}} \left[ ((f(\mathbf{w}_i) - \bar{f}(\mathbf{w}_i)) - (\hat{f}(\mathbf{w}_i) - \bar{f}(\mathbf{w}_i)))^2 \right] + \sigma^2 \\
&= \underbrace{\mathbb{E}_{\mathcal{D}} \left[ (f(\mathbf{w}_i) - \bar{f}(\mathbf{w}_i))^2 \right]}_{\text{Bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} \left[ (\hat{f}(\mathbf{w}_i) - \bar{f}(\mathbf{w}_i))^2 \right]}_{\text{Variance}} + \underbrace{\sigma^2}_{\text{Irreducible Error}}.
\end{aligned}$$

Here,  $\mathcal{D}$  denotes the distribution over datasets defined by the feature space. Independence between the noise  $\mu_i \sim \mathcal{N}(0, \sigma^2)$  and the estimator ensures that  $\mathbb{E}_{\mathcal{D},\mu} [\mu_i(f(\mathbf{w}_i) - \hat{f}(\mathbf{w}_i))] = 0$ . We also define  $\bar{f}(\mathbf{w}_i) = \mathbb{E}_{\mathcal{D}}[\hat{f}(\mathbf{w}_i)]$  as the expected predictor across datasets.

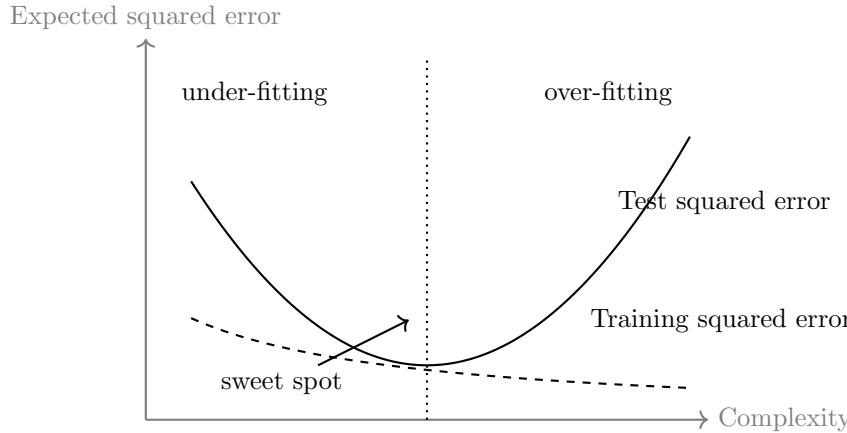
Thus, the ESE is composed of the squared bias, the variance of the prediction (which is a random variable due to data variation  $\mathcal{D}$ ), and the irreducible error. The key insight is that increasing model complexity, such as by including more inputs, typically reduces bias but increases variance. This trade-off can lead to overfitting, where models perform well on the training data but poorly on new, unseen data. There is, therefore, an optimal point at which the predictive error is minimized (see Figure 12.1).<sup>1</sup>

To avoid overfitting in machine learning, an important step called *cross-validation* is often employed. This involves splitting the dataset into multiple parts (called *folds*) and systematically training and testing models on these parts [195, 134]. The main goal is to evaluate how well models generalize to “unseen data”.

There is a compelling justification for cross-validation within Bayesian inference proposed by [41] in their section 6.1.6. The point of departure is to assume an  $\mathcal{M}$ -open view of nature, in which there exists a set of models  $\mathcal{M} = \{M_j : j \in J\}$  under comparison, but none of them represents the true data-generating process, which is assumed to be unknown. Nevertheless, we can compare the models in  $\mathcal{M}$  based on their posterior risk functions (see Chapter 1) without requiring the specification of a true model. In this framework, we select the model that minimizes the posterior expected loss. Unfortunately, we cannot explicitly compute this posterior expected loss due

---

<sup>1</sup>However, recent developments show that powerful modern machine learning methods, such as deep neural networks, often overfit and yet generalize remarkably well on unseen data. This phenomenon is known as *double descent* or *benign overfitting* [35, 28, 193].



**Notes:** Curves for training squared error (dashed line) and test squared error (solid line). The classical U-shaped error curve arising from the bias–variance trade-off.

**FIGURE 12.1**  
*Bias-variance trade-off* in machine learning.

to the lack of knowledge of the true posterior distribution under the  $\mathcal{M}$ -open assumption.<sup>2</sup>

Given the expected loss conditional on model  $j$  for a predictive problem, where the action  $a$  is chosen to minimize the posterior expected loss:

$$\mathbb{E}_{y_0}[L(Y_0, a) | M_j, \mathbf{y}] = \int_{\mathcal{Y}_0} L(y_0, a | M_j, \mathbf{y}) \pi(y_0 | \mathbf{y}) dy_0,$$

we note that there are  $N$  possible partitions of the dataset  $\mathbf{y} = \{y_1, y_2, \dots, y_N\}$  following a leave-one-out strategy, i.e.,  $\mathbf{y} = \{\mathbf{y}_{-k}, y_k\}$ , for  $k = 1, \dots, N$ , where  $\mathbf{y}_{-k}$  denotes the dataset excluding observation  $y_k$ . Assuming that  $\mathbf{y}$  is exchangeable, meaning its joint distribution is invariant to reordering, and that  $N$  is large, then  $\mathbf{y}_{-k}$  and  $y_k$  are good proxies for  $\mathbf{y}$  and  $y_0$ , respectively. Consequently,

$$\frac{1}{K} \sum_{k=1}^K L(y_k, a | M_j, \mathbf{y}_{-k}) \xrightarrow{P} \int_{\mathcal{Y}_0} L(y_0, a | M_j, \mathbf{y}) \pi(y_0 | \mathbf{y}) dy_0,$$

by the law of large numbers, as  $N \rightarrow \infty$  and  $K \rightarrow \infty$ .

Thus, we can select a model by minimizing the expected squared error

---

<sup>2</sup>This is not the case under an  $\mathcal{M}$ -closed view of nature, where one of the candidate models is assumed to be the true data-generating process. In that setting, the posterior distribution becomes a mixture distribution with mixing probabilities given by the posterior model probabilities (see Chapter 10).

based on its expected predictions  $\mathbb{E}[y_k \mid M_j, \mathbf{y}_{-k}]$ ; that is, we select the model with the lowest value of

$$\frac{1}{K} \sum_{k=1}^K (\mathbb{E}[y_k \mid M_j, \mathbf{y}_{-k}] - y_k)^2.$$

Note that if we want to compare two models based on their relative predictive accuracy using the log-score function [298], we select model  $j$  if

$$\int_{\mathcal{Y}_0} \log \frac{p(y_0 \mid M_j, \mathbf{y})}{p(y_0 \mid M_l, \mathbf{y})} \pi(y_0 \mid \mathbf{y}) dy_0 > 0.$$

However, we know that

$$\frac{1}{K} \sum_{k=1}^K \log \frac{p(y_k \mid M_j, \mathbf{y}_{-k})}{p(y_k \mid M_l, \mathbf{y}_{-k})} \xrightarrow{p} \int_{\mathcal{Y}_0} \log \frac{p(y_0 \mid M_j, \mathbf{y})}{p(y_0 \mid M_l, \mathbf{y})} \pi(y_0 \mid \mathbf{y}) dy_0,$$

by the law of large numbers as  $N \rightarrow \infty$  and  $K \rightarrow \infty$ .

This implies that we select model  $j$  over model  $l$  if

$$\prod_{k=1}^K \left( \frac{p(y_k \mid M_j, \mathbf{y}_{-k})}{p(y_k \mid M_l, \mathbf{y}_{-k})} \right)^{1/K} = \prod_{k=1}^K (B_{jl}(y_k, \mathbf{y}_{-k}))^{1/K} > 1,$$

where  $B_{jl}(y_k, \mathbf{y}_{-k})$  is the Bayes factor comparing model  $j$  to model  $l$ , based on the posterior  $\pi(\boldsymbol{\theta}_m \mid M_m, \mathbf{y}_{-k})$  and the predictive  $\pi(y_k \mid \boldsymbol{\theta}_m)$ , for  $m \in \{j, l\}$ .

Thus, under the log-score function, cross-validation reduces to the geometric average of Bayes factors that evaluate predictive performance based on the leave-one-out samples  $\mathbf{y}_{-k}$ .

## 12.2 Regularization

In this century, the amount of available data continues to grow. This means we have access to more covariates for prediction, and we can also generate additional inputs to enhance the predictive power of our models. As a result, we often encounter *wide* datasets, where the number of inputs may exceed the number of observations. Even in modest settings, we might have hundreds of inputs, and we can use them to identify which ones contribute to accurate predictions. However, we generally avoid using all inputs at once due to the risk of overfitting. Thus, we require a class of input selection or *regularization*.

In the standard linear regression setting,

$$\mathbf{y} = \mathbf{i}_N \boldsymbol{\beta}_0 + \mathbf{W} \boldsymbol{\beta} + \boldsymbol{\mu},$$

where  $\mathbf{i}_N$  is an  $N$ -dimensional vector of ones,  $\mathbf{W}$  is the  $N \times K$  design matrix of inputs, and  $\boldsymbol{\mu} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$ , there has been extensive development of techniques aimed at regularization within the Frequentist inferential framework. These include methods such as Ridge regression [207]; discrete subset selection techniques like best subset selection [152], forward selection, and backward stepwise selection [195]; as well as continuous subset selection approaches such as the LASSO [432], the elastic net [471], and OSCAR [53].

It is important to note, however, that Ridge regression does not perform variable selection; rather, it shrinks coefficient estimates toward zero without setting them exactly to zero.

Ridge regression and the LASSO can be viewed as special cases of a more general class of estimators known as *Bridge regression* [151], which also admits a Bayesian interpretation. Consider the following penalized least squares criterion in the linear regression setting:

$$\hat{\boldsymbol{\beta}}^{\text{Bridge}} = \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{k=1}^K \tilde{w}_{ik} \beta_k \right)^2 + \lambda \sum_{k=1}^K |\beta_k|^q \right\}, \quad q \geq 0,$$

where  $\tilde{w}_{ik}$  denotes the standardized inputs. Standardizing inputs is important in variable selection problems to avoid issues caused by differences in scale; otherwise, variables with larger magnitudes will be penalized less and disproportionately influence the regularization path.

Interpreting  $|\beta_k|^q$  as proportional to the negative log-prior density of  $\beta_k$ , the penalty shapes the contours of the prior distribution on the parameters. Specifically:

- $q = 0$  corresponds to best subset selection, where the penalty counts the number of nonzero coefficients.
- $q = 1$  yields the LASSO, which corresponds to a Laplace (double-exponential) prior.
- $q = 2$  yields ridge regression, which corresponds to a Gaussian prior.

In this light, best subset selection, the LASSO, and ridge regression can be viewed as maximum a posteriori (MAP) estimators under different priors centered at zero [330]. However, they are not Bayes estimators in the strict sense, since Bayes estimators are typically defined as the posterior *mean*. While ridge regression coincides with the posterior mean under a Gaussian prior [219], the LASSO and best subset selection yield posterior modes.

This distinction is important because the Bayesian framework naturally incorporates regularization through the use of proper priors, which helps mitigate overfitting. Specifically, when proper shrinkage priors are used, the posterior balances data likelihood and prior information, thus controlling model complexity.

Furthermore, empirical Bayes methods, where the marginal likelihood is

optimized, or cross-validation can be used to estimate the scale parameter of the prior covariance matrix for the regression coefficients. This scale parameter, in turn, determines the strength of regularization in ridge regression.

Note that regularization introduces bias into the parameter estimates because it constrains the model, shrinking the location parameters toward zero. However, it substantially reduces variance, as the estimates are prevented from varying excessively across samples. As a result, the mean squared error (MSE) of the estimates, which equals the sum of the squared bias and the variance, is often lower for regularization methods compared to ordinary least squares (OLS), which remains unbiased under the classical assumptions. This trade-off is particularly important when the goal is to identify causal effects, where unbiasedness may be preferred over predictive accuracy (see Chapter 13).

### 12.2.1 Bayesian LASSO

Given the popularity of the LASSO as a variable selection technique, we present its Bayesian formulation in this subsection [330]. The Gibbs sampler for the Bayesian LASSO exploits the representation of the Laplace distribution as a scale mixture of normals. This leads to the following hierarchical representation of the model:

$$\begin{aligned} \mathbf{y} | \beta_0, \boldsymbol{\beta}, \sigma^2, \mathbf{W} &\sim N(\mathbf{i}_N \beta_0 + \mathbf{W}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_N), \\ \boldsymbol{\beta} | \sigma^2, \tau_1^2, \dots, \tau_K^2 &\sim N(\mathbf{0}_K, \sigma^2 \mathbf{D}_\tau), \\ \tau_1^2, \dots, \tau_K^2 &\sim \prod_{k=1}^K \frac{\lambda^2}{2} \exp \left\{ -\frac{\lambda^2}{2} \tau_k^2 \right\}, \\ \sigma^2 &\sim \frac{1}{\sigma^2}, \\ \beta_0 &\sim c, \end{aligned}$$

where  $\mathbf{D}_\tau = \text{diag}(\tau_1^2, \dots, \tau_K^2)$  and  $c$  is a constant.

After integrating out  $\tau_1^2, \dots, \tau_K^2$ , the conditional prior of  $\boldsymbol{\beta} | \sigma^2$  is:

$$\pi(\boldsymbol{\beta} | \sigma^2) = \prod_{k=1}^K \frac{\lambda}{2\sqrt{\sigma^2}} \exp \left\{ -\frac{\lambda}{\sqrt{\sigma^2}} |\beta_k| \right\},$$

which implies that the log-prior is proportional to  $\lambda \sum_{k=1}^K |\beta_k|$ , matching the penalty term in the LASSO optimization problem.

The conditional posterior distributions for the Gibbs sampler are [330]:

$$\begin{aligned} \boldsymbol{\beta} | \sigma^2, \tau_1^2, \dots, \tau_K^2, \tilde{\mathbf{W}}, \tilde{\mathbf{y}} &\sim N(\boldsymbol{\beta}_n, \sigma^2 \mathbf{B}_n), \\ \sigma^2 | \boldsymbol{\beta}, \tau_1^2, \dots, \tau_K^2, \tilde{\mathbf{W}}, \tilde{\mathbf{y}} &\sim \text{Inverse-Gamma}(\alpha_n/2, \delta_n/2), \\ 1/\tau_k^2 | \boldsymbol{\beta}, \sigma^2 &\sim \text{Inverse-Gaussian}(\mu_{kn}, \lambda_n), \\ \beta_0 | \sigma^2, \tilde{\mathbf{W}}, \tilde{\mathbf{y}} &\sim N(\bar{y}, \sigma^2/N) \end{aligned}$$

where:

$$\begin{aligned}\boldsymbol{\beta}_n &= \mathbf{B}_n \tilde{\mathbf{W}}^\top \tilde{\mathbf{y}}, \\ \mathbf{B}_n &= \left( \tilde{\mathbf{W}}^\top \tilde{\mathbf{W}} + \mathbf{D}_\tau^{-1} \right)^{-1}, \\ \alpha_n &= (N - 1) + K, \\ \delta_n &= (\tilde{\mathbf{y}} - \tilde{\mathbf{W}}\boldsymbol{\beta})^\top (\tilde{\mathbf{y}} - \tilde{\mathbf{W}}\boldsymbol{\beta}) + \boldsymbol{\beta}^\top \mathbf{D}_\tau^{-1} \boldsymbol{\beta}, \\ \mu_{kn} &= \sqrt{\frac{\lambda^2 \sigma^2}{\beta_k^2}}, \\ \lambda_n &= \lambda^2,\end{aligned}$$

where  $\tilde{\mathbf{W}}$  is the matrix of standardized inputs and  $\tilde{\mathbf{y}}$  is the centered response vector.

Note that the posterior distribution of  $\boldsymbol{\tau}$  depends on the data through  $\boldsymbol{\beta}$  and  $\sigma^2$ , which is a typical feature of hierarchical models. In this formulation, we can interpret  $\tau_k$  as local shrinkage parameters, while  $\lambda$  acts as a global shrinkage parameter. Higher values of  $\tau_k$  and  $\lambda$  imply stronger shrinkage toward zero. [330] propose two approaches for specifying the global shrinkage parameter: empirical Bayes estimation or a fully Bayesian hierarchical specification, where  $\lambda^2$  is assigned a Gamma prior.

We should acknowledge that the Bayesian LASSO is more computationally expensive than the Frequentist LASSO. However, it provides credible intervals for the parameters automatically. In contrast, obtaining standard errors in the Frequentist LASSO is more challenging, particularly for parameters estimated to be exactly zero [254].

#### **Example: Simulation exercise to study the Bayesian LASSO performance**

We simulate the process  $y_i = \beta_0 + \sum_{k=1}^{10} \beta_k w_{ik} + \mu_i$ , where  $\beta_k \sim \mathcal{U}(-3, 3)$ ,  $\mu_i \sim \mathcal{N}(0, 1)$ , and  $w_{ik} \sim \mathcal{N}(0, 1)$ ,  $i = 1, 2, \dots, 500$ . Additionally, we generate 90 extra potential inputs from a standard normal distribution, which are included in the model specification. Our goal is to assess whether the Bayesian LASSO can successfully identify the truly relevant inputs.

We use the *bayesreg* package in **R** to perform the Bayesian LASSO, using 5,000 MCMC draws and 1,000 burn-in iterations. The following code illustrates the simulation exercise and compares the posterior means with the true population values.

The summary of the fit and the plot comparing the population parameters with the posterior means show that the Bayesian LASSO is able to identify the variables that are relevant in the data generating process. In Exercise 1, we propose programming the Gibbs sampler from scratch, assuming a hierarchical structure for the global shrinkage parameter, and comparing the results with those obtained using the *monomvn* package.

### R code. The Bayesian LASSO

```

1 rm(list = ls()); set.seed(10101)
2 library(bayesreg)
3 # Parameters
4 n <- 500 # sample size
5 p <- 100 # number of predictors
6 s <- 10 # number of non-zero coefficients
7 # Generate design matrix
8 X <- matrix(rnorm(n * p), nrow = n, ncol = p)
9 # True beta: first s coefficients are non-zero, rest are
# zero
10 beta_true <- c(runif(s, -3, 3), rep(0, p - s))
11 # Generate response with some noise
12 sigma <- 1
13 y <- X %*% beta_true + rnorm(n, sd = sigma)
14 df <- data.frame(X,y)
15 ### Using bayesreg ####
16 # Fit the model
17 fit <- bayesreg::bayesreg(y ~ X, data = df, model = "
gaussian", prior = "lasso",
18 n.samples = 5000, burnin = 1000)
19 # Check summary
20 summary(fit)
21 # Extract posterior means of beta
22 beta_post_mean <- rowMeans(fit$beta)
23 # Compare true vs estimated
24 plot(beta_true, beta_post_mean, pch = 19, col = "steelblue",
25 xlab = "True beta", ylab = "Posterior mean of beta",
26 main = "Bayesian LASSO Shrinkage")
27 abline(0, 1, col = "red", lty = 2)

```

#### 12.2.2 Stochastic search variable selection

Another well-known Bayesian strategy for variable selection in the presence of a large set of regressors (inputs) is *stochastic search variable selection* (SSVS) [168, 167]. SSVS is a particular case of the broader class of *spike-and-slab* priors, in which the prior distribution for the location parameters is specified as a hierarchical mixture that captures the uncertainty inherent in variable selection problems [219]. The hierarchical structure of the model is given by

$$\begin{aligned}
\mathbf{y} | \beta_0, \boldsymbol{\beta}, \sigma^2, \mathbf{W} &\sim N(\mathbf{i}_N \beta_0 + \mathbf{W}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_N), \\
\boldsymbol{\beta} | \sigma^2, \boldsymbol{\gamma} &\sim N(\mathbf{0}_K, \sigma^2 \mathbf{D}_{\gamma} \mathbf{R} \mathbf{D}_{\gamma}), \\
\sigma^2 &\sim \text{Inverse-Gamma} \left( \frac{v}{2}, \frac{v\lambda_{\gamma}}{2} \right),
\end{aligned}$$

$$\gamma_k \sim \text{Bernoulli}(p_k),$$

where  $p_k$  is the prior inclusion probability of regressor  $w_k$ , that is,  $P(\gamma_k = 1) = 1 - P(\gamma_k = 0) = p_k$ ,  $\mathbf{R}$  is a correlation matrix, and  $\mathbf{D}_\gamma$  is a diagonal matrix whose  $kk$ -th element is defined as

$$(\mathbf{D}_\gamma)_{kk} = \begin{cases} v_{0k}, & \text{if } \gamma_k = 0, \\ v_{1k}, & \text{if } \gamma_k = 1 \end{cases}.$$

This formulation implies that  $\beta_k \sim (1 - \gamma_k)N(0, v_{0k}) + \gamma_k N(0, v_{1k})$ , where  $v_{0k}$  and  $v_{1k}$  are chosen such that  $v_{0k}$  is small and  $v_{1k}$  is large. Therefore, when the data favors  $\gamma_k = 0$  over  $\gamma_k = 1$ , the corresponding  $\beta_k$  is shrunk toward zero, effectively excluding input  $k$  from the model. In this sense,  $N(0, v_{0k})$  is a spike prior concentrated at zero, while  $N(0, v_{1k})$  is a diffuse slab prior. A critical aspect of SSVS is the choice of the hyperparameters  $v_{0k}$  and  $v_{1k}$ , as they determine the amount of shrinkage applied to the regression coefficients (see [168, 167] for details).

The assumption  $\gamma_k \sim \text{Bernoulli}(p_k)$  implies that the prior on the inclusion indicators is given by  $\pi(\gamma) = \prod_{k=1}^K p_k^{\gamma_k} (1 - p_k)^{1 - \gamma_k}$ . This means that the inclusion of input  $k$  is independent of the inclusion of any other input  $j \neq k$ . A common choice is the uniform prior  $\pi(\gamma) = 2^{-K}$ , which corresponds to setting  $p_k = 1/2$ , giving each regressor an equal chance of being included [219].

A practical choice for the correlation matrix is to set  $\mathbf{R} \propto (\tilde{\mathbf{W}}^\top \tilde{\mathbf{W}})^{-1}$  [168]. Regarding the hyperparameters  $v$  and  $\lambda_\gamma$ , it is helpful to interpret  $\lambda_\gamma$  as a prior estimate of  $\sigma^2$ , and  $v$  as the prior sample size associated with this estimate. In the absence of prior information, [169] recommend setting  $\lambda_\gamma$  equal to the least squares estimate of the variance from the *saturated model*, that is, the model including all regressors, and  $v$  a small number, for instance, 0.01.

The conditional posterior distributions for the Gibbs sampler are [168]:

$$\begin{aligned} \boldsymbol{\beta} | \sigma^2, \gamma_1, \dots, \gamma_K, \tilde{\mathbf{W}}, \tilde{\mathbf{y}} &\sim N(\boldsymbol{\beta}_n, \mathbf{B}_n), \\ \sigma^2 | \boldsymbol{\beta}, \gamma_1, \dots, \gamma_K, \tilde{\mathbf{W}}, \tilde{\mathbf{y}} &\sim \text{Inverse-Gamma}(\alpha_n/2, \delta_n/2), \\ \gamma_k | \boldsymbol{\beta}, \sigma^2 &\sim \text{Bernoulli}(p_{kn}), \end{aligned}$$

where:

$$\begin{aligned} \boldsymbol{\beta}_n &= \sigma^{-2} \mathbf{B}_n \tilde{\mathbf{W}}^\top \tilde{\mathbf{y}}, \\ \mathbf{B}_n &= \left( \sigma^{-2} \tilde{\mathbf{W}}^\top \tilde{\mathbf{W}} + \mathbf{D}_\gamma^{-1} \mathbf{R}^{-1} \mathbf{D}_\gamma^{-1} \right)^{-1}, \\ \alpha_n &= N + v, \\ \delta_n &= (\tilde{\mathbf{y}} - \tilde{\mathbf{W}} \boldsymbol{\beta})^\top (\tilde{\mathbf{y}} - \tilde{\mathbf{W}} \boldsymbol{\beta}) + v \lambda_\gamma, \end{aligned}$$

$$p_{kn} = \frac{\pi(\boldsymbol{\beta} | \boldsymbol{\gamma}_{-k}, \gamma_k = 1) \times p_k}{\pi(\boldsymbol{\beta} | \boldsymbol{\gamma}_{-k}, \gamma_k = 1) \times p_k + \pi(\boldsymbol{\beta} | \boldsymbol{\gamma}_{-k}, \gamma_k = 0) \times (1 - p_k)},$$

where  $\tilde{\mathbf{W}}$  is the matrix of standardized inputs,  $\tilde{\mathbf{y}}$  is the centered response vector,  $\boldsymbol{\gamma}_{-k}$  denotes the vector composed of  $\gamma_1, \gamma_2, \dots, \gamma_K$  excluding  $\gamma_k$ , and  $\pi(\boldsymbol{\beta} | \boldsymbol{\gamma}_{-k}, \gamma_k = \delta)$  is the posterior density of  $\boldsymbol{\beta}$  evaluated at  $\boldsymbol{\gamma}_{-k}$  and  $\gamma_k = \delta$ , where  $\delta \in \{0, 1\}$ .

In general, it is wise to consider the inclusion of regressors jointly due to potential correlations among them; that is, the marginal frequency of  $\gamma_k = 1$  should be interpreted with caution. SSVS is more effective at identifying a good set of potential models rather than selecting a single best model.

#### Example: Simulation exercise to study SSVS performance

Let's use the simulation setting from the previous example to evaluate the performance of SSVS in uncovering the data-generating process. In particular, we use the *BoomSpikeSlab* package to implement this example.

The analysis is performed using 5,000 posterior draws and the default prior. However, the package allows the user to modify the default prior via the *SpikeSlabPrior* function.

The results show that the posterior inclusion probabilities for regressors 2 through 10 are 100%, and the model with the highest posterior probability (94%) includes all of these nine variables. However, the true data-generating process, which also includes regressor 1, receives a posterior model probability of 0%. This is because the population coefficient of this regressor is essentially zero. The plot comparing the posterior means with the true population parameters indicates good performance of SSVS. In general, Bayesian methods for variable selection perform well, and the choice of the most suitable method largely depends on the prior specification [327].

#### R code. Stochastic search variable selection

```

1 rm(list = ls()); set.seed(10101)
2 library(BoomSpikeSlab)
3 library(dplyr); library(tibble)
4 # Parameters
5 n <- 500 # sample size
6 k <- 100 # number of predictors
7 s <- 10 # number of non-zero coefficients
8 # Generate design matrix
9 X <- matrix(rnorm(n * k), nrow = n, ncol = k)
10 # True beta: first s coefficients are non-zero, rest are
11 # zero
11 beta_true <- c(runif(s, -3, 3), rep(0, k - s))

```

### R code. Stochastic search variable selection

```

1 # Generate response with some noise
2 sigma <- 1
3 y <- X %*% beta_true + rnorm(n, sd = sigma)
4 df <- data.frame(X,y)
5 ##### Using BoomSpikeSlab #####
6 #Scale regressors
7 W <- scale(X); yh <- y - mean(y)
8 prior <- SpikeSlabPrior(W, yh,
9 expected.model.size = ncol(W)/2, # expect 50 nonzero
   predictors
10 prior.df = .01, # weaker prior than the default
11 prior.information.weight = .01,
12 diagonal.shrinkage = 0) # shrink to zero
13 niter <- 5000
14 #####Estimate model#####
15 SSBoomNew <- lm.spike(yh ~ W - 1, niter = niter, prior =
   prior)
16 Models <- SSBoomNew$beta != 0
17 PIP <- colMeans(SSBoomNew$beta != 0)
18 # Convert the logical matrix to a data frame and then to a
   tibble
19 df <- as.data.frame(Models); df_tbl <- as_tibble(df)
20 # Count identical rows
21 row_counts <- df_tbl %>% count(across(everything()), name =
   "frequency") %>% arrange(desc(frequency))
22 sum(row_counts[1:100,101])
23 # Ensure your vector and matrix are logical
24 trueModel <- c(rep(1, 10), rep(0, 90)) == 1 # convert to
   logical if needed
25 # Assume your matrix is named 'mat'
26 matching_rows <- apply(row_counts[,-101], 1, function(row)
   all(row == trueModel))
27 # Get indices (row numbers) where the match is TRUE
28 row_counts[which(matching_rows), 101]
29 # Coefficients
30 SummarySS <- summary(coda::mcmc(SSBoomNew$beta))
31 # Extract posterior means of beta
32 beta_post_mean <- SummarySS$statistics[, 1]
33 # Compare true vs estimated
34 plot(beta_true, beta_post_mean, pch = 19, col = "steelblue",
   xlab = "True beta", ylab = "Posterior mean of beta",
   main = "SSVS Shrinkage")
35 abline(0, 1, col = "red", lty = 2)
36

```

The examples and exercises presented thus far have considered scenarios in which the number of inputs is smaller than the number of observations

$(K < N)$ . In Exercise 4, we challenge the Bayesian LASSO and SSVS in a setting where the number of inputs exceeds the sample size ( $K > N$ ). As you will observe in that experiment, both the Bayesian LASSO and SSVS perform well. However, the Bayesian LASSO requires more time to produce results compared to SSVS in this exercise. [381] propose a connection between the LASSO and spike-and-slab priors for variable selection in linear models, offering oracle properties and optimal posterior concentration even in high-dimensional settings where  $K > N$ .

In addition, there are other Bayesian methods for regularization, such as the spike-and-slab approach proposed by [219] and non-local priors introduced by [230], which can be implemented using the **R** packages *spikeslab* and *mombf*, respectively.

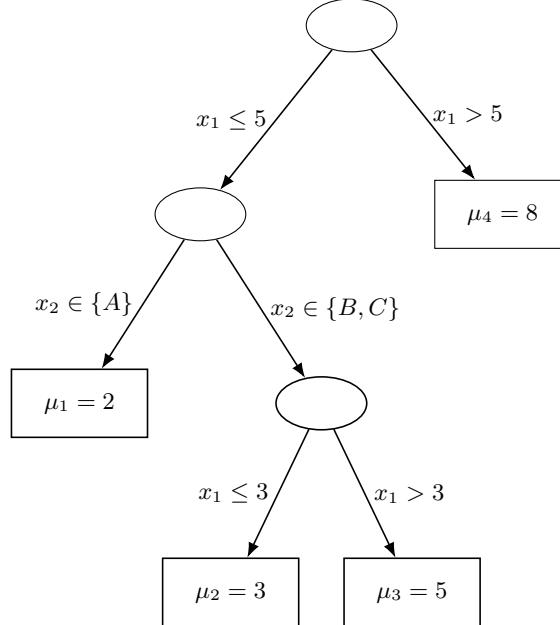
### 12.3 Bayesian additive regression trees

A Classification and Regression Tree (CART) is a method used to predict outcomes based on inputs without assuming a parametric model. It is referred to as a classification tree when the outcome variable is qualitative, and as a regression tree when the outcome variable is quantitative. The method works by recursively partitioning the dataset into smaller, more homogeneous subsets using decision rules based on the input variables. For regression tasks, CART selects splits that minimize prediction error, typically measured by the sum of squared residuals. For classification problems, it aims to create the purest possible groups, using criteria such as Gini impurity or entropy. The result is a tree-like structure in which each internal node represents a decision based on one variable, and each leaf node corresponds to a prediction. CART was popularized in the statistical community by [61].

Figure 12.2 displays a binary regression tree with two variables: one continuous ( $x_1$ ) and one categorical ( $x_2 \in A, B, C$ ). The tree has seven nodes in total, four of which are terminal nodes ( $B = 4$ ), dividing the input space  $\mathbf{x} = (x_1, x_2)$  into four non-overlapping regions. Each internal node indicates the splitting variable and rule, while each terminal node shows the value  $\theta_b$ , representing the conditional mean of  $y$  given  $\mathbf{x}$ .

The first split is based on the rule  $x_1 \leq 5$  (left) versus  $x_1 > 5$  (right). The leftmost terminal node corresponds to  $x_2 \in A$  with  $\mu_1 = 2$ . The second terminal node, with  $\mu_2 = 3$ , is defined by  $x_1 \leq 3$  and  $x_2 \in B, C$ . The third node assigns  $\mu_3 = 5$  for  $3 < x_1 \leq 5$  and  $x_2 \in B, C$ . Finally, the rightmost terminal node assigns  $\mu_4 = 8$  for all observations with  $x_1 > 5$ .

We can view a CART model as specifying the conditional distribution of  $y$  given the vector of features  $\mathbf{x} = [x_1, x_2, \dots, x_K]^\top$ . There are two main components: (i) the binary tree structure  $T$ , which consists of a sequence of binary decision rules of the form  $x_k \in A$  versus  $x_k \notin A$ , where  $A$  is a



**Notes:** Decision tree with one continuous variable ( $x_1$ ) and one categorical variable ( $x_2 \in \{A, B, C\}$ ). The tree has four terminal nodes, each associated with a value denoted by  $\mu_b$ .

**FIGURE 12.2**

Decision tree.

subset of the range of  $x_k$ , and  $B$  terminal nodes that define a non-overlapping and exhaustive partition of the input space; and (ii) the parameter vector  $\boldsymbol{\theta} = [\mu_1, \mu_2, \dots, \mu_B]^\top$  and  $\sigma^2$ , where each  $\mu_b$  corresponds to the parameter associated with terminal node  $b$ , and  $\sigma^2$  is the variance. Consequently, the response  $y \mid \mathbf{x} \sim p(y \mid \mu_b, \sigma^2)$  if  $\mathbf{x}$  belongs to the region associated with terminal node  $b$ , where  $p$  denotes a parametric distribution indexed by  $\mu_b$  and  $\sigma^2$ .

Assuming that  $y_{bi}$  is independently and identically distributed within each terminal node and independently across nodes, for  $b = 1, 2, \dots, B$  and  $i = 1, 2, \dots, n_b$ , we have:

$$\begin{aligned} p(\mathbf{y} \mid \mathbf{x}, T, \boldsymbol{\theta}, \sigma^2) &= \prod_{b=1}^B p(\mathbf{y}_b \mid \mu_b, \sigma^2) \\ &= \prod_{b=1}^B \prod_{i=1}^{n_b} p(y_{bi} \mid \mu_b, \sigma^2), \end{aligned}$$

where  $\mathbf{y}_b = [y_{b1}, y_{b2}, \dots, y_{bn_b}]^\top$  is the set of observations in the terminal node  $b$ .

[99] introduced a Bayesian formulation of CART models, in which inference is carried out through a combination of prior specification on the binary tree structure  $T$  and the parameters  $\boldsymbol{\theta}, \sigma^2 \mid T$ , along with a stochastic search strategy based on a Metropolis-Hastings algorithm. The transition kernels used in the algorithm include operations such as growing, pruning, changing, and swapping tree branches, and candidate trees are evaluated based on their marginal likelihood. This approach enables exploration of a richer set of potential tree models and offers a more flexible and effective alternative to the greedy algorithms commonly used in classical CART.

While CART is a simple yet powerful tool, it is prone to overfitting. To mitigate this, ensemble methods such as boosting, bagging, and random forests are often used. *Boosting* combines multiple weak trees sequentially, each correcting the errors of its predecessor, to create a strong predictive model [146]. *Bagging* builds multiple models on bootstrapped datasets and averages their predictions to reduce variance [59], and *random forests* extend bagging by using decision trees and adding random input selection at each split to increase model diversity [60]. Although a single tree might overfit and generalize poorly, aggregating many randomized trees typically yields more accurate and stable predictions.

[100] introduced Bayesian Additive Regression Trees (BART). The starting point is the model

$$y_i = f(\mathbf{x}_i) + \mu_i,$$

where  $\mu_i \sim N(0, \sigma^2)$ .

Thus, the conditional expectation  $\mathbb{E}[y_i \mid \mathbf{x}_i] = f(\mathbf{x}_i)$  is approximated as

$$\begin{aligned} f(\mathbf{x}_i) &\approx h(\mathbf{x}_i) \\ &= \sum_{j=1}^J g_j(\mathbf{x}_i \mid T_j, \boldsymbol{\theta}_j), \end{aligned}$$

that is,  $f(\mathbf{x}_i)$  is approximated by the sum of  $J$  regression trees. Each tree is defined by a structure  $T_j$  and a corresponding set of terminal node parameters  $\boldsymbol{\theta}_j$ , where  $g_j(\mathbf{x}_i \mid T_j, \boldsymbol{\theta}_j)$  denotes the function that assigns the value  $\mu_{bj} \in \boldsymbol{\theta}_j$  to  $\mathbf{x}_i$  according to the partition defined by  $T_j$ .

The main idea is to construct this sum-of-trees model by imposing a prior that regularizes the fit, ensuring that the individual contribution of each tree remains small. Thus, each tree  $g_j$  explains a small and distinct portion of  $f$ . This is achieved through Bayesian backfitting MCMC [194], where successive fits to the residuals are added. In this sense, BART can be viewed as a Bayesian version of boosting.

[100] use the following prior structure

$$\pi(\{(T_1, \boldsymbol{\theta}_1), (T_2, \boldsymbol{\theta}_2), \dots, (T_J, \boldsymbol{\theta}_J), \sigma^2\}) = \left[ \prod_{j=1}^J \pi(T_j, \boldsymbol{\theta}_j) \right] \pi(\sigma)$$

$$\begin{aligned}
&= \left[ \prod_{j=1}^J \pi(\boldsymbol{\theta}_j \mid T_j) \pi(T_j) \right] \pi(\sigma) \\
&= \left[ \prod_{j=1}^J \prod_{b=1}^B \pi(\mu_{bj} \mid T_j) \pi(T_j) \right] \pi(\sigma).
\end{aligned}$$

The prior for the binary tree structure has three components: (i) the probability that a node at depth  $d = 0, 1, \dots$  is nonterminal, given by  $\alpha(1+d)^{-\beta}$ , where  $\alpha \in (0, 1)$  and  $\beta \in [0, \infty)$ ; the default values are  $\alpha = 0.95$  and  $\beta = 2$ ; (ii) a uniform prior over the set of available regressors to define the distribution of splitting variable assignments at each interior node; and (iii) a uniform prior over the discrete set of available splitting values to define the splitting rule assignment, conditional on the chosen splitting variable.

The prior for the terminal node parameters,  $\pi(\mu_{bj} \mid T_j)$ , is specified as  $N(0, 0.5/(k\sqrt{J}))$ . The observed values of  $y$  are scaled and shifted to lie within the range  $y_{\min} = -0.5$  to  $y_{\max} = 0.5$ , and the default value  $k = 2$  ensures that

$$P_{\pi(\mu_{bj} \mid T_j)} (\mathbb{E}[y \mid \mathbf{x}] \in (-0.5, 0.5)) = 0.95.$$

Note that this prior shrinks the effect of individual trees toward zero, ensuring that each tree contributes only a small amount to the overall prediction. Moreover, although the dependent variable is transformed, there is no need to transform the input variables, since the tree-splitting rules are invariant to monotonic transformations of the regressors.

The prior for  $\sigma^2$  is specified as  $\pi(\sigma^2) \sim v\lambda/\chi_v^2$ . [100] recommend setting  $v = 3$ , and choosing  $\lambda$  such that  $P(\sigma < \hat{\sigma}) = q$ ,  $q = 0.9$ , where  $\hat{\sigma}$  is an estimate of the residual standard deviation from a saturated linear model, i.e., a model including all available regressors when  $K < N$ , or the standard deviation of  $y$  when  $K \geq N$ .

Finally, regarding the number of trees  $J$ , the default value is 200. As  $J$  increases from 1, BART's predictive performance typically improves substantially until it reaches a plateau, after which performance may degrade very slowly. However, if the primary goal is variable selection, using a smaller  $J$  is preferable, as it facilitates identifying the most important regressors by enhancing the internal competition among variables within a limited number of trees.

To sum up, the default hyperparameters are  $(\alpha, \beta, k, J, v, q) = (0.95, 2, 2, 200, 3, 0.9)$ ; however, cross-validation can be used to tune these hyperparameters if desired. BART's predictive performance appears to be relatively robust to the choice of hyperparameters, provided they are set to sensible values, except in cases where  $K > N$ , in which cross-validated tuning tends to yield better performance, albeit at the cost of increased computational time.

Given this specification, we can use a Gibbs sampler that cycles through

the  $J$  trees. At each iteration, we sample from the conditional posterior distribution

$$\pi(T_j, \boldsymbol{\theta}_j | R_j, \sigma) = \pi(T_j | R_j, \sigma) \times \pi(\boldsymbol{\theta}_j | T_j, R_j, \sigma),$$

where  $R_j = y - \sum_{k \neq j} g_k(\mathbf{x} | T_k, \boldsymbol{\theta}_k)$  represents the residuals excluding the contribution of the  $j$ -th tree.

The posterior distribution  $\pi(T_j | R_j, \sigma)$  is explored using a Metropolis-Hastings algorithm, where the candidate tree is generated by one of the following moves: (i) growing a terminal node with probability 0.25; (ii) pruning a pair of terminal nodes with probability 0.25; (iii) changing a nonterminal splitting rule with probability 0.4; or (iv) swapping a rule between a parent and child node with probability 0.1 [99].

The posterior distribution of  $\boldsymbol{\theta}_j$  is the product of the posterior distributions  $\pi(\mu_{jb} | T_j, R_j, \sigma)$ , which are Gaussian. The posterior distribution of  $\sigma$  is inverse-gamma. The posterior draws of  $\mu_{bj}$  are then used to update the residuals  $R_{j+1}$ , allowing the sampler to proceed to the next tree in the cycle. The number of iterations in the Gibbs sampler does not need to be very large; for instance, 200 burn-in iterations and 1,000 post-burn-in iterations typically work well.

As we obtain posterior draws from the sum-of-trees model, we can compute point estimates at each  $\mathbf{x}_i$  using the posterior mean,  $\mathbb{E}[y | \mathbf{x}_i] = f(\mathbf{x}_i)$ . Additionally, pointwise measures of uncertainty can be derived from the quantiles of the posterior draws. We can also identify the most relevant predictors by tracking the relative frequency with which each regressor appears in the sum-of-trees model across iterations. Moreover, we can perform inference on functionals of  $f$ , such as the *partial dependence functions* [148], which quantify the marginal effects of the regressors.

Specifically, if we are interested in the effect of  $\mathbf{x}_s$  on  $y$ , while marginalizing over the remaining variables  $\mathbf{x}_c$ , such that  $\mathbf{x} = [\mathbf{x}_s^\top \ \mathbf{x}_c^\top]^\top$ , the partial dependence function is defined as

$$f(\mathbf{x}_s) = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_s, \mathbf{x}_{ic}),$$

where  $\mathbf{x}_{ic}$  denotes the  $i$ -th observed value of  $\mathbf{x}_c$  in the dataset.

Note that the calculation of the partial dependence function assumes that a subset of the variables is held fixed while averaging over the remainder. This assumption may be questionable when strong dependence exists among input variables, as fixing some variables while varying others may result in unrealistic or implausible combinations. Therefore, caution is warranted when interpreting the results.

[272] extended BART models to high-dimensional settings for prediction and variable selection, while [205, 189] applied them to causal inference. The asymptotic properties of BART models have been studied by [382, 383, 375].

**Example: Simulation exercise to study BART performance**

We use the *BART* package [419] in the **R** software environment to perform estimation, prediction, inference, and marginal analysis using Bayesian Additive Regression Trees. In addition to modeling continuous outcomes, this package also supports dichotomous, categorical, and time-to-event outcomes.

We adopt the simulation setting proposed by [147], which is also used by [100]:

$$y_i = 10 \sin(\pi x_{i1} x_{i2}) + 20(x_{i3} - 0.5)^2 + 10x_{i4} + 5x_{i5} + \mu_i,$$

where  $\mu_i \sim N(0, 1)$ , for  $i = 1, 2, \dots, 500$ .

We set the hyperparameters to  $(\alpha, \beta, k, J, v, q) = (0.95, 2, 2, 200, 3, 0.9)$ , with a burn-in of 100 iterations, a thinning parameter of 1, and 1,000 post-burn-in MCMC iterations. We analyze the trace plot of the posterior draws of  $\sigma$  to assess convergence, compare the true values of  $y$  with the posterior mean and 95% predictive intervals for both the training and test sets (using 80% of the data for training and 20% for testing), and visualize the relative importance of the regressors across different values of  $J = 10, 20, 50, 100, 200$ .

**R code. Bayesian additive regression trees**

```

1 rm(list = ls()); set.seed(10101)
2 library(BART); library(tidyverse)
3 library(ggplot2); library(dplyr)
4 N <- 500; K <- 10
5 # Simulate the dataset
6 MeanFunct <- function(x){
7   f <- 10*sin(pi*x[,1]*x[,2]) + 20*(x[,3]-.5)^2+10*x[,4]+5*x
8   [,5]
9   return(f)
10 }
11 sig2 <- 1
12 e <- rnorm(N, 0, sig2^0.5)
13 X <- matrix(runif(N*K), N, K)
14 y <- MeanFunct(X[,1:5]) + e
15 # Train and test
16 c <- 0.8
17 Ntrain <- floor(c * N)
18 Ntest <- N - Ntrain
19 X.train <- X[1:Ntrain, ]
20 y.train <- y[1:Ntrain]
21 X.test <- X[(Ntrain+1):N, ]
22 y.test <- y[(Ntrain+1):N]
23 # Hyperparameters
24 alpha <- 0.95; beta <- 2; k <- 2
25 v <- 3; q <- 0.9; J <- 200

```

### R code. Bayesian additive regression trees

```

1 # MCMC parameters
2 MCMCiter <- 1000; burnin <- 100; thinning <- 1
3 # Estimate BART
4 BARTfit <- wbart(x.train = X.train, y.train = y.train, x.
+   test = X.test, base = alpha,
5 power = beta, k = k, sigdf = v, sigquant = q, ntree = J,
6 ndpost = MCMCiter, nskip = burnin, keepevery = thinning)
7 # Trace plot sigma
8 keep <- seq(burnin + 1, MCMCiter + burnin, thinning)
9 df_sigma <- data.frame(iteration = 1:length(keep), sigma =
+   BARTfit$sigma[keep])
10 ggplot(df_sigma, aes(x = iteration, y = sigma)) + geom_line(
+   color = "steelblue") + labs(title = "Trace Plot of Sigma",
+   x = "Iteration", y = expression(sigma)) + theme_
+   minimal()
11 # Prediction plot training
12 train_preds <- data.frame( y_true = y.train, mean = apply(
+   BARTfit$yhat.train, 2, mean), lower = apply(BARTfit$yhat.
+   train, 2, quantile, 0.025), upper = apply(BARTfit$yhat.
+   train, 2, quantile, 0.975)) %>% arrange(y_true) %>%
+   mutate(index = row_number())
13 ggplot(train_preds, aes(x = index)) + geom_ribbon(aes(ymin =
+   lower, ymax = upper, fill = "95% Interval"), alpha =
+   0.4, show.legend = TRUE) + geom_line(aes(y = mean, color
+   = "Predicted Mean")) + geom_line(aes(y = y_true, color
+   = "True y")) + scale_color_manual(name = "Line", values
+   = c("Predicted Mean" = "blue", "True y" = "black")) +
+   scale_fill_manual(name = "Interval", values = c("95%
+   Interval" = "lightblue")) + labs(title = "Training Data:
+   Ordered Predictions with 95% Intervals", x = "Ordered
+   Index", y = "y") + theme_minimal()
14 # Prediction plot test
15 test_preds <- data.frame( y_true = y.test, mean = apply(
+   BARTfit$yhat.test, 2, mean), lower = apply(BARTfit$yhat.
+   test, 2, quantile, 0.025), upper = apply(BARTfit$yhat.
+   test, 2, quantile, 0.975)) %>% arrange(y_true) %>%
+   mutate(index = row_number())
16 ggplot(test_preds, aes(x = index)) + geom_ribbon(aes(ymin =
+   lower, ymax = upper, fill = "95% Interval"), alpha =
+   0.4, show.legend = TRUE) + geom_line(aes(y = mean, color
+   = "Predicted Mean")) + geom_line(aes(y = y_true, color
+   = "True y")) + scale_color_manual(name = "Line", values
+   = c("Predicted Mean" = "blue", "True y" = "black")) +
+   scale_fill_manual(name = "Interval", values = c("95%
+   Interval" = "lightblue")) + labs(title = "Test Data:
+   Ordered Predictions with 95% Intervals",
17 x = "Ordered Index", y = "y") + theme_minimal()
18 # Relevant regressors
19 Js <- c(10, 20, 50, 100, 200)
20 VarImportance <- matrix(0, length(Js), K)
21 l <- 1

```

### R code. Bayesian additive regression trees

```

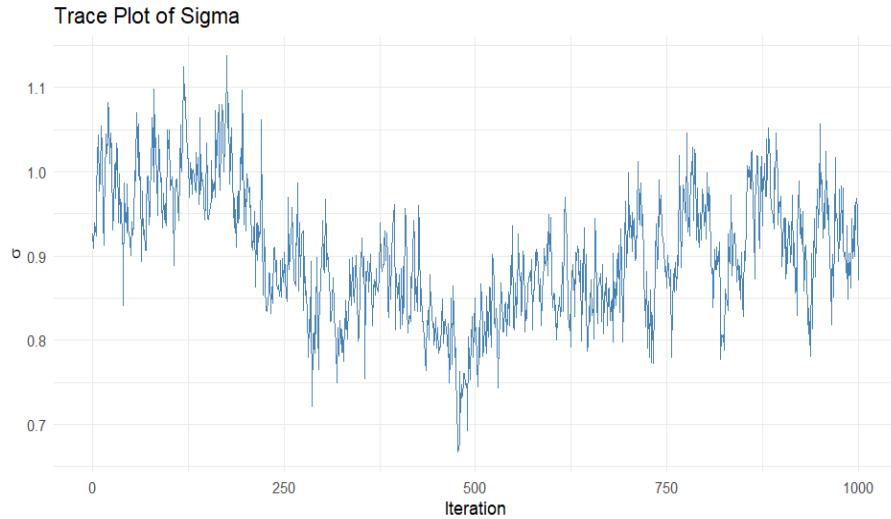
1 for (j in Js){
2   BARTfit <- wbart(x.train = X.train, y.train = y.train, x.
+     test = X.test, base = alpha,
3   power = beta, k = k, sigdf = v, sigquant = q, ntree = j,
4   ndpost = MCMCiter, nskip = burnin, keepevery = thinning)
5   VarImportance[1, ] <- BARTfit[["varcount.mean"]]/j
6   l <- l + 1
7 }
8 rownames(VarImportance) <- c("10", "20", "50", "100", "200")
9 colnames(VarImportance) <- as.character(1:10)
10 importance_df <- as.data.frame(VarImportance) %>% mutate(
+   trees = rownames(.)) %>%
11 pivot_longer(cols = -trees, names_to = "variable", values_to
+   = "percent_used")
12 importance_df$variable <- as.numeric(importance_df$variable)
13 importance_df$trees <- factor(importance_df$trees, levels =
+   c("10", "20", "50", "100", "200"))
14 ggplot(importance_df, aes(x = variable, y = percent_used,
+   color = trees, linetype = trees)) + geom_line() + geom_
+   point() + scale_color_manual(values = c("10" = "red", "
+   20" = "green", "50" = "blue", "100" = "cyan", "200" = "
+   magenta")) + scale_x_continuous(breaks = 1:10) + labs(x
+   = "variable", y = "percent used", color = "#trees",
+   linetype = "#trees") + theme_minimal()

```

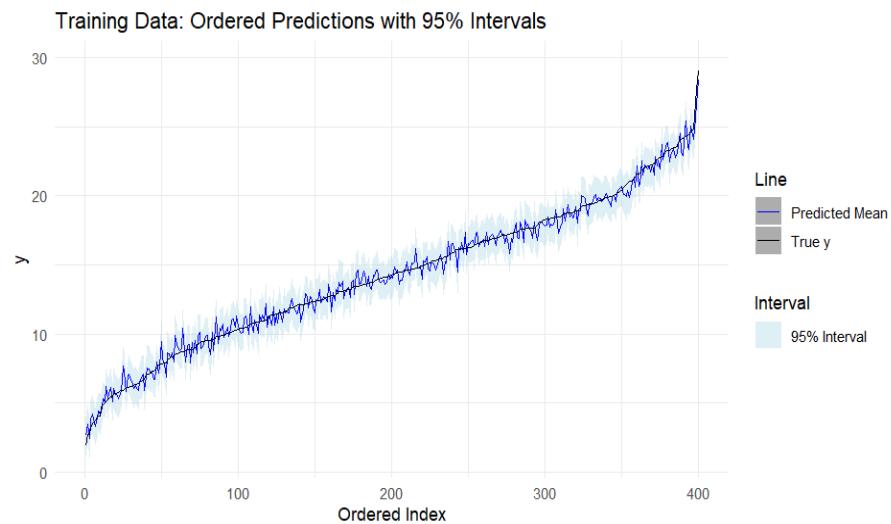
Figure 12.3 displays the trace plot of  $\sigma$ , which appears to have reached a stationary distribution. However, the posterior draws are slightly lower than the true value (1).

Figures 12.4 and 12.5 compare the true values of  $y$  with the posterior mean and 95% predictive intervals. The BART model performs well in both sets, and the intervals in the test set are notably wider than those in the training set.

Figure 12.6 shows the relative frequency with which each variable is used in the trees, a measure of variable relevance. When the number of trees is small, the algorithm more clearly identifies the most relevant predictors (variables 1–5). As the number of trees increases, this discrimination gradually disappears.

**FIGURE 12.3**

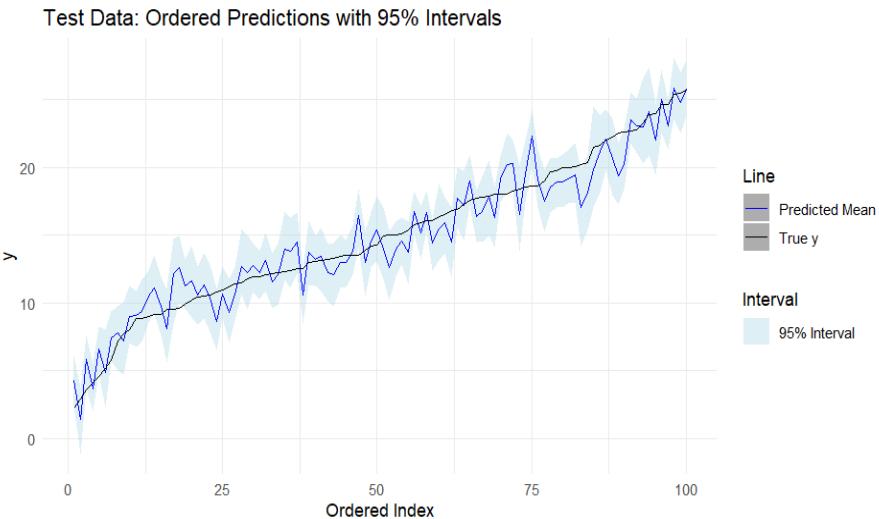
Trace plot of the posterior draws of  $\sigma$  in the BART simulation.

**FIGURE 12.4**

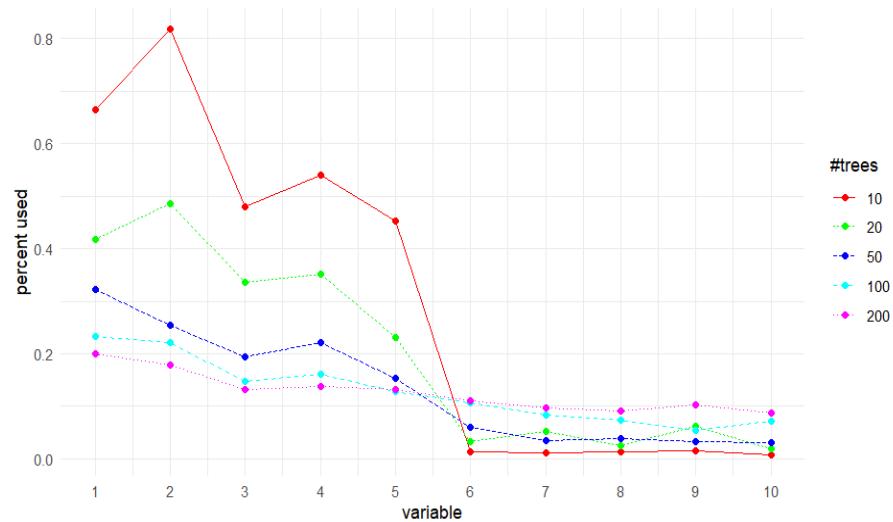
Training data: true  $y$  versus posterior mean and 95% predictive intervals.

## 12.4 Gaussian process

A Gaussian Process (GP) is an infinite collection of random variables, any finite subset of which follows a joint Gaussian distribution. A GP is fully

**FIGURE 12.5**

Test data: true  $y$  versus posterior mean and 95% predictive intervals.

**FIGURE 12.6**

Relative frequency of variable usage in the BART trees for different numbers of trees.

specified by its mean function and covariance function, that is,

$$f(\mathbf{x}) \sim GP(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')),$$

where  $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$  and  $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$ . It is common to assume  $m(\mathbf{x}) = 0$  to simplify calculations, although this is not required.

Perhaps the most commonly used covariance function in Gaussian Processes is the *squared exponential kernel* (or *radial basis function*) [220], defined as

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2l^2}\|\mathbf{x} - \mathbf{x}'\|^2\right),$$

where  $\sigma_f^2$  is the signal variance, which controls the vertical variation (amplitude) of the function,  $l$  is the length-scale parameter, which determines how quickly the function varies with features distance, and  $\|\mathbf{x} - \mathbf{x}'\|^2$  is the squared Euclidean distance between the feature vectors  $\mathbf{x}$  and  $\mathbf{x}'$ .

The squared exponential kernel implies that the function is infinitely differentiable, leading to very smooth function draws. While this smoothness may be desirable in some applications, it can be too restrictive in others. Alternative kernels like the Matérn class allow for more flexibility by controlling the degree of differentiability [364].

A GP can be interpreted as a prior distribution over a space of functions. The starting point in working with GPs is the specification of this prior before any data are observed. The following code illustrates five sample paths drawn from a GP with a squared exponential kernel, assuming a signal variance  $\sigma_f^2 = 1$  and a length-scale  $l = 0.2$ , evaluated over a grid of input values  $x \in [0, 1]$ . A small *jitter term* is added to the covariance matrix to ensure numerical stability during simulation. Figure 12.7 displays the five realizations drawn from the Gaussian Process.

### R code. Simulation of five realizations: Gaussian process

```

1 rm(list = ls()); set.seed(10101)
2 library(ggplot2); library(dplyr)
3 library(tidyr); library(MASS)
4 # Simulation setup
5 n <- 100
6 x <- seq(0, 1, length.out = n)
7 sigma_f <- 1; l <- 0.2; sigma_n <- 1e-8
8 # Squared Exponential Kernel function
9 SE_kernel <- function(x1, x2, sigma_f, l) {
10   outer(x1, x2, function(a, b) sigma_f^2 * exp(-0.5 * (a - b)^2 / l^2))
11 }
```

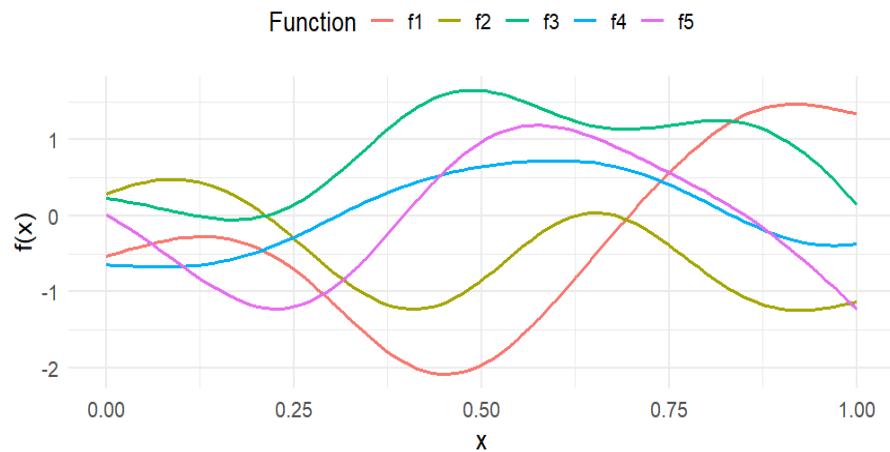
**R code. Simulation of five realizations: Gaussian process**

```

1 K <- SE_kernel(x, x, sigma_f, 1) + diag(sigma_n, n)
2 samples <- mvtnorm(n = 5, mu = rep(0, n), Sigma = K)
3 # Transpose and rename columns to f1, f2, ..., f5
4 samples_t <- t(samples)
5 colnames(samples_t) <- paste0("f", 1:5)
6 # Convert to tidy data frame
7 df <- data.frame(x = x, samples_t) |>
8 pivot_longer(cols = -x, names_to = "draw", values_to =
  "value")
9 # Plot
10 ggplot(df, aes(x = x, y = value, color = draw)) + geom_line(
  linewidth = 1) +
11 labs( title = "Simulated Gaussian Process Draws", x = "x",
  y = "f(x)", color = "Function" ) + theme_minimal(base_
  size = 14) + theme(legend.position = "top")
12

```

Simulated Gaussian Process Draws



**FIGURE 12.7**  
Simulation: Gaussian process.

Thus, for any finite set of feature points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ , the corresponding

function values follow a multivariate Gaussian distribution:

$$\mathbf{f} = \begin{bmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_N) \end{bmatrix} \sim N(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X})),$$

where the  $(i, j)$ -th entry of the covariance matrix  $\mathbf{K}(\mathbf{X}, \mathbf{X})$  is given by  $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ .

If we are interested in the properties of a function evaluated at a finite set of input points  $\{(f_i, x_i)\}_{i=1}^N$ , inference can be performed using only those points, effectively disregarding the uncountably infinite values the function may take elsewhere.

The following code illustrates how to perform inference for a GP given four observed points  $\{(f_i, x_i)\}_{i=1}^4$ , assuming that the true underlying process is

$$f_i = \sin(2\pi x_i).$$

The inference is based on the properties of the conditional Gaussian distribution (see below). Figure 12.8 shows that the posterior mean (solid blue line) interpolates the observed points (red dots). Moreover, the level of uncertainty (light blue shaded area) increases in regions that are farther from the observed inputs, where the posterior mean tends to deviate more from the true underlying function (dashed green line).

In situations where the input locations can be selected, such as in experimental designs, *active learning strategies* can be employed to choose the points that minimize predictive uncertainty. This is typically achieved by optimizing an *acquisition function* that quantifies the expected informativeness of candidate locations [408].

Consequently, GPs play a central role in *Bayesian optimization*, a stochastic method for finding the maximum of expensive or unknown objective functions. In this approach, a prior is placed over the objective function, which is then updated using observed data to form a posterior distribution over possible functions. This posterior guides the selection of new input points by balancing exploration and exploitation through the acquisition function [63].

**R code. Simulation:  $f_i = \sin(2\pi x_i)$  and Gaussian process**

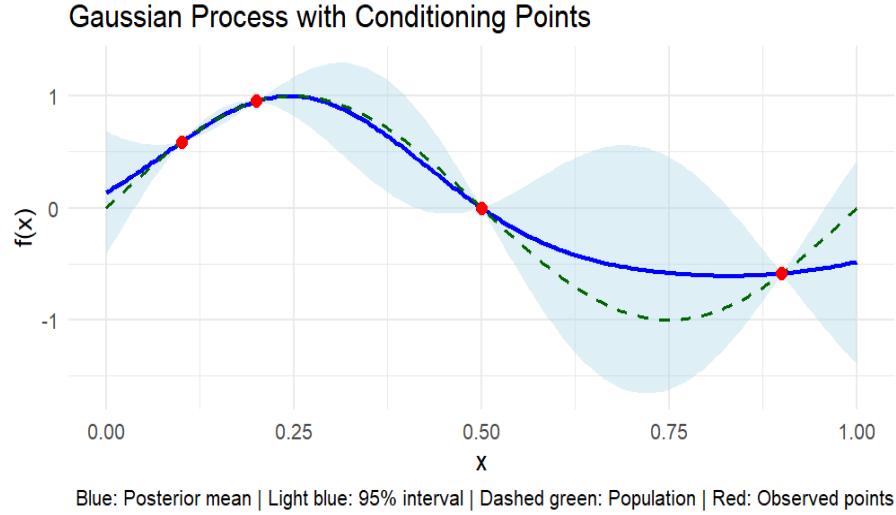
```
1 rm(list = ls()); set.seed(10101)
2 library(ggplot2); library(MASS)
```

**R code. Simulation:  $f_i = \sin(2\pi x_i)$  and Gaussian process**

```

1 # Define the squared exponential kernel
2 SE_kernel <- function(x1, x2, sigma_f, l) {
3   outer(x1, x2, function(a, b) sigma_f^2 * exp(-0.5 * (a - b
4     )^2 / l^2))
5 }
6 # Define the input space and observed points
7 x_star <- seq(0, 1, length.out = 200)
8 x0 <- c(0.1, 0.2, 0.5, 0.9)
9 y0 <- sin(2 * pi * x0)
10 # Hyperparameters
11 sigma_f <- 1
12 l <- 0.2
13 sigma_n <- 1e-8 # Jitter term for stability
14 # Compute covariance matrices
15 K_x0x0 <- SE_kernel(x0, x0, sigma_f, l) + diag(sigma_n,
16   length(x0))
17 K_xstarx0 <- SE_kernel(x_star, x0, sigma_f, l)
18 K_xstarxstar <- SE_kernel(x_star, x_star, sigma_f, l) + diag
19   (sigma_n, length(x_star))
20 # Compute posterior mean and covariance
21 K_inv <- solve(K_x0x0)
22 posterior_mean <- K_xstarx0 %*% K_inv %*% y0
23 posterior_cov <- K_xstarxstar - K_xstarx0 %*% K_inv %*% t(K_
24   xstarx0)
25 # Sample from the posterior
26 sample_draw <- sin(2 * pi * x_star)
27 # Compute 95% intervals
28 posterior_sd <- sqrt(diag(posterior_cov))
29 lower <- posterior_mean - 1.96 * posterior_sd
30 upper <- posterior_mean + 1.96 * posterior_sd
31 # Data frame for plotting
32 df <- data.frame(
33   x = x_star,
34   mean = posterior_mean,
35   lower = lower,
36   upper = upper,
37   sample = sample_draw
38 )
39 obs <- data.frame(x = x0, y = y0)
40 # Plot
41 ggplot(df, aes(x = x)) + geom_ribbon(aes(ymin = lower, ymax
42   = upper), fill = "lightblue", alpha = 0.4) + geom_line(
43   aes(y = mean), color = "blue", linewidth = 1.2) + geom_
44   line(aes(y = sample), color = "darkgreen", linewidth =
45   1, linetype = "dashed") + geom_point(data = obs, aes(x =
46   x, y = y), color = "red", size = 3) + labs( title = "
47   Gaussian Process with Conditioning Points", x = "x", y =
48   "f(x)", caption = "Blue: Posterior mean | Light blue:
49   95% interval | Dashed green: Population | Red: Observed
50   points" ) + theme_minimal(base_size = 14)

```



**FIGURE 12.8**  
Simulation:  $f_i = \sin(2\pi x_i)$  and Gaussian process

In practice, we have an observed dataset  $\{(y_i, \mathbf{x}_i)\}_{i=1}^N$  such that

$$y_i = f(\mathbf{x}_i) + \mu_i,$$

where  $\mu_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$ . This means that  $y_i$  is a noisy observation of  $f(\mathbf{x}_i)$ .

Thus, the marginal distribution of the observed outputs is

$$\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N),$$

where  $\mathbf{K}(\mathbf{X}, \mathbf{X})$  is the covariance matrix generated by the GP kernel evaluated at the training inputs.

Note that this implies the log marginal likelihood is given by

$$\log p(\mathbf{y} \mid \mathbf{X}) = -\frac{1}{2}\mathbf{y}^\top (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K} + \sigma^2 \mathbf{I}_N| - \frac{N}{2} \log 2\pi.$$

We can adopt an empirical Bayes approach to estimate the hyperparameters of the GP prior by maximizing the log marginal likelihood with respect to the kernel parameters (e.g.,  $\sigma_f^2$ ,  $l$ ) and the noise variance  $\sigma^2$ .

To make predictions at a new set of features  $\mathbf{X}_*$ , we consider the joint distribution:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}) & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right).$$

Using the conditional distribution of a multivariate Gaussian, the *posterior predictive distribution* [364] is:

$$\mathbf{f}_* \mid \mathbf{y} \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)),$$

where

$$\bar{\mathbf{f}}_* = \mathbb{E}[\mathbf{f}_* | \mathbf{y}, \mathbf{X}, \mathbf{X}_*] = \mathbf{K}(\mathbf{X}_*, \mathbf{X})[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) - \mathbf{K}(\mathbf{X}_*, \mathbf{X})[\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N]^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}_*).$$

Therefore, GP regression provides a flexible and efficient nonparametric framework for predicting unobserved responses, with accuracy that improves as more data become available. GPs are widely used due to their favorable computational properties, including the availability of closed-form expressions, and posterior consistency under mild conditions [101, 426]. Moreover, predictive performance can be further enhanced by incorporating derivative information, as the derivative of a GP is itself a GP [417, 220].

However, a major limitation of GPs is the need to invert an  $N \times N$  covariance matrix, which requires  $O(N^3)$  computational operations, making them computationally expensive for large datasets. To address this, several scalable methods have been proposed that reduce the computational burden. For instance, [451, 154, 338] develop algorithms that reduce complexity to  $O(N)$ .

#### **Example: Simulation exercise to study GP performance**

We simulate the process

$$f_i = \sin(2\pi x_{i1}) + \cos(2\pi x_{i2}) + \sin(x_{i1}x_{i2}),$$

where  $x_{i1}$  and  $x_{i2}$  are independently drawn from a uniform distribution on the interval  $[0, 1]$ , for  $i = 1, 2, \dots, 100$ .

We use the *DiceKriging* package in **R** to estimate and make predictions using a GP. This package applies maximum likelihood estimation to infer the length-scale parameters ( $l_k$ ) and the signal variance ( $\sigma_f^2$ ). Note that there are two separate length-scale parameters, one for each input variable. The following code illustrates how to carry out this example, and Figure 12.9 displays a 3D plot with the observed points and the posterior mean surface. The package also provides pointwise credible intervals for the predictions.

***R code. Simulation:***  
 $f_i = \sin(2\pi x_{i1}) + \cos(2\pi x_{i2}) + \sin(x_{i1}x_{i2})$  ***and Gaussian process***

---

```

1 # Load required packages
2 library(DiceKriging); library(rgl)
3 # Simulate training data
4 set.seed(10101); n_train <- 100

```

---

**R code. Simulation:**

$f_i = \sin(2\pi x_{i1}) + \cos(2\pi x_{i2}) + \sin(x_{i1}x_{i2})$  and Gaussian process

```

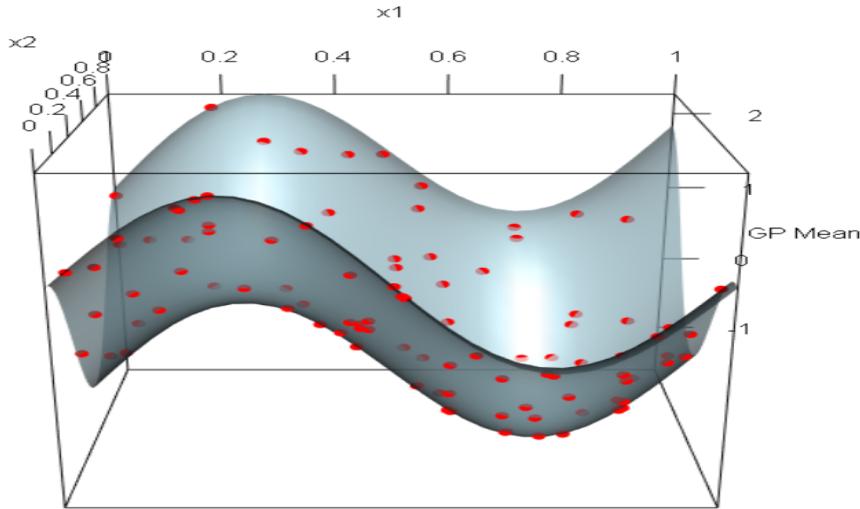
1 x1 <- runif(n_train); x2 <- runif(n_train)
2 X_train <- data.frame(x1 = x1, x2 = x2)
3 # True function without noise
4 f_train <- sin(2 * pi * X_train$x1) + cos(2 * pi * X_train$ 
      x2) + sin(X_train$x1 * X_train$x2)
5 # Fit Gaussian Process
6 fit_km <- km(design = X_train, response = f_train, covtype =
      "gauss", nugget = 1e-10)
7 # Prediction grid
8 grid_points <- 30
9 x1_seq <- seq(0, 1, length.out = grid_points)
10 x2_seq <- seq(0, 1, length.out = grid_points)
11 grid <- expand.grid(x1 = x1_seq, x2 = x2_seq)
12 # Predict GP surface
13 pred <- predict(fit_km, newdata = grid, type = "UK")
14 z_pred <- matrix(pred$mean, nrow = grid_points, ncol = grid_
      points)
15 # Plot
16 persp3d(x = x1_seq, y = x2_seq, z = z_pred, col = "lightblue"
      , alpha = 0.7, xlab = "x1", ylab = "x2", zlab = "GP
      Mean")
17 points3d(x = X_train$x1, y = X_train$x2, z = f_train, col =
      "red", size = 8)
18 fit_km@covariance@range.val # length-scale
19 fit_km@covariance@sd2 # Signal variance

```

A limitation of the *DiceKriging* package is that it is designed for deterministic simulations and, consequently, does not estimate the noise variance. Therefore, in Exercise 7, we ask to simulate the process

$$f_i = \sin(2\pi x_{i1}) + \cos(2\pi x_{i2}) + \sin(x_{i1}x_{i2}) + \mu_i,$$

where  $\mu_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 0.1^2)$ , and to use an empirical Bayes approach to estimate the hyperparameters. These estimated hyperparameters should then be used to perform GP prediction.

**FIGURE 12.9**

Simulation:  $f_i = \sin(2\pi x_{i1}) + \cos(2\pi x_{i2}) + \sin(x_{i1}x_{i2})$  and Gaussian process

## 12.5 Tall data problems

In this section, we explore several methods developed to perform Bayesian inference when the sample size is large, particularly when there is a large number of observational units, commonly referred to as *tall datasets*.

Bayesian inference in such settings is computationally demanding because each iteration of an MCMC algorithm requires evaluating the likelihood function over all  $N$  observations. For large  $N$ , this renders standard MCMC methods prohibitively expensive. Recent efforts have focused on developing scalable Monte Carlo algorithms that significantly reduce the computational cost compared to standard approaches. One alternative is to use *Variational Bayes* (see Chapter 14); however, it can be challenging to implement and may exhibit limitations in uncertainty quantification, particularly for the joint posterior distribution. Another alternative is the *Integrated Nested Laplace Approximation* (INLA, see Chapter 14); however, its computational cost grows exponentially with the dimension of the parameter space.

In scenarios where observations are assumed to be independent, two main frameworks have been proposed to scale MCMC algorithms: *divide-and-conquer approaches* and *subsampling-based algorithms*. Divide-and-conquer methods partition the dataset into disjoint subsets, run MCMC independently on each batch to obtain subposterior distributions, and then combine them to approximate the full posterior distribution. Subsampling-based algorithms, on the other

hand, aim to reduce the number of data points used to evaluate the likelihood at each iteration, often relying on *pseudo-marginal* MCMC methods [15]. The key idea of pseudo-marginal MCMC is to augment the model with latent variables such that the sample average of the likelihood, computed over draws from these latent variables, provides an unbiased estimator of the marginal likelihood. This approach is particularly valuable when the marginal likelihood is not available in closed form. Moreover, the same principles can be adapted to reduce the computational burden of evaluating the log-likelihood. For an excellent review of divide-and-conquer and subsampling-based approaches, see [27].

### 12.5.1 Divide-and-conquer methods

In *divide-and-conquer* methods, the main idea is to partition the dataset and distribute the subsets across multiple computing machines/cores. An independent MCMC algorithm is then executed on each subset to obtain a corresponding *subposterior* distribution. The central challenge lies in accurately and efficiently combining these subposteriors into a single approximation of the full posterior distribution. Several approaches have been proposed to address this issue. For instance, [212, 401, 366, 402] introduce the *Consensus Monte Carlo* algorithm; [446] develop a method based on the *Weierstrass sampler* for parallelizing MCMC; [310] propose using the *geometric median of posterior distributions*; and [460] suggest combining *rescaled subposteriors*.

In divide-and-conquer methods, the dataset is partitioned into  $B$  disjoint batches  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_B$ , and the posterior is rewritten using the identity:

$$\pi(\boldsymbol{\theta} | \mathbf{y}) \propto \prod_{b=1}^B \pi(\boldsymbol{\theta})^{1/B} p(\mathbf{y}_b | \boldsymbol{\theta}),$$

which implies that the full posterior is proportional to the product of appropriately rescaled subposteriors.

*Consensus Monte Carlo* (CMC) operates by running separate Monte Carlo algorithms on each subset in parallel, and then averaging the resulting posterior draws. Specifically, given samples  $\boldsymbol{\theta}_b^{(s)}$ , for  $b = 1, 2, \dots, B$  and  $s = 1, 2, \dots, S$ , obtained independently from each batch  $\mathbf{y}_b$ , the  $s$ -th draw from the consensus posterior is computed as:

$$\boldsymbol{\theta}^{(s)} = \left( \sum_{b=1}^B \mathbf{w}_b \right)^{-1} \sum_{b=1}^B \mathbf{w}_b \boldsymbol{\theta}_b^{(s)},$$

where the optimal weight is the inverse covariance matrix of the subposterior, i.e.,  $\mathbf{w}_b = \text{Var}^{-1}(\boldsymbol{\theta} | \mathbf{y}_b)$ . In practice, one may use the marginal variances of each parameter to simplify the computation, which can still yield good performance.

When each subposterior  $\pi_b(\boldsymbol{\theta} | \mathbf{y}_b)$  is Gaussian, the full posterior  $\pi(\boldsymbol{\theta} | \mathbf{y})$

is also Gaussian, and can be recovered exactly by combining the subposterior distributions using simple rules based on their means and covariances [401, 402]. In the non-Gaussian case, standard asymptotic results in Bayesian inference (see Chapter 1) imply that the posterior distributions converge to a Gaussian distribution as the batch size increases. Alternative merging procedures that are more robust to non-Gaussianity have also been proposed [319, 311]; however, it remains difficult to quantify the approximation error in these approaches. Moreover, this procedure is limited to continuous parameter spaces and may exhibit small-sample bias; that is, when the dataset is divided into small batches, the subposterior distributions may be biased. In such cases, jackknife bias corrections are recommended to reduce the overall approximation error.

In particular, we perform CMC following Algorithm A33 [401]. Next, we compute the CMC posterior repeatedly, each time leaving out one of the  $B$  subsets. Let  $\pi_{-b}(\boldsymbol{\theta} | \mathbf{y})$  denote the resulting posterior when subset  $b$  is excluded. The average of these leave-one-out posteriors is denoted by

$$\bar{\pi}_{-b}(\boldsymbol{\theta} | \mathbf{y}) = \frac{1}{B} \sum_{b=1}^B \pi_{-b}(\boldsymbol{\theta} | \mathbf{y}).$$

Then, the jackknife bias-corrected posterior is given by

$$\pi_{jk}(\boldsymbol{\theta} | \mathbf{y}) = B \cdot \pi_{CMC}(\boldsymbol{\theta} | \mathbf{y}) - (B - 1) \cdot \bar{\pi}_{-b}(\boldsymbol{\theta} | \mathbf{y}),$$

where  $\pi_{CMC}(\boldsymbol{\theta} | \mathbf{y})$  is the original CMC posterior based on all  $B$  subsets.

---

**Algorithm A33** Consensus Monte Carlo algorithm

---

- 1: Divide the dataset into  $B$  disjoint batches  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_B$
  - 2: Run  $B$  separate MCMC algorithms to sample  $\boldsymbol{\theta}_b^{(s)} \sim \pi(\boldsymbol{\theta} | \mathbf{y}_b)$ ,  $b = 1, 2, \dots, B$ , and  $s = 1, 2, \dots, S$  using the prior distribution  $\pi(\boldsymbol{\theta})^{1/B}$
  - 3: Combine the posterior draws using  $\boldsymbol{\theta}^{(s)} = \left( \sum_{b=1}^B \mathbf{w}_b \right)^{-1} \sum_{b=1}^B \mathbf{w}_b \boldsymbol{\theta}_b^{(s)}$  using  $\mathbf{w}_b = \text{Var}^{-1}(\boldsymbol{\theta} | \mathbf{y}_b)$
- 

The main difficulty is how to effectively merge the subposterior distributions, especially when their supports are not well-aligned. This misalignment can lead to poor scalability with an increasing number of batches. Moreover, most theoretical guarantees for these methods are asymptotic in the size of each batch, which may limit their performance in practice [27].

### 12.5.2 Subsampling-based algorithms

An alternative to divide-and-conquer methods is to avoid evaluating the likelihood over all observations, which requires  $O(N)$  operations. Instead, the likelihood is approximated using a smaller subset of observations,  $n \ll N$ ,

in order to reduce the computational burden of the algorithm. The starting point is the log-likelihood function for  $N$  independent observations:

$$\log p(\mathbf{y} \mid \boldsymbol{\theta}) = \sum_{i=1}^N \log p(y_i \mid \boldsymbol{\theta}).$$

The literature has focused on the log-likelihood because it is a sum over independent contributions, which is analogous to the problem of estimating a population total.

A class of subsampling methods relies on estimating the marginal likelihood via the *pseudo-marginal* approach. Examples include the confidence sampler [26], the Firefly Monte Carlo algorithm [284], whose relationship to subsampling MCMC is formally established in [27], and approaches using data-expanded and parameter-expanded control variates [343].

The intuition behind the pseudo-marginal method is straightforward: introduce a set of auxiliary random variables  $\mathbf{z} \sim p(\mathbf{z})$ , such that the marginal likelihood can be written as an expectation with respect to  $\mathbf{z}$ :

$$\mathbb{E}_{\mathbf{z}}[p(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{z})] = \int_{\mathcal{Z}} p(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{z}) p(\mathbf{z}) d\mathbf{z} = p(\mathbf{y} \mid \boldsymbol{\theta}).$$

This implies that

$$\hat{p}(\mathbf{y} \mid \boldsymbol{\theta}) = \frac{1}{S} \sum_{s=1}^S p(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{z}^{(s)})$$

is an unbiased estimator of the marginal likelihood,  $\mathbf{z}^{(s)} \sim p(\mathbf{z})$ .

As a result, the pseudo-marginal method enables exact simulation-based inference for  $p(\mathbf{y} \mid \boldsymbol{\theta})$  in settings where the likelihood cannot be evaluated analytically [15], for instance, in non-linear random effects models (see also approximate methods such as *approximate Bayesian computation* and *Bayesian synthetic likelihood* in Chapter 14). [15] show that replacing the likelihood with an unbiased and positive estimator within the Metropolis–Hastings (MH) algorithm yields samples from the correct posterior distribution.

The pseudo-marginal likelihood approach can also be applied in settings where the sample size is so large that evaluating the full likelihood at each iteration of an MCMC algorithm becomes computationally prohibitive. In such cases, the likelihood can be approximated using a small subset of observations,  $n \ll N$ . The choice of the subset size  $n$  is particularly important, as it directly affects the variance of the likelihood estimator, which in turn is critical to ensuring an efficient Metropolis–Hastings (MH) algorithm.

In particular, a likelihood estimator with high variance may result in an accepted draw that overestimates the likelihood. As a consequence, subsequent proposals are unlikely to be accepted, causing the algorithm to become stuck and leading to a very low acceptance rate. Therefore, the choice of  $n$  determines the computational efficiency of the algorithm: a small  $n$  increases

the estimator's variance, which reduces the acceptance rate, whereas a large  $n$  increases the number of likelihood evaluations per iteration.

[344] recommend targeting a likelihood estimator variance between 1 and 3.3 to optimize computational efficiency, as supported by the findings of [337].

Let  $\ell_i(y_i | \boldsymbol{\theta}) = \log p(y_i | \boldsymbol{\theta})$  denote the contribution of the  $i$ -th observation to the log-likelihood, and let  $z_1, \dots, z_N$  be latent binary variables such that  $z_i = 1$  indicates that  $y_i$  is included in a subsample of size  $n$ , selected without replacement. Then, an unbiased estimator of the log-likelihood is given by

$$\hat{\ell}(\mathbf{y} | \boldsymbol{\theta}) = \frac{N}{n} \sum_{i=1}^N \ell_i(y_i | \boldsymbol{\theta}) z_i.$$

However, note that what we require is an unbiased estimator of the likelihood, not the log-likelihood. Consequently, a bias correction is needed:

$$\hat{p}(\mathbf{y} | \boldsymbol{\theta}) = \exp \left\{ \hat{\ell}(\mathbf{y} | \boldsymbol{\theta}) - \frac{1}{2} \sigma_{\hat{\ell}}^2(\boldsymbol{\theta}) \right\},$$

where  $\sigma_{\hat{\ell}}^2(\boldsymbol{\theta})$  denotes the variance of  $\hat{\ell}(\mathbf{y} | \boldsymbol{\theta})$  [75]. This correction is exact if  $\sigma_{\hat{\ell}}^2(\boldsymbol{\theta})$  is known and  $\hat{\ell}(\mathbf{y} | \boldsymbol{\theta})$  follows a normal distribution.

Given the importance of controlling the variance of the log-likelihood estimator in subsampling methods, and the limitations of simple random sampling in achieving low variability, [343] propose a highly efficient, unbiased estimator of the log-likelihood based on *control variates*, specifically through data-expanded and parameter-expanded control variates. The key idea is to construct a function  $q_i(\boldsymbol{\theta})$  that is highly correlated with the log-likelihood contribution  $\ell_i(y_i | \boldsymbol{\theta})$ , thereby stabilizing the log-likelihood estimator. In particular, [343] introduce a *difference estimator* of the form:

$$\hat{\ell}_{DE}(\mathbf{y} | \boldsymbol{\theta}, \mathbf{z}) = \sum_{i=1}^N q_i(\boldsymbol{\theta}) + \frac{N}{n} \sum_{i:z_i=1} (\ell_i(y_i | \boldsymbol{\theta}) - q_i(\boldsymbol{\theta})).$$

This estimator  $\hat{\ell}_{DE}(\mathbf{y} | \boldsymbol{\theta}, \mathbf{z})$  is unbiased for the full log-likelihood  $\log p(\mathbf{y} | \boldsymbol{\theta})$ .

[343] propose constructing  $q_i(\boldsymbol{\theta})$  using a second-order Taylor expansion of the log-likelihood around a central value of  $\boldsymbol{\theta}$ , such as the mode. An alternative approach is to perform a second-order Taylor expansion around the nearest centroid of each observation, where the centroids are obtained from a pre-clustering of the data.

The first approach can perform poorly when the current draw  $\boldsymbol{\theta}$  is far from the central expansion point, leading to inaccurate approximations. The second approach encounters difficulties in high-dimensional settings due to the curse of dimensionality: many observations will be far from their assigned centroid. To address these issues, the authors propose an adaptive strategy: initialize the algorithm using data-expanded control variates, and switch to

parameter-expanded control variates once the sampler reaches a region closer to the center of the parameter space.

It is important to note that this strategy targets an approximation to the posterior distribution, due to the small bias introduced by the difference estimator. However, this bias diminishes rapidly and has a negligible impact on the quality of the posterior inference.

Once a good estimator of the log-likelihood is obtained, meaning it has low variance, the likelihood can be recovered using the appropriate bias correction. This corrected likelihood estimator is then used within the acceptance probability of the Metropolis–Hastings algorithm (see Section 4.1.2), resulting in the so-called pseudo-marginal Metropolis–Hastings (PMMH) method. This strategy significantly reduces computational cost in tall data settings.

Another class of subsampling methods, which does not rely on the pseudo-marginal likelihood, consists of stochastic gradient MCMC algorithms. These methods begin based on ideas from stochastic gradient descent [370] and the Langevin diffusion-based stochastic differential equations.

The point of departure is the unnormalized posterior distribution:

$$\begin{aligned}\pi(\boldsymbol{\theta} \mid \mathbf{y}) &\propto \pi(\boldsymbol{\theta}) \prod_{i=1}^N p(y_i \mid \boldsymbol{\theta}) \\ &= \exp \left\{ \sum_{i=1}^N \left[ \frac{1}{N} \log \pi(\boldsymbol{\theta}) + \log p(y_i \mid \boldsymbol{\theta}) \right] \right\} \\ &= \exp \left\{ - \sum_{i=1}^N U_i(\boldsymbol{\theta}) \right\} \\ &= \exp \{-U(\boldsymbol{\theta})\},\end{aligned}$$

where  $\boldsymbol{\theta} \in \mathbb{R}^K$ ,  $U_i(\boldsymbol{\theta}) = -\frac{1}{N} \log \pi(\boldsymbol{\theta}) - \log p(y_i \mid \boldsymbol{\theta})$ , and  $U(\boldsymbol{\theta}) = \sum_{i=1}^N U_i(\boldsymbol{\theta})$  is assumed to be continuous and differentiable almost everywhere.

The advantage of this formulation is that, under mild regularity conditions [374], the Langevin diffusion process

$$d\boldsymbol{\theta}(s) = -\frac{1}{2} \nabla U(\boldsymbol{\theta}(s)) ds + d\mathbf{B}_s,$$

has  $\pi(\boldsymbol{\theta} \mid \mathbf{y})$  as its stationary distribution. Here,  $\nabla U(\boldsymbol{\theta}(s))$  is the drift term, and  $\mathbf{B}_s$  is a  $K$ -dimensional Brownian motion.<sup>3</sup>

Using an Euler-Maruyama discretization of the Langevin diffusion gives a proposal draw from the posterior:

$$\boldsymbol{\theta}^{(c)} = \boldsymbol{\theta}^{(s)} - \frac{\epsilon}{2} \nabla U(\boldsymbol{\theta}^{(s)}) + \boldsymbol{\psi},$$

---

<sup>3</sup>A Brownian motion is a continuous-time stochastic process that starts at zero, has independent increments with  $B(s) - B(t) \sim N(0, s-t)$ , and is continuous almost surely but nowhere differentiable.

where  $\psi \sim \mathcal{N}(\mathbf{0}, \epsilon \mathbf{I}_K)$  and  $\epsilon > 0$  is a suitably chosen step size (learning rate). This proposal is used within a Metropolis–Hastings algorithm (see Section 4.1.2) to correct for the discretization error introduced by the Euler approximation. This method is known as the *Metropolis-adjusted Langevin algorithm* (MALA) [374].

A simpler variant, known as the *unadjusted Langevin algorithm* (ULA), omits the acceptance step. As a result, ULA produces a biased approximation of the posterior distribution. However, a major computational bottleneck in both MALA and ULA is the requirement to evaluate the full gradient  $\nabla U(\boldsymbol{\theta}) = \sum_{i=1}^N \nabla U_i(\boldsymbol{\theta})$  at every iteration, which becomes computationally prohibitive when  $N$  is large.

To overcome this limitation, [449] proposed the *Stochastic Gradient Langevin Dynamics* (SGLD), which replaces the full gradient with an unbiased estimate computed using a mini-batch of data. Given a random sample of size  $n \ll N$ , the stochastic gradient estimate at iteration  $s$  is:

$$\hat{\nabla}U(\boldsymbol{\theta})^{(n)} = \frac{N}{n} \sum_{i \in \mathcal{S}_n} \nabla U_i(\boldsymbol{\theta}), \quad (12.1)$$

where  $\mathcal{S}_n \subset \{1, 2, \dots, N\}$  is a randomly selected subset of dimension  $n$  without replacement.

Therefore,

$$\boldsymbol{\theta}^{(s+1)} = \boldsymbol{\theta}^{(s)} - \frac{\epsilon_s}{2} \hat{\nabla}U(\boldsymbol{\theta}^{(s)})^{(n)} + \psi_s,$$

such that  $\sum_{s=1}^{\infty} \epsilon_s = \infty$  and  $\sum_{s=1}^{\infty} \epsilon_s^2 < \infty$ . These conditions guarantee almost sure convergence: the former condition allows the algorithm to continue exploring the parameter space (no premature convergence), and the latter ensures that the cumulative noise does not explode.

[430] formally show that, under verifiable assumptions, the SGLD algorithm is consistent. That is, given a test function  $\phi(\boldsymbol{\theta}) : \mathbb{R}^K \rightarrow \mathbb{R}$ ,

$$\lim_{S \rightarrow \infty} \frac{\epsilon_1 \phi(\boldsymbol{\theta}_1) + \epsilon_2 \phi(\boldsymbol{\theta}_2) + \dots + \epsilon_S \phi(\boldsymbol{\theta}_S)}{\sum_{s=1}^S \epsilon_s} = \int_{\mathbb{R}^K} \phi(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

Moreover, the algorithm satisfies a central limit theorem:  $\lim_{S \rightarrow \infty} \pi_S(\phi(\boldsymbol{\theta})) = \pi(\phi(\boldsymbol{\theta}))$ , and its asymptotic bias-variance decomposition is characterized by a functional of  $\epsilon_s$ , such that the optimal step size that minimizes the asymptotic mean squared error is proportional to  $s^{-1/3}$ . In the common practice of using a constant step size, it has been shown that the optimal choice to minimize the asymptotic mean squared error is of order  $S^{-1/3}$  [442]. However, we recommend tuning this parameter based on the specific application, guided by the theoretical results presented here.

Importantly, this iterative process does not require the computation of acceptance probabilities, which significantly reduces the computational burden. Empirical evidence suggests that SGLD often outperforms the Metropolis–Hastings algorithm when applied to large datasets under a fixed computational budget [268].

Algorithm A34 summarizes the SGLD procedure [321].

---

**Algorithm A34** Stochastic gradient Langevin dynamic algorithm

---

- 1: Set  $\boldsymbol{\theta}^{(0)}$  and the step size schedule  $\epsilon_s$
  - 2: **for**  $s = 1, \dots, S$  **do**
  - 3:     Draw  $\mathcal{S}_n$  of size  $n$  from  $i = \{1, 2, \dots, N\}$  without replacement
  - 4:     Calculate  $\hat{\nabla}U(\boldsymbol{\theta})^{(n)}$  using Equation 12.1
  - 5:     Draw  $\boldsymbol{\psi}_s \sim N(\mathbf{0}, \epsilon_s \mathbf{I}_K)$
  - 6:     Update  $\boldsymbol{\theta}^{(s+1)} \leftarrow \boldsymbol{\theta}^{(s)} - \frac{\epsilon_s}{2} \hat{\nabla}U(\boldsymbol{\theta}^{(s)})^{(n)} + \boldsymbol{\psi}_s$
  - 7: **end for**
- 

A critical component of the SGLD algorithm is the estimation of the stochastic gradient (Equation 12.1), particularly because high variability in this estimator can lead to algorithmic instability, a challenge also encountered in pseudo-marginal methods, as described previously. To mitigate this issue, the literature also employs *control variates* to reduce the variance of the estimator. The core idea is to construct a simple function  $u_i(\boldsymbol{\theta})$  that is highly correlated with  $\nabla U_i(\boldsymbol{\theta})$  and has a known expectation. This correlation allows the fluctuations in  $u_i(\boldsymbol{\theta})$  to “cancel out” some of the noise in  $\nabla U_i(\boldsymbol{\theta})$ , thereby stabilizing the stochastic gradient estimates. Specifically,

$$\sum_{i=1}^N \nabla U_i(\boldsymbol{\theta}) = \sum_{i=1}^N u_i(\boldsymbol{\theta}) + \sum_{i=1}^N (\nabla U_i(\boldsymbol{\theta}) - u_i(\boldsymbol{\theta})),$$

which leads to the following unbiased estimator:

$$\sum_{i=1}^N u_i(\boldsymbol{\theta}) + \frac{N}{n} \sum_{i \in \mathcal{S}_n} (\nabla U_i(\boldsymbol{\theta}) - u_i(\boldsymbol{\theta})).$$

To construct effective control variates, one common strategy is to first approximate the posterior mode  $\hat{\boldsymbol{\theta}}$  using stochastic gradient descent (SGD), which serves as the initialization point for Algorithm A34. SGD proceeds via a stochastic approximation of the gradient:

$$\boldsymbol{\theta}^{(s+1)} = \boldsymbol{\theta}^{(s)} - \epsilon_s \frac{1}{n} \sum_{i \in \mathcal{S}_n} \nabla U_i(\boldsymbol{\theta}^{(s)}).$$

This approximation introduces stochasticity into the updates but significantly reduces computational cost.

Two commonly used learning rate (or step size) schedules are  $\epsilon_s = s^{-\kappa}$  and  $\epsilon_s = \epsilon_0 / (1 + s/\tau)^\kappa$ , where  $\epsilon_0$  is the initial learning rate,  $\tau$  is a stability constant that slows down early decay (larger values lead to more stable early behavior), and  $\kappa \in (0.5, 1]$  controls the long-run decay rate. If  $\kappa$  is too large, the learning rate decays too quickly and the algorithm may stagnate. Conversely, if  $\kappa$  is too small, the algorithm may remain unstable or fail to converge.

An important distinction to note is that SGLD operates with gradient *sums*, while SGD typically uses *averages*. This distinction affects how step sizes and noise scaling should be interpreted in practice.

After convergence, we obtain a reliable estimate of the posterior mode  $\hat{\boldsymbol{\theta}}$ . Based on this, we define the control variate as  $u_i(\boldsymbol{\theta}) = \nabla U_i(\hat{\boldsymbol{\theta}})$ . The resulting control variate estimator of the gradient is:

$$\hat{\nabla}_{\text{cv}} U(\boldsymbol{\theta}) = \sum_{i=1}^N \nabla U_i(\hat{\boldsymbol{\theta}}) + \frac{N}{n} \sum_{i \in \mathcal{S}_n} (\nabla U_i(\boldsymbol{\theta}) - \nabla U_i(\hat{\boldsymbol{\theta}})).$$

This expression is used in place of  $\hat{\nabla}U(\boldsymbol{\theta})^{(n)}$  from Equation 12.1 within Algorithm A34.

It has been shown that the computational cost per effective sample size of SGLD algorithm with control variates is  $O(1)$ , rather than  $O(N)$  [321]. However, it is also known that the approximate posterior distribution produced by SGLD, despite capturing the correct mode, tends to exhibit higher variance compared to the true target posterior.

[78] propose the stochastic gradient Hamiltonian Monte Carlo (SGHMC) algorithm, which performs better asymptotically than SGLD in terms of the Wasserstein distance to the target distribution, especially in high-dimensional settings (with large  $K$ ). Similarly, [283] propose an alternative diffusion process whose stationary distribution is the desired posterior. See Table 1 in [321] for a comprehensive list of stochastic gradient MCMC (SGMCMC) methods that fit in the general diffusion process proposed by [283].

[321] also discuss convergence diagnostics for SGMCMC methods. They argue that standard diagnostics, such as those discussed in Section 4.4, are inadequate for assessing SGMCMC behavior, as these methods target asymptotically biased approximations of the posterior. Consequently, a growing body of research focuses on designing suitable diagnostics for SGMCMC, with approaches based on Stein's discrepancy being among the most prominent (see [321] for details).

A limitation of SGMCMC algorithms is that they are applicable only to a restricted class of models, specifically, those with fixed-dimensional parameter spaces in  $\mathcal{R}^K$  and differentiable log-posterior densities. To broaden their applicability, [418] propose an extended stochastic gradient MCMC framework capable of addressing more general large-scale Bayesian inference problems, including those involving dimension-jumping and missing data.

### **Example: Simulation exercise to study the performance of CMC and SGLD**

In this example, we follow the logistic regression simulation setup introduced by [321]:

$$P(y_i = 1 \mid \boldsymbol{\beta}, \mathbf{x}_i) = \frac{\exp \{ \mathbf{x}_i^\top \boldsymbol{\beta} \}}{1 + \exp \{ \mathbf{x}_i^\top \boldsymbol{\beta} \}},$$

with log-likelihood function given by:

$$\begin{aligned}\log p(\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{x}) &= \sum_{i=1}^N y_i (\mathbf{x}_i^\top \boldsymbol{\beta} - \log(1 + \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\})) + (1 - y_i) (-\log(1 + \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\})) \\ &= \sum_{i=1}^N y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \log(1 + \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}).\end{aligned}$$

Note that this implies that the gradient vector is

$$\nabla \log p(\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{x}) = \sum_{i=1}^N \left( y_i - \frac{\exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}}{1 + \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}} \right) \mathbf{x}_i.$$

We assume a prior distribution for  $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, 10\mathbf{I}_K)$ , leading to the log-prior:

$$\log \pi(\boldsymbol{\beta}) = -\frac{K}{2} \log(2\pi) - \frac{1}{2} \log(|10\mathbf{I}_K|) - \frac{1}{2} \boldsymbol{\beta}^\top (10^{-1}\mathbf{I}_K) \boldsymbol{\beta}.$$

This means that the gradient vector of the log-prior is

$$\nabla \log \pi(\boldsymbol{\beta}) = -\frac{1}{10} \boldsymbol{\beta}.$$

Also note that

$$\begin{aligned}\pi(\boldsymbol{\beta})^{1/B} &\propto \left\{ \exp\left(-\frac{1}{2 \cdot 10} \boldsymbol{\beta}^\top \boldsymbol{\beta}\right) \right\}^{1/B} \\ &= \exp\left(-\frac{1}{2 \cdot 10 \cdot B} \boldsymbol{\beta}^\top \boldsymbol{\beta}\right).\end{aligned}$$

This implies that, when implementing CMC, the prior variance must be scaled by the number of batches  $B$ . In other words, each subposterior should use a prior with variance  $10 \cdot B$  in this example, so that the product of the  $B$  subposterior reconstrucst the correct full posterior.

We set  $K = 10$ ,  $\boldsymbol{\beta} = 0.5 \cdot \mathbf{i}_K$ , and  $N = 10^5$ . The covariates  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ , where the covariance matrix  $\boldsymbol{\Sigma}^{(i,j)} = U[-\rho, \rho]^{|i-j|}$  with  $\rho = 0.4$ , and  $\mathbf{i}_K$  denotes a  $K$ -dimensional vector of ones.

We run 2,000 MCMC iterations initialized at zero, and 500 as burn-in. We scale the regressors beforehand, as this is generally recommended to improve numerical stability and convergence. The following code simulates the model, and sets the hyperparameters of the algorithms.

### R code. Simulation: Logit model

```

1 rm(list = ls()); set.seed(10101)
2 library(mvtnorm); library(MCMCpack)
3 library(ggplot2); library(dplyr)
4 library(parallel); library(GGally)
5 ##### Generate correlated covariates
6 genCovMat <- function(K, rho = 0.4) {
7   Sigma0 <- matrix(1, K, K)
8   for (i in 2:K) {
9     for (j in 1:(i - 1)) {
10       Sigma0[i, j] <- runif(1, -rho, rho)^(i - j)
11     }
12   }
13   Sigma0 <- Sigma0 * t(Sigma0)
14   diag(Sigma0) <- 1
15   return(Sigma0)
16 }
17 ##### Simulate logistic regression data
18 simulate_logit_data <- function(K, N, beta_true) {
19   Sigma0 <- genCovMat(K)
20   X <- rmvnorm(N, mean = rep(0, K), sigma = Sigma0)
21   linpred <- X %*% beta_true
22   p <- 1 / (1 + exp(-linpred))
23   y <- rbinom(N, 1, p)
24   list(y = y, X = X)
25 }
26 ##### Parameters
27 K <- 10; N <- 100000
28 beta_true <- rep(0.5, K); B <- 5
29 batch_prop <- 0.01; Prior_prec <- 0.1
30 n_iter <- 2000; burnin <- 500
31 stepsize <- 1e-4
32 k_target1 <- 4 # beta5
33 k_target2 <- 5 # beta5
34 ks <- k_target1:k_target2
35 ##### Simulate data
36 sim_data <- simulate_logit_data(K, N, beta_true)
37 y <- sim_data$y
38 X <- scale(sim_data$X)
39 ##### Run MCMCpack logit
40 df <- as.data.frame(X)
41 colnames(df) <- paste0("X", 1:K)
42 df$y <- y
43 formula <- as.formula(paste("y ~", paste(colnames(df)[1:K],
44 collapse = " + "), "-1")))
44 posterior_mh <- MCMClogit(formula, data = df, b0 = 0, B0 =
45 Prior_prec, burnin = burnin, mcmc = n_iter)
45 full_posterior <- as.matrix(posterior_mh)[, 1:K]

```

The following code implements Algorithm A33, running five parallel MCMC chains and combining the resulting subposterior using three different weighting schemes: equal weights, weights based on marginal variances, and weights based on the full covariance matrices.

### R code. Consensus Monte Carlo (CMC): Logit regression

```

1  ### Split data
2 batch_ids <- split(1:N, sort(rep(1:B, length.out = N)))
3 ### Function to run MCMC on a subset
4 mcmc_batch <- function(batch_index, X, y, n_iter, burnin) {
5   ids <- batch_ids[[batch_index]]
6   X_b <- X[ids, ]
7   y_b <- y[ids]
8   mcmc_out <- MCMClogit(y_b ~ X_b - 1, burnin = burnin, mcmc
9   = n_iter, verbose = 0, b0 = 0, B0 = Prior_prec * (1/B))
9   return(mcmc_out)
10 }
11 ### Run in parallel
12 cl <- makeCluster(B)
13 clusterExport(cl, c("X", "y", "batch_ids", "n_iter", "burnin",
14   "mcmc_batch", "Prior_prec", "B"))
14 clusterEvalQ(cl, library(MCMCpack))
15 chains <- parLapply(cl, 1:B, function(b) mcmc_batch(b, X, y,
16   n_iter, burnin))
16 stopCluster(cl)
17 ### Stack MCMC results
18 posteriors <- lapply(chains, function(x) x[, 1:K]) # remove
18   intercept if added
19 ### CMC posteriors
20 equal_cmc <- Reduce("+", posteriors) / B
21 invvar_cmc <- {
22   vars <- lapply(posteriors, function(x) apply(x, 2, var))
23   weights <- lapply(vars, function(v) 1 / v)
24   weights_sum <- Reduce("+", weights)
25   weighted_post <- Reduce("+", Map(function(x, w) sweep(x,
25     2, w, "*"), posteriors, weights))
26   sweep(weighted_post, 2, weights_sum, "/")
27 }
```

### R code. Consensus Monte Carlo (CMC): Logit regression

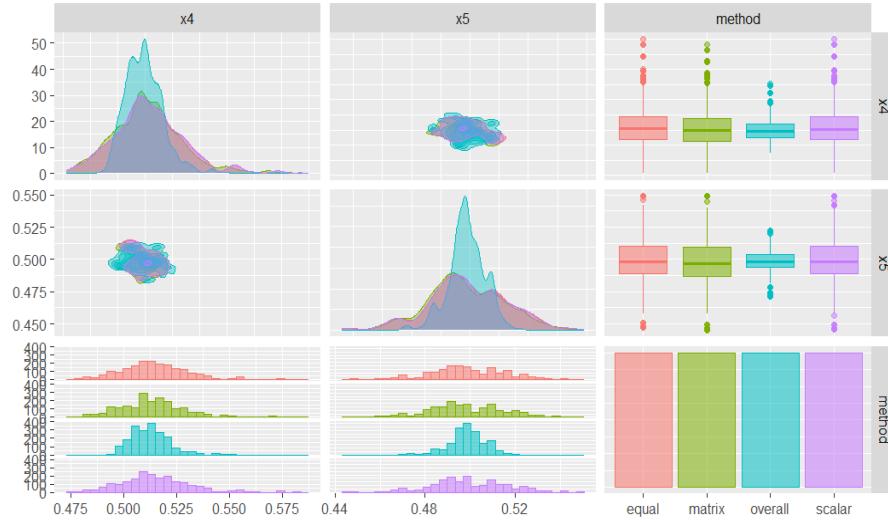
```

1 invmat_cmc <- {
2   covs <- lapply(posteriors, cov)
3   invs <- lapply(covs, solve)
4   weight_sum <- Reduce("+", invs)
5   consensus <- matrix(NA, nrow = n_iter, ncol = K)
6   for (i in 1:n_iter) {
7     draws <- lapply(posteriors, function(p) matrix(p[i, ],
8       ncol = 1))
8     weighted_sum <- Reduce("+", Map(function(w, d) w %*% d,
9       invs, draws))
9     consensus[i, ] <- as.vector(solve(weight_sum, weighted_
10      sum)))
10   }
11   consensus
12 }
13 # Combine all for plotting
14 build_df <- function(mat, method) {
15   df <- as.data.frame(mat)
16   colnames(df) <- paste0("x", ks)
17   df$method <- method
18   return(df)
19 }
20 df_full <- build_df(full_posterior[,ks], "overall")
21 df_equal <- build_df(equal_cmc[,ks], "equal")
22 df_scalar <- build_df(invvar_cmc[,ks], "scalar")
23 df_matrix <- build_df(invmat_cmc[,ks], "matrix")
24 df_plot <- rbind(df_full, df_matrix, df_scalar, df_equal)
25 # Plot
26 ggpairs(df_plot, aes(color = method, fill = method, alpha =
27   0.4), upper = list(continuous = GGally::wrap("density",
28   alpha = 0.4)), lower = list(continuous = GGally::wrap("density",
29   alpha = 0.4)), diag = list(continuous = GGally
30   ::wrap("densityDiag", alpha = 0.4)))

```

By running the code, you can verify that the computational time of the CMC algorithm is lower than that of the Metropolis–Hastings algorithm. Figure 12.10 shows the posterior distributions of  $\beta_4$  and  $\beta_5$ . We observe that all three weighting schemes perform reasonably well, yielding posterior modes similar to those obtained from the full-data MCMC algorithm. However, the consensus Monte Carlo (CMC) methods produce more dispersed draws, particularly when using equal weights. In contrast, the weighting schemes based on marginal variances and the full covariance matrices yield comparable and more concentrated posterior distributions.

To implement the SGLD algorithm (Algorithm A34), we set  $n = 0.01 \cdot N$ ,

**FIGURE 12.10**

Simulation: Posterior distributions from consensus Monte Carlo in logit simulation exercise

and the step size to  $1 \times 10^{-4}$ . The following code illustrates how to implement the SGLD algorithm.<sup>4</sup>

In Exercise 8, you are asked to implement the control variate version of SGLD. Begin by running 1,500 SGD iterations to locate the posterior mode. This mode should then be used as the initial value for a subsequent run of 1,000 SGLD iterations.

By running the code, you can verify that the computational time of the SGLD algorithm is lower than that of the Metropolis–Hastings algorithm. Figure 12.11 shows the posterior distributions of the fifth location parameter obtained from SGLD and Metropolis–Hastings. We observe that both modes are centered around the true population value; however, the SGLD distribution exhibits greater dispersion compared to the Metropolis–Hastings distribution.

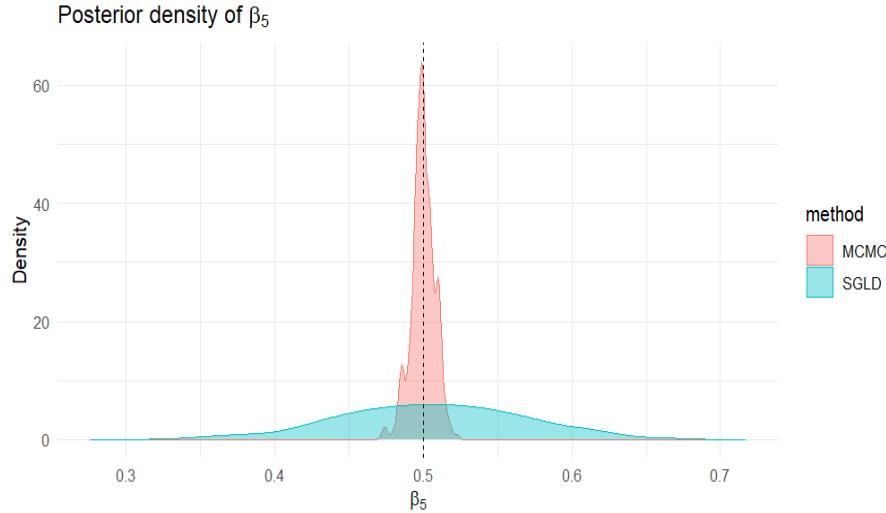
<sup>4</sup>There is an **R** package called *sgmcmc*, developed by [24], which provides implementations of various stochastic gradient MCMC methods, including SGLD and SGHMC. However, this package depends on version 1 of the *tensorflow* package, while the current version is 2, and *sgmcmc* has not been updated on CRAN. We attempted to install the package from its GitHub repository using the `devtools::install_github("STOR-i/sgmcmc")` command, but encountered compatibility issues due to conflicting TensorFlow versions.

### R code. Stochastic Gradient Langevin Dynamic (SGLD): Logit regression

```

1 SGLD_step <- function(beta, y, X, stepsize, batch_size,
2   prior_var = 10) {
3   N <- nrow(X); K <- length(beta)
4   ids <- sample(1:N, size = batch_size, replace = FALSE)
5   grad <- rep(0, K)
6   for (i in ids) {
7     xi <- X[i, ]
8     eta <- sum(xi * beta)
9     pi <- 1 / (1 + exp(-eta))
10    grad_i <- -(y[i] - pi) * xi
11    grad <- grad + grad_i
12  }
13  grad <- grad / batch_size * N
14  grad <- grad + beta / prior_var # gradient of log-prior
15  noise <- rnorm(K, 0, sqrt(stepsize))
16  beta_new <- beta - 0.5 * stepsize * grad + noise
17  return(beta_new)
18 }
19 #### SGLD algorithm
20 run_SGLD <- function(y, X, stepsize, batch_prop, n_iter,
21   burnin, beta_init = NULL) {
22   N <- nrow(X)
23   K <- ncol(X)
24   batch_size <- round(batch_prop * N)
25   beta_mat <- matrix(0, n_iter + burnin, K)
26   beta_mat[1, ] <- if (is.null(beta_init)) rep(0, K) else
27     beta_init
28   for (s in 2:(n_iter + burnin)) {
29     beta_mat[s, ] <- SGLD_step(beta_mat[s - 1, ], y, X,
30       stepsize, batch_size)
31   }
32   beta_mat[(burnin + 1):(n_iter + burnin), ]
33 }
34 #### Run SGLD
35 posterior_sgld <- run_SGLD(y = y, X = X, stepsize, batch_
36   prop, n_iter, burnin)
37 #### Compare densities for beta5
38 df_plot <- data.frame(value = c(posterior_sgld[, k_target2],
39   posterior_mh[, k_target2]), method = rep(c("SGLD", "MCMC"),
40   each = n_iter))
41 ggplot(df_plot, aes(x = value, fill = method, color = method
42 )) + geom_density(alpha = 0.4) + geom_vline(xintercept =
43   beta_true[k_target2], linetype = "dashed", color =
44   "black") + labs(title = expression(paste("Posterior
45   density of ", beta[5])), x = expression(beta[5]), y =
46   "Density") + theme_minimal()

```

**FIGURE 12.11**

Simulation: Posterior distributions from stochastic gradient Langevin Dynamic logit simulation exercise

## 12.6 Summary

In this chapter, we introduced several Bayesian machine learning methods designed to address the challenges posed by *wide* and *tall* data. However, the field of Bayesian machine learning is rapidly evolving, and the material presented here should be viewed as an introductory overview. Many important topics were not covered but are highly relevant, such as Bayesian neural networks [318] and neural posterior estimation [329, 280, 183]. Other key approaches, such as Variational Bayes, particularly in its stochastic implementations, which are rooted in machine learning and offer scalable solutions for tall data, are introduced in Chapter 14 [443].

## 12.7 Exercises

### 1. Simulation exercise: the Bayesian LASSO continues

Program the Gibbs sampler for the Bayesian LASSO from scratch, assuming a hierarchical structure for the global shrinkage parameter, where both the shape and rate parameters are set to 1. Per-

form inference using this sampler in the Bayesian LASSO simulation exercise and compare the results with those obtained using the *monomvn* package.

2. [226] employ SSVS to identify the main drivers of civil conflict in the post-Cold War era, considering a set of 35 potential determinants across 175 countries worldwide. We use a subset of their dataset provided in *Conflict.csv*, where the dependent variable is *conflictcw*, a binary indicator of civil conflict. Perform SSVS using the *BoomSpikeSlab* package, specifically the *lm.spike* function, to identify the best subset of models.
3. [437] proposes an SSVS approach for binary response models. Use the dataset *Conflict.csv*, where the dependent variable is *conflictcw*, to perform SSVS using the *BoomSpikeSlab* package, specifically the *logit.spike* function, in order to identify the best subset of models. Compare the results with those obtained in Exercise 2.

#### 4. Example: Simulation exercise $K > N$

Use the simulation setting from the Bayesian LASSO and SSVS examples, but now assume there are 600 inputs. This setup implies that the number of inputs exceeds the sample size. In such a scenario, there is no unique solution to the least squares estimator because the determinant of  $\mathbf{W}^\top \mathbf{W}$  is zero. This means the matrix is not invertible, and consequently, standard inference procedures based on the least squares estimator cannot be applied. On the other hand, Bayesian inference in this setup is well-defined because the prior helps regularize the problem, which is a key motivation for these methods.

#### 5. Simulation exercise: the BART model continues

Compute Friedman's partial dependence functions [148] for all variables in the BART model simulation example, and plot the posterior mean along with the 95% credible intervals.

6. [100] presents BART probit for classification. This method can be implemented using the *BART* package through the function *pbart*. Use the file *Conflict.csv*, where the dependent variable is *conflictcw*, to perform BART probit, implementing *k-fold* cross-validation to select the threshold that maximizes the sum of the true positive and true negative rates. Additionally, identify the most important predictors by evaluating different numbers of trees.

#### 7. Simulation exercise: The Gaussian Process simulation continues

Simulate the process

$$f_i = \sin(2\pi x_{i1}) + \cos(2\pi x_{i2}) + \sin(x_{i1}x_{i2}) + \mu_i,$$

where  $\mu_i \stackrel{\text{i.i.d.}}{\sim} N(0, 0.1^2)$ ,  $x_{ik} \sim U(0, 1)$  for  $k = 1, 2$ , and the sample size is 500.

Define a grid of 20 evenly spaced values between 0 and 1 for each covariate  $x_{ik}$ , and use this grid to perform prediction.

Estimate the hyperparameters of the Gaussian Process by maximizing the log marginal likelihood. Then, use the *km* function from the *DiceKriging* package to fit the Gaussian Process, fixing the noise variance at the value that maximizes the log marginal likelihood.

Finally, use the fitted model to predict the outputs on the grid points, and produce a 3D plot showing the predicted surface along with the training data points.

#### 8. Simulation exercise: Stochastic gradient MCMC continues

Program from scratch the stochastic gradient Langevin dynamic algorithm for the logit simulation exercise implementing the control variate version performing 1,500 stochastic gradient descent iterations to locate the posterior mode, which should be then used as the initial value for 1,000 subsequent MCMC iterations using a step size set to  $1 \times 10^{-4}$ .

9. Perform the simulation according to the model  $y_i = 1 - 2x_{i1} + 0.5x_{i2} + \mu_i$ , where  $\mu_i \sim N(0, 1)$ , the sample size is 100,000, and the covariates  $\mathbf{x}_i \sim N(\mathbf{0}, \mathbf{I}_2)$ . Use 5,000 MCMC iterations and a batch size of 1,000 to implement the SGLD algorithm. Set a learning rate schedule that yields sensible results.

Assume independent priors  $\pi(\boldsymbol{\beta}, \sigma^2) = \pi(\boldsymbol{\beta}) \times \pi(\sigma^2)$ , with  $\boldsymbol{\beta} \sim N(\mathbf{0}, \mathbf{I}_3)$  and  $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ , where  $\alpha_0 = \delta_0 = 0.01$ .

# 13

---

## Causal inference

---

In this chapter, we present some Bayesian methods to perform inference on *causal* effects. The point of departure is *identification conditions* which are the set of assumptions that allow us to express the *estimand* —the causal or statistical quantity of interest— as a unique function of the observable data distribution. These identification conditions are conceptually distinct from the econometric or statistical framework used to perform inference once those conditions are satisfied. In other words, the assumptions necessary for identification do not constrain whether we apply a Frequentist or a Bayesian approach for estimation and inference.

Identification addresses the question: *Can the causal effect be expressed as a function of the observable data distribution under certain assumptions?* Once the causal effect is identified in terms of the observed data distribution, we can proceed with statistical inference using either Frequentist methods or Bayesian methods. These inferential frameworks differ in how they estimate and quantify uncertainty but operate on the same identified causal effect. Thus, the identification assumptions are logically prior and independent of whether we adopt a Frequentist or a Bayesian inferential paradigm.

Thus, we present the underlying identification assumptions employed in popular strategies such as randomized controlled trials, conditional independence, and instrumental variables, among others, as well as the Bayesian inferential framework used to estimate causal effects in these settings.

We should be clear that this is only an introduction to causal inference, and we do not aim to provide an in-depth discussion. There are excellent texts on causal inference, such as [19, 18, 218, 204, 109, 84]. We recommend that readers consult these references for a deeper understanding of some of the concepts and tools introduced in this chapter.

---

### 13.1 Identification setting

The aim is to identify the *causal* or *treatment effect* of a particular regressor or treatment on an outcome. For example, this could involve estimating the causal effect of a labor training program on individuals' earnings or unemployment, or the price elasticity of demand. The starting point is the *potential*

*outcomes* framework [385], which defines a *counterfactual scenario* describing what would happen to the outcome variable if the treatment status or level were different (commonly referred to as the *Rubin causal model* [210]).

In this context, *treatment status* refers to a binary treatment case with two potential outcomes; for instance, what would my earnings be if I had participated in the training program? Conversely, *treatment level* applies to cases where there is a potential outcome for each level of the treatment; for example, what would demand be if the price had increased by 10%?

Following the potential outcomes notation in the binary treatment case, let  $D_i = 1$  and  $D_i = 0$  indicate the treatment status, corresponding to *treated* and *control* units, respectively. The potential outcomes  $Y_i(1)$  and  $Y_i(0)$  represent the outcome for unit  $i = 1, 2, \dots, N$  under treatment and control, respectively. For instance,  $Y_i(1)$  denotes the employment status an individual would have if she/he would have participated in the training program, regardless of their actual participation, whereas  $Y_i(0)$  denotes the employment status if the individual would not have attended the program. The individual-level treatment effect is then defined as

$$\tau_i = Y_i(1) - Y_i(0).$$

The potential outcomes framework can also be extended to situations where the treatment is continuous [214]. In this case, let  $Y_i(x_s)$  denote the outcome for unit  $i$  under the counterfactual scenario where the treatment variable  $X_s$  takes the value  $x_s$ . The “treatment effect” of changing  $X_s$  from  $x_s$  to  $x'_s$ , for example, the effect of increasing the price by 10% on demand, is given by

$$\beta_{si} = Y_i(x'_s) - Y_i(x_s).$$

When the change is infinitesimal, the causal effect can be expressed as a marginal effect:

$$\beta_{si} = \frac{\partial Y_i(x_s)}{\partial x_s}.$$

This setting is more complex than the binary treatment case because there is a potential outcome for each possible value of  $x_s$  [177]. Therefore, we begin by considering the binary treatment case. However, the *fundamental problem of causal inference* [210] remains: it is impossible to observe the same unit under different *treatment statuses* simultaneously. Consequently, we must learn about causal effects by comparing the outcomes of treated and untreated units.

Bayesian inference for causal effects [386] is more direct than procedures based on Fisher’s *p*-value approach, which relies on the logic of stochastic proof by contradiction, or Neyman’s randomization-based inference, which is grounded in the idea of repeated sampling and the construction of confidence intervals for treatment effects [390].

Bayesian inference for causal effects treats the potential outcomes as random variables and involves computing the posterior predictive distribution to evaluate treatments not received, conditional on the observed responses to treatments actually received. This approach yields the posterior distribution of the causal estimands as a function of both the observed outcomes and the unobserved potential outcomes, which are treated as missing data and handled through data augmentation methods [428].

In addition, it is important to note that the likelihood function does not contain information about the correlation between potential outcomes due to *the fundamental problem of causal inference*. Consequently, most of the classical literature on treatment effects has focused on average treatment effects, which require only the marginal distributions of the potential outcomes  $Y(1)$  and  $Y(0)$  [200]. Nevertheless, it is also possible to estimate *quantile treatment effects* [3, 83].

There are also Bayesian proposals for performing inference on the correlation between potential outcomes [249, 200]. These approaches allow recovery of the joint distribution of the potential outcomes, thereby enabling inference beyond average treatment effects. However, this requires learning from the prior rather than the data, or imposing additional structure on the causal model.

A Bayesian framework further makes it possible to estimate *distributional treatment effects* with uncertainty quantification, as formally defined by [1]. See also [200] for related proposals. These effects also characterize the entire distribution of potential outcomes under treatment and control, rather than focusing solely on averages, thereby capturing heterogeneity. This is particularly relevant for policy, as it enables estimation of the proportion of the population that benefits from a social program or the probability that a treated individual experiences a positive effect [200, 355].

Furthermore, under a Bayesian framework, the impact of adding or relaxing identification assumptions can be assessed by examining how the distributional treatment effects change [217].

---

## 13.2 Randomized controlled trial (RCT)

Two of the most relevant estimands in the causal inference literature are the average treatment effect (ATE),

$$\text{ATE} = \mathbb{E}[Y_i(1) - Y_i(0)],$$

and the average treatment effect on the treated (ATT),

$$\text{ATT} = \mathbb{E}[Y_i(1) - Y_i(0) \mid D_i = 1].$$

However, these estimands are not directly identifiable without additional

assumptions because we never observe the same unit under both treatment states simultaneously. Therefore, identification restrictions are required to express these causal quantities in terms of the observed data.

Note that the observed mean difference between treated and untreated units can be decomposed as:

$$\begin{aligned} \mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0] &= \underbrace{\mathbb{E}[Y_i(1) - Y_i(0) | D_i = 1]}_{\text{ATT}} \\ &\quad + \underbrace{\left( \mathbb{E}[Y_i(0) | D_i = 1] - \mathbb{E}[Y_i(0) | D_i = 0] \right)}_{\text{Selection Bias}}, \end{aligned}$$

that is, the observed mean difference equals the ATT plus the *selection bias*, which captures the difference between the expected potential outcome of treated individuals had they not been treated (unobserved) and the expected outcome of untreated individuals (observed).

The sign of the selection bias depends on the nature of self-selection into treatment. Two illustrative examples are:

**TABLE 13.1**

Illustrative examples of selection bias

Scenario	$\mathbb{E}[Y(0)   D = 1]$	$\mathbb{E}[Y(0)   D = 0]$	Selection Bias
<i>Training Program</i>	80	70	+10
<i>Remedial Education</i>	60	70	-10

In the first scenario, motivated individuals self-select into a training program and would have higher earnings even without participation, producing a positive selection bias and causing the observed mean difference to overestimate the ATT. In the second scenario, disadvantaged individuals are more likely to enroll in a remedial education program and would have lower outcomes without the program, resulting in a negative selection bias and causing the observed mean difference to underestimate the ATT.

Note that if we assume that *assignment to treatment is random*,

$$\{Y_i(1), Y_i(0)\} \perp D_i,$$

then the selection bias becomes zero,

$$\mathbb{E}[Y_i(0) | D_i = 1] - \mathbb{E}[Y_i(0) | D_i = 0] = 0,$$

which means there is no selection bias under random assignment.

Moreover, under random assignment, ATT equals ATE:

$$\tau = \underbrace{\mathbb{E}[Y_i(1) | D_i = 1] - \mathbb{E}[Y_i(0) | D_i = 1]}_{\text{ATT}} = \underbrace{\mathbb{E}[Y_i(1) - Y_i(0)]}_{\text{ATE}}.$$

Consequently,

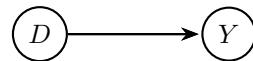
$$\tau = \mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0],$$

that is, the average causal effect is function of the data under random assignment to treatment.

*Randomized controlled trials* (RCTs) are characterized by *random assignment to treatment*. The identification of causal effects in RCTs additionally relies on two key assumptions: *overlap*, which requires that every unit has a positive probability of being assigned to both treatment and control groups ( $0 < P(D_i = 1) < 1$ ), and the *Stable Unit Treatment Value Assumption* (SUTVA). The latter consists of two components: (i) *no interference*, meaning one unit's potential outcome is unaffected by another unit's treatment, and (ii) *consistency*, meaning the observed outcome equals the potential outcome corresponding to the received treatment.

We can represent the causal mechanism in an RCT using *causal diagrams* [332, 333]. These diagrams provide a powerful framework for visualizing and analyzing the underlying structures that enable causal identification. It is important to note that causal diagrams do not impose specific functional forms, such as those assumed in linear regression models, and are closely related to structural equation models (SEMs). SEMs offer an alternative perspective on causal inference in econometrics, which is closely connected to the potential outcomes approach [187, 293, 214].

In particular, we can depict the causal mechanism using a *Directed Acyclic Graph (DAG)*, which is a graphical representation of a set of variables and their assumed causal relationships. A DAG consists of *nodes*, representing variables, and *directed edges (arrows)*, representing direct causal effects from one variable to another. The graph is *acyclic*, meaning it contains no directed cycles: starting from any node and following the arrows, it is impossible to return to the same node. A defining property of causal DAGs is that, conditional on its direct causes, any variable on the DAG is independent of any other variable for which it is not a cause (*causal Markov assumption*). See [204] for a nice introduction, and the package *dagitty* in **R** software to perform visualization and analysis of DAGs. Figure 13.1 illustrates the causal structure underlying RCTs.

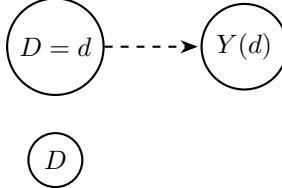


**FIGURE 13.1**

Directed Acyclic Graph (DAG) implied by a Randomized Controlled Trial (RCT).

The counterfactual representation of the DAG, known as the *Single-World Intervention Graph* (SWIG), is shown in Figure 13.2. In this figure, the treatment variable  $D$  is split to reflect the intervention: we fix  $D$  at a specific value  $d$  (the intervention), and replace the outcome  $Y$  with its counterfactual

representation  $Y(d)$ , which denotes the value of  $Y$  if  $D$  were set to  $d$ . The natural value of  $D$  is also shown but becomes irrelevant for the outcome once the intervention is imposed.



**FIGURE 13.2**

Single-World Intervention Graph (SWIG) for intervention  $do(D = d)$ :  $Y$  is replaced by the counterfactual  $Y(d)$ , and  $D$  is split into the fixed value  $D = d$  and its natural value  $D$ . The dashed arrow indicates the fixed causal assignment.

Note that we can express the observed outcome in terms of the potential outcomes as:

$$Y_i = [Y_i(1) - Y_i(0)]D_i + Y_i(0),$$

which shows that the observed outcome equals the control potential outcome plus the treatment effect if treated.

Under the assumption of a constant treatment effect and a linear *Conditional Expectation Function* (CEF)  $\mathbb{E}[Y_i | D_i]$ , this relationship can be represented as a linear regression model [19]:

$$Y_i = \underbrace{\beta_0}_{\mathbb{E}[Y_i(0)]} + \underbrace{\tau}_{\text{constant treatment effect}} D_i + \underbrace{\mu_i}_{Y_i(0) - \mathbb{E}[Y_i(0)]},$$

where  $\tau$  represents the common treatment effect across all units.

Under random assignment, we have

$$\mathbb{E}[\mu_i | D_i] = \mathbb{E}[Y_i(0) | D_i] - \mathbb{E}[Y_i(0)] = \mathbb{E}[Y_i(0)] - \mathbb{E}[Y_i(0)] = 0,$$

that is, the error term is *mean-independent of treatment status*. This implies that a regression of  $Y_i$  on  $D_i$  consistently estimates the causal effect under random assignment.

### Example: Simple example

Bayesian inference for causal effects in a completely randomized experiment without covariates can be illustrated using the normal model in [389]. Assume that  $Y_i(d) \sim N(\mu_d, \sigma_d^2)$  for  $d \in \{1, 0\}$ . Then, the likelihood function given a random sample and independent treatment assignment is:

$$p(\mathbf{y} | \mathbf{d}; \mu_1, \mu_0, \sigma_1^2, \sigma_0^2) = \prod_{i:d_i=1} \phi(y_i | \mu_1, \sigma_1^2) \prod_{i:d_i=0} \phi(y_i | \mu_0, \sigma_0^2)$$

$$= \prod_i \phi(y_i | \mu_1, \sigma_1^2)^{d_i} \phi(y_i | \mu_0, \sigma_0^2)^{1-d_i},$$

where  $\mathbf{y}$  and  $\mathbf{d}$  collect the observations of the outcomes and treatments, and  $\phi(\cdot)$  denotes the normal density function.

Assume conjugate priors:

$$\mu_d | \sigma_d^2 \sim N\left(\mu_{0d}, \frac{\sigma_d^2}{\beta_{0d}}\right), \quad \sigma_d^2 \sim IG\left(\frac{\alpha_{0d}}{2}, \frac{\delta_{0d}}{2}\right).$$

From the normal/inverse-gamma model in Chapter 3, the posterior conditional distributions are:

$$\mu_d | \sigma_d^2, \{y_i, d_i\}_{i:d_i=d} \sim N\left(\mu_{nd}, \frac{\sigma_d^2}{\beta_{nd}}\right),$$

where

$$\mu_{nd} = \frac{\beta_{0d}\mu_{0d} + N_d\bar{y}_d}{\beta_{0d} + N_d}, \quad \beta_{nd} = \beta_{0d} + N_d,$$

$\bar{y}_d$  and  $N_d$  denote the sample mean and sample size of group  $d$ . For the variance:

$$\sigma_d^2 | \{y_i, d_i\}_{i:d_i=d} \sim IG\left(\frac{\alpha_{nd}}{2}, \frac{\delta_{nd}}{2}\right),$$

where

$$\alpha_{nd} = \alpha_{0d} + N_d, \quad \delta_{nd} = \sum_{i:d_i=d} (y_i - \bar{y}_d)^2 + \delta_{0d} + \frac{\beta_{0d}N_d}{\beta_{0d} + N_d}(\bar{y}_d - \mu_{0d})^2.$$

In addition, the predictive distribution is given by

$$Y(d) | \mathbf{y} \sim t\left(\mu_{nd}, \frac{(\beta_{nd}+1)\delta_{nd}}{\beta_{nd}\alpha_{nd}}, \alpha_{nd}\right).$$

Therefore, given the definition of the average treatment effect,

$$\begin{aligned} \text{ATE} &= \mathbb{E}[Y(1) - Y(0)] \\ &= \mathbb{E}[Y(1)] - \mathbb{E}[Y(0)] \\ &= \int_{\mathcal{R}} y(1)f_{Y(1)}(y(1))dy(1) - \int_{\mathcal{R}} y(0)f_{Y(0)}(y(0))dy(0) \\ &= \mu_1 - \mu_0. \end{aligned}$$

Note that this expectation is taken with respect to the population distribution of the potential outcomes, not with respect to the posterior distribution.<sup>1</sup> Thus, the posterior distribution of the ATE is given by

$$\pi(\text{ATE} | \mathbf{y}) = \pi(\mu_1 - \mu_0 | \mathbf{y}).$$

---

<sup>1</sup>Other functions of the potential outcomes, such as  $P(Y(1) \geq Y(0))$ , require the joint distribution of the potential outcomes rather than just the marginal distributions, as is the case for the ATE.

We can obtain this posterior distribution by simulation, because the difference of two independent Student's  $t$ -distributed random variables does not follow a Student's  $t$  distribution. However, as the degrees of freedom increase, each posterior distribution of  $\mu_d$  converges to a normal distribution, and the difference of two normals is also normal. Consequently, an approximate point estimate of the ATE is  $\mu_{n1} - \mu_{n0}$ , which asymptotically equals  $\bar{y}_1 - \bar{y}_0$  under non-informative priors. This is the classical *estimator* of the ATE.

Finally, note that we can also obtain the posterior distribution of the ATE using the simple linear regression framework by assuming non-informative prior distributions, in which case the posterior mean of the slope parameter coincides with the maximum likelihood estimator (see Exercise 1).

Let's assume that  $\mu_1 = 15$ ,  $\mu_0 = 10$ ,  $\sigma_1^2 = 4$ ,  $\sigma_0^2 = 2$ , and  $N_1 = N_0 = 100$ , and compute the posterior distribution of the ATE. The following code illustrates this procedure, while Figures 13.3 and 13.4 display the predictive densities of the potential outcomes and the posterior distribution of the ATE. The posterior mean is 4.60, and the 95% credible interval is (4.10, 5.12). Note that the population value is  $\mu_1 - \mu_0 = 5$ .

### R code. Simple example: Average treatment effects

```

1 set.seed(10101)
2 # Parameters
3 mu1 <- 15; mu0 <- 10
4 sigma1 <- sqrt(4); sigma0 <- sqrt(2)
5 N1 <- 100; N0 <- 100
6 # Simulate data
7 y1 <- rnorm(N1, mu1, sigma1); y0 <- rnorm(N0, mu0, sigma0)
8 # Prior hyperparameters
9 alpha0 <- 0.01; delta0 <- 0.01
10 simulate_t <- function(y) {
11   N <- length(y); ybar <- mean(y)
12   sse <- sum((y - ybar)^2)
13   alpha_n <- alpha0 + N; delta_n <- sse + delta0
14   scale2 <- ((N + 1) * delta_n) / (N * alpha_n)
15   df <- alpha_n
16   loc <- ybar; scale <- sqrt(scale2)
17   rt(1, df = df) * scale + loc
18 }
19 # Posterior predictive draws
20 ppd1 <- replicate(1000, simulate_t(y1))
21 ppd0 <- replicate(1000, simulate_t(y0))
22 # Plot
23 hist(ppd1, col = rgb(1, 0, 0, 0.4), freq = FALSE, main = "
  Posterior Predictive",
24 xlab = "Y", xlim = c(5, 25))
```

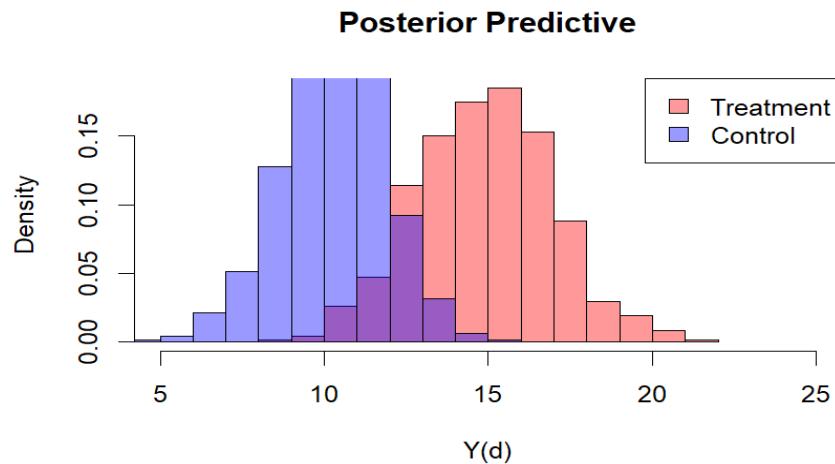
### R code. Simple example: Average treatment effects

```

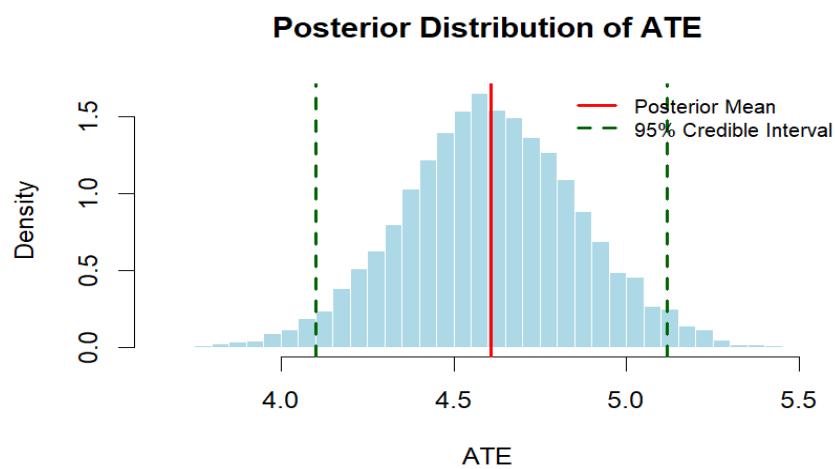
1 hist(ppd0, col = rgb(0, 0, 1, 0.4), freq = FALSE, add = TRUE
      )
2 legend("topright", legend = c("Treatment", "Control"), fill
      = c(rgb(1, 0, 0, 0.4), rgb(0, 0, 1, 0.4)))
3 # Posterior distribution of ATE
4 posterior_mu <- function(y) {
5   N <- length(y); ybar <- mean(y)
6   sse <- sum((y - ybar)^2)
7   alpha_n <- alpha0 + N
8   delta_n <- sse + delta0
9   # Posterior variance for mean
10  var_mu <- delta_n / (alpha_n * N)
11  df <- alpha_n
12  loc <- ybar; scale <- sqrt(var_mu)
13  rt(1, df = df) * scale + loc
14 }
15 # Posterior draws for parameters
16 n_draws <- 10000
17 mu1_draws <- replicate(n_draws, posterior_mu(y1))
18 mu0_draws <- replicate(n_draws, posterior_mu(y0))
19 # Parameter uncertainty: difference of means
20 ate_draws <- mu1_draws - mu0_draws
21 summary(coda::mcmc(ate_draws))
22 # Summaries
23 ate_mean <- mean(ate_draws)
24 ci <- quantile(ate_draws, c(0.025, 0.975))
25 cat("Posterior mean of ATE:", ate_mean, "\n")
26 cat("95% Credible Interval:", ci, "\n")
27 # Plot posterior distribution of ATE
28 hist(ate_draws, breaks = 50, freq = FALSE,
29 main = "Posterior Distribution of ATE",
30 xlab = "ATE (Y(1) - Y(0))", col = "lightblue", border =
      "white")
31 abline(v = ate_mean, col = "red", lwd = 2)
32 abline(v = ci, col = "darkgreen", lty = 2, lwd = 2)
33 legend("topright", legend = c("Posterior Mean", "95%
      Credible Interval"), col = c("red", "darkgreen"), lwd =
      2, lty = c(1, 2), bty = "n", cex = 0.8) # Smaller
      legend using cex

```

An RCT is considered the gold standard for identifying causal effects because it provides the strongest basis for satisfying the key identification assumption in causal inference: *independence between treatment assignment and potential outcomes*. However, RCTs may face challenges such as *non-compliance*, which occurs when individuals do not adhere to their assigned treatment in an experimental study. This is important because, in many real-



**FIGURE 13.3**  
Predictive distributions: Potential outcomes.



**FIGURE 13.4**  
Posterior distribution: Average treatment effects.

world settings, some individuals do not take the treatment they were assigned, leading to deviations from the ideal randomized controlled trial design. In addition, RCTs are sometimes infeasible due to ethical, logistical, or financial constraints. Moreover, they can be too narrowly focused or too localized to provide general conclusions about what works [114]. Thus, the *external validity* of causal effects from RCTs may be questionable.

It is also important to note that RCTs primarily identify the mean of the treatment effect distribution but do not capture other features, such as the median or higher-order moments. These additional aspects of the distribution of treatment effects can be highly relevant for policymakers and stakeholders. See [114] for a detailed discussion of other potential shortcomings of RCTs.

### 13.3 Conditional independence assumption (CIA)

In practice, researchers often work with *observational data*, where the treatment status is not randomly assigned because units actively choose the treatment they receive. In an RCT, the assignment mechanism is determined by chance, whereas in observational studies it is driven by choice. In these situations, we can identify the causal effect if the *conditional independence assumption* (CIA) holds. This assumption states that the potential outcomes are independent of the treatment status, conditional on a set of observed *pre-treatment variables*:

$$\{Y_i(1), Y_i(0)\} \perp D_i \mid \mathbf{X}_i.$$

This means that, conditional on the *pre-treatment variables*  $\mathbf{X}_i$ , treatment assignment is as good as random. This property is known as *unconfoundedness given  $\mathbf{X}_i$* , or equivalently, the *no unmeasured confounders* assumption. When this condition is combined with the requirement that, for all possible values of  $\mathbf{X}_i$ , there is a positive probability of receiving each treatment level ( $0 < P(D_i = 1 \mid \mathbf{X}_i) < 1$ ), the joint condition is referred to as *strong ignorability*. Identification of average causal effects in this setting also requires the SUTVA.

Note that under the CIA,

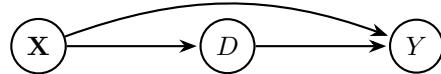
$$\begin{aligned} \text{ATE} &= \mathbb{E}[Y_i(1)] - \mathbb{E}[Y_i(0)] \\ &= \mathbb{E}_{\mathbf{X}} \{\mathbb{E}[Y_i(1) \mid \mathbf{X}_i] - \mathbb{E}[Y_i(0) \mid \mathbf{X}_i]\} \\ &= \mathbb{E}_{\mathbf{X}} \{\mathbb{E}[Y_i(1) \mid \mathbf{X}_i, D_i = 1] - \mathbb{E}[Y_i(0) \mid \mathbf{X}_i, D_i = 0]\} \\ &= \mathbb{E}_{\mathbf{X}} \{\mathbb{E}[Y_i \mid \mathbf{X}_i, D_i = 1] - \mathbb{E}[Y_i \mid \mathbf{X}_i, D_i = 0]\} \\ &= \mathbb{E}_{\mathbf{X}} \{\tau(\mathbf{X}_i)\}, \end{aligned}$$

where the second equality follows from the law of iterated expectations, the third uses CIA, and

$$\tau(\mathbf{X}_i) = \mathbb{E}[Y_i \mid \mathbf{X}_i, D_i = 1] - \mathbb{E}[Y_i \mid \mathbf{X}_i, D_i = 0]$$

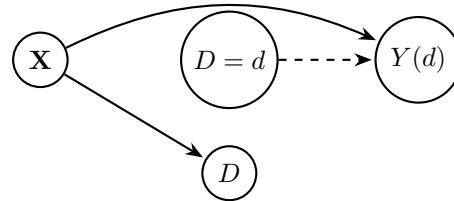
is the *conditional average treatment effect* (CATE) at covariate value  $\mathbf{X}_i$ , which captures treatment effect heterogeneity across different values of  $\mathbf{X}$ . Therefore, the ATE is the expectation of the CATE over the distribution of  $\mathbf{X}$ .

Figure 13.5 illustrates the causal structure underlying the CIA. We follow the convention that time flows from left to right. Here,  $\mathbf{X}$  - a vector of pre-treatment variables or *confounders*- influences both the treatment ( $D$ ) and the outcome ( $Y$ ). Under the CIA, we can identify the causal effect of  $D$  on  $Y$  by adjusting for  $\mathbf{X}$  because this causal structure satisfies the *back-door criterion*. This criterion states that a set of variables  $\mathbf{X}$  satisfies the condition for identifying the effect of  $D$  on  $Y$  if no variable in  $\mathbf{X}$  is a descendant of  $D$  and  $\mathbf{X}$  blocks every back-door path from  $D$  to  $Y$ . A *back-door path* is any path from  $D$  to  $Y$  that begins with an arrow pointing into  $D$ . In Exercise 3, we ask to construct the DAG in Figure 13.5, verify that it is acyclic, and check whether the causal effect of  $D$  on  $Y$  is identifiable by controlling for  $\mathbf{X}$  using the package *dagitty*. In addition, Figure 13.6 displays the SWIG.



**FIGURE 13.5**

Directed Acyclic Graph (DAG) implied by the Conditional Independence Assumption (CIA).



**FIGURE 13.6**

Single-World Intervention Graph (SWIG) for  $do(D = d)$ : The treatment  $D$  is split into the fixed intervention  $D = d$  and its natural value. The outcome is replaced by  $Y(d)$ .  $\mathbf{X}$  still influences  $Y(d)$ , so adjustment for  $\mathbf{X}$  is needed for identification.

The simple regression framework of the previous section can be extended to a multiple linear regression model by including the *pre-treatment variables*; therefore, if the CEF is linear in the treatment and pre-treatment variables:

$$Y_i = \beta_0 + \tau D_i + \mathbf{X}_i^\top \boldsymbol{\beta} + \mu_i.$$

Note that, by assumption in RCTs,  $D_i$  and  $\mathbf{X}_i$  are independent. Thus, the

identification of  $\tau$  is not affected under random assignment; however, including  $\mathbf{X}_i$  helps explain part of the variability in  $Y_i$ , thereby improving the precision of the estimates.

On the other hand, if the CIA is satisfied, the error term in the linear regression is defined as

$$\mu_i = Y_i(0) - \mathbb{E}[Y_i(0) | \mathbf{X}_i].$$

Thus,

$$\mathbb{E}[\mu_i | D_i, \mathbf{X}_i] = \mathbb{E}[Y_i(0) - \mathbb{E}[Y_i(0) | \mathbf{X}_i] | D_i, \mathbf{X}_i] = \mathbb{E}[Y_i(0) | D_i, \mathbf{X}_i] - \mathbb{E}[Y_i(0) | \mathbf{X}_i].$$

By the CIA,  $Y_i(0) \perp D_i | \mathbf{X}_i$ , so:

$$\mathbb{E}[Y_i(0) | D_i, \mathbf{X}_i] = \mathbb{E}[Y_i(0) | \mathbf{X}_i].$$

Therefore:

$$\mathbb{E}[\mu_i | D_i, \mathbf{X}_i] = 0.$$

This condition is known as *conditional mean independence*: after controlling for  $\mathbf{X}_i$ , treatment assignment is independent of unobserved determinants of the outcome. It justifies unbiased estimation of  $\tau$  by regression adjustment in observational studies.

#### **Example: Treatment effect of 401(k) eligibility on net financial assets**

We study the average treatment effect of eligibility (*e401*) for participation in the 401(k) retirement savings plan in the United States on net financial assets (*net\_tfa*) using the dataset *401k.csv*. The rationale for the exogeneity of 401(k) eligibility is that it becomes exogenous after conditioning on observable characteristics related to job choice, which may correlate with whether a firm offers this retirement plan [84].

Accordingly, we control for the following covariates: age (*age*), income (*inc*), family size (*fsize*), years of education (*educ*), a marital status indicator (*marr*), a two-earner status indicator (*twoearn*), a defined benefit pension status indicator (*db*), an IRA participation indicator (*pira*), and a home ownership indicator (*hown*). Under this specification, the key assumption is that eligibility is conditionally independent of net financial assets given these covariates [84].

We can use the framework in Section 6.1 to estimate the average treatment effect of 401k eligibility on net financial assets. The code 13.3 shows how to get the posterior distribution of this effect. Figure 13.7 displays the posterior distribution of the treatment effect, the 95% credible interval is (USD 3,473, USD 8,371), and the posterior mean is USD 5,903.

**R code. Treatment effect: 401(k) eligibility on net financial assets**

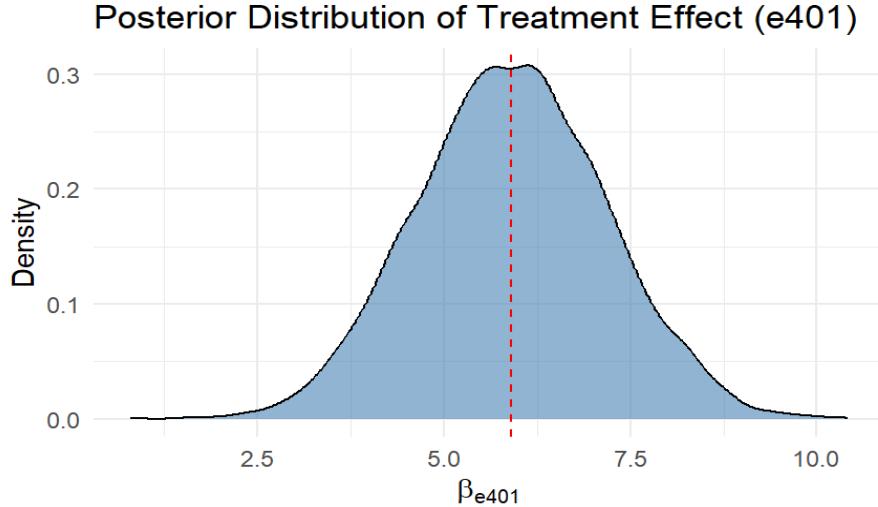
---

```

1 # 401k: Treatment effects
2 rm(list = ls())
3 set.seed(10101)
4
5 library(MCMCpack)
6 library(coda)
7 library(ggplot2)
8
9 mydata <- read.csv("https://raw.githubusercontent.com/
  BEsmarter-consultancy/BSTApp/refsrc/heads/master/DataApp/
  401k.csv", sep = ",", header = TRUE, quote = "")
10 attach(mydata )
11 y <- net_tfa
12 # net_tfa: net financial assets
13 # Regressors quantity including intercept
14 X <- cbind(e401, age, inc, fsize, educ, marr, twoearn, db,
  pira, hown, 1)
15 # e401: 401k eligibility
16 # age, income, family size, years of education, marital
  status indicator, two-earner status indicator, defined
  benefit pension status indicator, IRA participation
  indicator, and home ownership indicator.
17
18 # Posterior distributions using packages: MCMCpack sets the
  model in terms of the precision matrix. We use default
  parameters
19 posterior <- MCMCpack::MCMCregress(y~X-1)
20 summary(coda::mcmc(posterior))
21
22 # Extract posterior samples for e401 (first coefficient)
23 beta_e401 <- posterior[, 1]
24
25 # Convert to data frame for ggplot
26 df_posterior <- data.frame(beta_e401)
27
28 # Plot with ggplot
29 ggplot(df_posterior, aes(x = beta_e401)) +
30 geom_density(fill = "steelblue", alpha = 0.6) +
31 geom_vline(xintercept = mean(beta_e401), color = "red",
  linetype = "dashed") +
32 labs(
33 title = "Posterior Distribution of Treatment Effect (e401)",
34 x = expression(beta[e401]),
35 y = "Density"
36 ) + theme_minimal(base_size = 14)

```

---

**FIGURE 13.7**

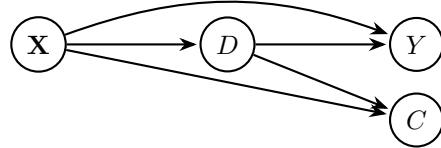
Posterior distribution: Treatment effect 401k eligibility on net financial assets.

Observe that the identification strategy relies on conditioning on  $\mathbf{X}$ . However, adding more controls is not always safe because some variables, known as *bad controls*, can introduce bias if included in the adjustment set [19]. One important type of bad control is a *collider*, a variable that is caused by (or is a common effect of) two or more other variables in a DAG. Colliders play a critical role because conditioning on them can create spurious associations between their causes, leading to what is known as *collider bias* (or selection bias).

Figure 13.8 illustrates this situation. Here,  $C$  is a common effect of  $D$  and  $\mathbf{X}$ . Conditioning on  $C$  opens an additional path between  $D$  and  $Y$ , creating a spurious association that violates the back-door criterion. In particular, the variable  $C$  is a collider on the path  $D \rightarrow C \leftarrow \mathbf{X}$ . By default, a collider *blocks* the flow of association along the path on which it lies, so this path is closed when  $C$  is not conditioned on. However, conditioning on  $C$  (or on any descendant of  $C$ ) opens this path, creating a spurious association between  $D$  and  $\mathbf{X}$ . Since  $\mathbf{X}$  also affects  $Y$ , this induces bias in estimating the causal effect of  $D$  on  $Y$ .

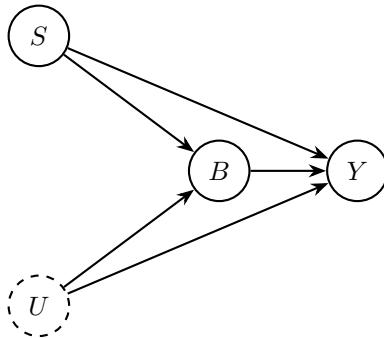
#### Example: Birth-weight “paradox”, collider bias

The collider bias illustrated in DAG 13.9, based on [84], reflects the well-known Birth-weight “paradox”. This phenomenon arises when conditioning on an intermediate variable (birth weight,  $B$ ) induces a spurious association between the binary indicator of smoking ( $S$ ) and infant mortality ( $Y$ ). In this setting,  $U$  represents unobserved factors that affect both birth weight and

**FIGURE 13.8**

Directed Acyclic Graph (DAG): The additional collider  $C$  caused by both  $D$  and  $\mathbf{X}$  would open a path between  $D$  and  $\mathbf{X}$ , inducing bias.

infant mortality. Conditioning on  $B$  creates collider bias because it opens a back-door path through the unobserved factors  $U$  (dashed).

**FIGURE 13.9**

DAG illustrating the birth-weight paradox:  $S$  (smoking),  $B$  (birth weight),  $Y$  (infant mortality), and  $U$  (other unobserved health factors).

The code in 13.3 illustrates the bias by performing 100 replications, and Figure 13.10 shows the distribution of posterior means across these replications. In this setting, the total effect of smoking on infant mortality is 2, consisting of a direct effect of 1 and an indirect effect through birth weight, also equal to 1.

#### *R code. Collider bias: Birth weight “paradox”*

```

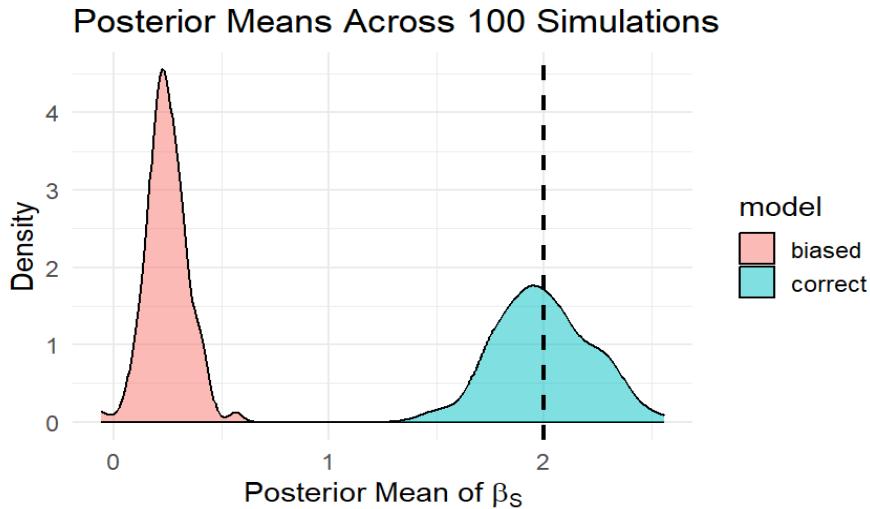
1 rm(list = ls()); set.seed(10101)
2 library(MCMCpack); library(dplyr); library(ggplot2)
3 # Parameters
4 n <- 1000; true_effect <- 2; replications <- 100
  
```

### R code. Collider bias: Birth weight “paradox”

```

1 # Store results
2 results <- data.frame(rep = 1:replications, correct =
  numeric(replications), biased = numeric(replications))
3 for (r in 1:replications) {
4   # Simulate data under DAG: S -> B -> Y, S -> Y, U -> B, U
   -> Y
5   U <- rnorm(n, 0, 1); S <- rbinom(n, prob = 0.7, size = 1)
6   B <- S + U + rnorm(n)
7   Y <- S + B + 1.5*U + rnorm(n)
8   # Correct model: does NOT condition on collider B
9   model_correct <- MCMCregress(Y ~ S, burnin = 1000, mcmc =
  3000, verbose = FALSE)
10  # Biased model: conditions on collider B
11  model_biased <- MCMCregress(Y ~ S + B, burnin = 1000, mcmc =
  3000, verbose = FALSE)
12  # Posterior means for S
13  results$correct[r] <- mean(as.matrix(model_correct)[, "S"])
14  results$biased[r] <- mean(as.matrix(model_biased)[, "S"])
15 }
16 # Compute bias
17 results <- results %>% mutate(
18   bias_correct = correct - true_effect,
19   bias_biased = biased - true_effect
20 )
21 # Average bias and SD
22 avg_bias <- results %>%
23 summarise(
24   mean_correct = mean(bias_correct),
25   mean_biased = mean(bias_biased),
26   sd_correct = sd(bias_correct),
27   sd_biased = sd(bias_biased)
28 )
29 print(avg_bias)
30 # Visualization: distribution of posterior means across 100
  simulations
31 df_long <- results %>%
32 select(rep, correct, biased) %>%
33 tidyr::pivot_longer(cols = c(correct, biased), names_to =
  "model", values_to = "estimate")
34 ggplot(df_long, aes(x = estimate, fill = model)) + geom_
  density(alpha = 0.5) + geom_vline(xintercept = true_
  effect, color = "black", linetype = "dashed", linewidth =
  1) + labs(title = "Posterior Means Across 100
  Simulations", x = expression(paste("Posterior Mean of ",
  beta[S])), y = "Density") + theme_minimal(base_size =
  14)

```

**FIGURE 13.10**

Posterior distribution: Posterior means of the collider bias example.

Collider (selection) bias is one possible source of bias in the identification of causal effects. The well-known Heckman's sample selection problem in econometrics [198] can be interpreted as a form of collider bias because restricting the analysis to selected observations (e.g., those with positive wages) conditions on a variable that is a common effect of observed and unobserved factors, thereby opening a non-causal path and creating bias. In other words, the sample is no longer representative of the population, which undermines the identification of the causal effect.

Other well-known sources of bias in econometrics include the omission of common causes affecting both the treatment and the outcome, that is, *omission of correlated relevant regressors*, measurement error in regressors leading to *attenuation bias* (where the estimated causal effect is biased toward zero), and *simultaneous causality*, which often arises in systems of equations (see [455] for details). We will discuss these additional sources of bias later.

## 13.4 Instrumental variables (IV)

In many real-world situations, we face biases such as those described in the previous paragraph, arising from measurement errors, omission of relevant variables, and similar issues. Even in these cases, it is still possible to identify

the causal effect. One common strategy is to use a set of *instrumental variables* ( $\mathbf{Z}$ ) that satisfy two key conditions:

1. *Relevance*, meaning the instruments are correlated with the treatment:

$$\mathbf{Z}_i \not\perp D_i.$$

2. *Exogeneity*, which in the linear model requires:

$$\mathbb{E}[\mu_i | \mathbf{Z}_i] = 0,$$

and, in terms of potential outcomes, entails:

- *Exclusion restriction*: instruments affect the outcome only through the treatment:

$$Y_i(D_i = d, \mathbf{Z}_i = \mathbf{z}) = Y_i(D_i = d, \mathbf{Z}_i = \mathbf{z}') = Y_i(D_i = d), \quad d \in \{0, 1\},$$

- *Marginal exchangeability*: instruments are independent of the potential outcomes:

$$\mathbf{Z}_i \perp \{Y_i(1), Y_i(0)\}.$$

Relevance is testable, typically by checking whether instruments significantly predict the endogenous variable [421, 107, 244, 425]. However, *weak instruments*, those that exhibit only a weak association with the treatment, can lead to serious consequences, including a high level of uncertainty and the potential amplification of even small biases when estimating the causal effect.

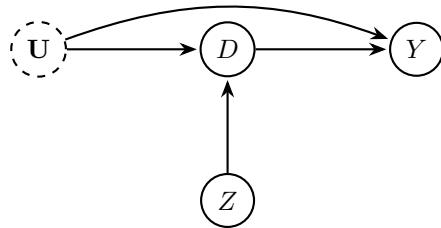
Exogeneity is fundamentally untestable [215]. However, when the number of instruments exceeds the number of endogenous variables, the Sargan test (or its robust version, the Hansen  $J$ -test) can be used to assess the validity of overidentifying restrictions. This test should not be misinterpreted as a direct test of the exclusion restriction. Instead, it evaluates whether the overidentifying restrictions implied by the IV model hold. Specifically, the null hypothesis states that all instruments are jointly uncorrelated with the structural error term,  $\mathbb{E}[\mathbf{Z}^\top \mu] = \mathbf{0}$  [395, 192].

If the test rejects, it indicates that at least one instrument is invalid, but it does not reveal whether the violation arises from a failure of the exclusion restriction, correlation with unobserved confounders, or model misspecification. Conversely, if the test fails to reject, it only suggests that the sample moment conditions do not provide evidence against instrument validity under the maintained model. This outcome does not prove that the exclusion restrictions hold: the test may have low power, a mix of valid and invalid instruments could pass if violations “cancel out”, and if the model is misspecified, a “passing” test is uninformative.

In general, the exclusion restriction concerns unobservables (no correlation with the error term). Since the error term is not observed, the exclusion

restriction can never be proven from the data; we can only look for evidence against it. Therefore, while the Sargan test is informative about overall instrument validity, it is not a separate or definitive test of the exclusion restriction [455, 197].

Figure 13.11 illustrates a situation where pre-treatment variables that influence both the treatment and the outcome are partially observed or measured with error. In such cases, an instrument can help identify the causal effect of the treatment on the outcome.



**FIGURE 13.11**

Directed Acyclic Graph (DAG) showing an instrumental variable  $Z$  used to identify the causal effect of  $D$  on  $Y$  when some confounders  $\mathbf{U}$  (dashed) are unobserved or measured with error.

However, the two conditions for instruments (relevance and exogeneity) only imply *partial identification* of the Local Average Treatment Effect (LATE; see below, and [288, 289, 290, 291] for detailed treatments of partial identification). Thus, it is necessary to impose a third condition to achieve *point identification* of a particular treatment effect using IV. This condition is *monotonicity* which asserts that the instrument moves all units' treatment decisions in the same direction, that is,

$$D_i(1) \geq D_i(0), i = 1, 2, \dots, N,$$

assuming a binary instrument  $Z_i = \{0, 1\}$ .

Note that monotonicity is also fundamentally untestable [215], as it refers to unobserved potential outcomes. However, its validity can be argued in the context of each application. We can examine whether the instrument shifts treatment in the expected direction. Moreover, if rich covariates are available, we can stratify by them and verify that, within each subgroup, the instrument increases (or at least does not decrease) the probability of treatment. For instance, we can analyze the cumulative distribution functions of the outcome when the binary instrumental variable takes values of 0 and 1; if the functions do not cross, this would suggest plausibility of the monotonicity assumption [359].

In the case where  $Z_i$  denotes the assignment to treatment and  $D_i$  the actual treatment status, there are four potential groups: *always-takers* ( $D = 1 | Z = 1$  and  $D = 1 | Z = 0$ ), *never-takers* ( $D = 0 | Z = 1$  and  $D = 0 | Z = 0$ ),

*compliers* ( $D = 1 \mid Z = 1$  and  $D = 0 \mid Z = 0$ ), and *defiers* ( $D = 0 \mid Z = 1$  and  $D = 1 \mid Z = 0$ ). Note that these groups are not identified; for instance, we do not know whether an individual with  $Z_i = 1$  and  $D_i = 1$  is a complier or an always taker. The monotonicity assumption rules out defiers, meaning that the instrument (assignment to treatment) either does not change treatment status or only increases it.

Table 13.2 displays the distribution of unit types when defiers are ruled out by the monotonicity assumption.

**TABLE 13.2**  
Distribution of compliance types by instrument and treatment

Instrument	Treatment	
	$D_i = 0$	$D_i = 1$
$Z_i = 0$	Compliers, Never-takers	Always-takers
$Z_i = 1$	Never-takers	Compliers, Always-takers

From this table we can identify the proportions of compliance types ( $C$ ) as follows. The probability of always-takers (a) is

$$P(C_i = a) = P(D_i = 1 \mid Z_i = 0),$$

and the probability of never-takers (n) is

$$P(C_i = n) = P(D_i = 0 \mid Z_i = 1).$$

For compliers (c), note that

$$P(D_i = 1 \mid Z_i = 1) = P(C_i = c) + P(C_i = a),$$

so that

$$P(C_i = c) = P(D_i = 1 \mid Z_i = 1) - P(D_i = 1 \mid Z_i = 0).$$

Equivalently, using  $D_i = 0$ ,

$$P(D_i = 0 \mid Z_i = 0) = P(C_i = n) + P(C_i = c),$$

so that

$$P(C_i = c) = P(D_i = 0 \mid Z_i = 0) - P(D_i = 0 \mid Z_i = 1).$$

Therefore, under relevance, exogeneity, monotonicity, and SUTVA, IV methods allow identification of the causal effect for the subgroup of *compliers*, that is, individuals who take the treatment when encouraged by the instrument and do not take it otherwise [215, 16]. This estimand is known as the *Local Average Treatment Effect* (LATE) because it applies to a specific subpopulation rather than the entire population [204]:

$$\tau_{LATE} = \mathbb{E}[Y_i(1) - Y_i(0) \mid D_i(1) = 1, D_i(0) = 0].$$

We define  $Y_i(z_i, D_i(z))$  to be the outcome for unit  $i$  if exposed to treatment  $D_i(z)$  after being assigned to treatment  $z$ . Therefore, the *Intention-to-Treat* (ITT) causal effect of  $Z_i$  on  $Y_i$  is  $\mathbb{E}[Y_i(1, D_i(1)) - Y_i(0, D_i(0))]$ , and can be expressed as:

$$\begin{aligned} & \underbrace{\mathbb{E}[Y_i(1, D_i(1)) - Y_i(0, D_i(0)) \mid D_i(1) = 1, D_i(0) = 1]}_{\text{always-takers}} \times P(D_i(1) = 1, D_i(0) = 1) \\ & + \underbrace{\mathbb{E}[Y_i(1, D_i(1)) - Y_i(0, D_i(0)) \mid D_i(1) = 0, D_i(0) = 0]}_{\text{never-takers}} \times P(D_i(1) = 0, D_i(0) = 0) \\ & + \underbrace{\mathbb{E}[Y_i(1, D_i(1)) - Y_i(0, D_i(0)) \mid D_i(1) = 1, D_i(0) = 0]}_{\text{compliers}} \times P(D_i(1) = 1, D_i(0) = 0) \\ & + \underbrace{\mathbb{E}[Y_i(1, D_i(1)) - Y_i(0, D_i(0)) \mid D_i(1) = 0, D_i(0) = 1]}_{\text{defiers}} \times P(D_i(1) = 0, D_i(0) = 1), \end{aligned}$$

that is, ITT is a mixture over the four potential groups.

We know that the ITT effect of  $Z_i$  on  $Y_i$  is zero for always-takers and never-takers due to the exclusion restriction, i.e.,  $Z_i$  affects  $Y_i$  only through  $D_i$ , and  $D_i$  does not vary with  $Z_i$  for always-takers and never-takers. Moreover, the monotonicity assumption rules out the existence of defiers. Hence, only compliers contribute to the effect of  $Z_i$  on  $Y_i$ , then  $\mathbb{E}[Y_i(1, D_i(1)) - Y_i(0, D_i(0))]$  is

$$\mathbb{E}[Y_i(1, D_i(1)) - Y_i(0, D_i(0)) \mid D_i(1) = 1, D_i(0) = 0] \cdot P(D_i(1) = 1, D_i(0) = 0).$$

For compliers,  $D_i = Z_i$ , so we can simplify  $Y_i(1, D_i(1)) = Y_i(1, 1) = Y_i(1)$  and  $Y_i(0, D_i(0)) = Y_i(0, 0) = Y_i(0)$ . Thus,

$$\mathbb{E}[Y_i(1) - Y_i(0)] = \mathbb{E}[Y_i(1) - Y_i(0) \mid D_i(1) = 1, D_i(0) = 0] \cdot P(D_i(1) = 1, D_i(0) = 0),$$

which implies that

$$\tau_{LATE} = \mathbb{E}[Y_i(1) - Y_i(0) \mid D_i(1) = 1, D_i(0) = 0] = \frac{\mathbb{E}[Y_i(1) - Y_i(0)]}{P(D_i(1) = 1, D_i(0) = 0)}.$$

Under the assumption that the instrument  $Z_i$  is independent of the potential outcomes (i.e., marginal exchangeability), we have

$$\tau_{LATE} = \mathbb{E}[Y_i(1) - Y_i(0) \mid D_i(1) = 1, D_i(0) = 0] = \frac{\mathbb{E}[Y_i \mid Z_i = 1] - \mathbb{E}[Y_i \mid Z_i = 0]}{P(D_i(1) = 1, D_i(0) = 0)}.$$

We know that  $P(D_i(1) = 1, D_i(0) = 0) = P(D_i = 1 \mid Z_i = 1) - P(D_i =$

$1 \mid Z_i = 0$ ), and taking into account that  $\mathbb{E}[D_i \mid Z_i] = P(D_i = 1 \mid Z_i)$ , we obtain

$$\tau_{LATE} = \frac{\mathbb{E}[Y_i \mid Z_i = 1] - \mathbb{E}[Y_i \mid Z_i = 0]}{\mathbb{E}[D_i \mid Z_i = 1] - \mathbb{E}[D_i \mid Z_i = 0]}.$$

This is the standard instrumental variables (IV) estimand, often referred to as the *Wald estimator*. Note that the *relevance condition* is required to ensure identification, i.e.,

$$\mathbb{E}[D_i \mid Z_i = 1] - \mathbb{E}[D_i \mid Z_i = 0] \neq 0,$$

that is, the instrument is strong enough that some units actually comply, this means, there is a proportion of compliers in the population.

The IV estimand represents the ratio of the intention-to-treat effect on the outcome to the intention-to-treat effect on the treatment (i.e., the proportion of compliers). The instrument induces random variation in treatment: the numerator captures how outcomes respond to the instrument, and the denominator captures how treatment responds to the instrument. Dividing the two isolates the causal effect for compliers.

Note that the LATE is not the global average causal effect, and has been subject to several criticisms. First, the proportion of compliers in the population may be small, raising concerns about the policy relevance of this estimand. Observe that while we cannot identify all four principal strata (compliers, never-takers, always-takers, and defiers), the proportion of compliers is identifiable. Second, the monotonicity assumption is not always plausible, particularly in observational studies where instruments may induce heterogeneous behavioral responses. Third, partitioning the population into four strata can become ill-defined in settings where instruments arise from different mechanisms or individual-specific criteria [114, 204]. For arguments in defense of LATE against these critiques, see [17].

**Example: A simple setting for point and set identified LATE**

Assume that we see the results display in Table 13.3.

**TABLE 13.3**

Distribution of units by type, instrument, treatment, and outcomes

Type	$Z_{obs}$	$D_{obs}$	$Y_{obs}$	Units
Compliers, Never-takers	0	0	3	50
Never-takers	1	0	4	20
Compliers, Always-takers	1	1	7	100
Always-takers	0	1	5	30

Let's calculate the LATE assuming exclusion restrictions, and without assuming exclusion restrictions.

We first estimate the *first stage and reduced form*.

$$E[D | Z = 0] = \frac{30}{80} = 0.37, \quad E[D | Z = 1] = \frac{100}{120} = 0.83,$$

$$\Delta_D = E[D | Z = 1] - E[D | Z = 0] = 0.46.$$

$$E[Y | Z = 0] = \frac{3 \cdot 50 + 5 \cdot 30}{80} = 3.75, \quad E[Y | Z = 1] = \frac{4 \cdot 20 + 7 \cdot 100}{120} = 6.50,$$

$$\Delta_Y = E[Y | Z = 1] - E[Y | Z = 0] = 2.75.$$

From these we also recover the compliance shares:

$$\eta_a = P(D = 1 | Z = 0) = \frac{3}{8} = 0.37, \quad \eta_c = \Delta_D = \frac{11}{24} = 0.46,$$

thus,

$$\eta_n = 1 - \eta_a - \eta_c = \frac{1}{6} = 0.17.$$

Then, we can *point-identified* the LATE under exclusion restriction,

$$\text{LATE} = \frac{\Delta_Y}{\Delta_D} = \frac{2.75}{0.46} = 6.$$

We achieved point identification by imposing the exclusion restriction, which requires that treatment assignment is unrelated to potential outcomes for never-takers and always-takers.<sup>2</sup> This, is

$$E[Y | Z = 1, C = n] = E[Y | Z = 0, C = n],$$

and

$$E[Y | Z = 1, C = a] = E[Y | Z = 0, C = a].$$

Without exclusion restrictions, we should take into account that

$$\delta_n = E[Y | Z = 1, C = n] - E[Y | Z = 0, C = n] \neq 0,$$

$$\delta_a = E[Y | Z = 1, C = a] - E[Y | Z = 0, C = a] \neq 0.$$

Then

$$\Delta_Y = \eta_c \cdot \text{LATE} + \eta_n \delta_n + \eta_a \delta_a,$$

Thus,

$$\text{LATE} = \frac{\Delta_Y - \eta_n \delta_n - \eta_a \delta_a}{\eta_c}.$$

From the table we know

$$E[Y | Z = 1, C = n] = 4, \quad E[Y | Z = 0, C = a] = 5.$$

---

<sup>2</sup>Note that under monotonicity, defiers are ruled out, whereas for compliers, the exclusion restriction is not a separate assumption: by construction, their potential outcomes are indexed only by treatment status, since the instrument deterministically shifts them between  $D_i = 0$  and  $D_i = 1$ . Thus, the exclusion restriction applies only to never-takers and always-takers, whose treatment status does not vary with the instrument.

We assume that the bounding outcomes come from the observed support  $Y \in [3, 7]$ . Then,

$$\delta_n \in [-3, 1], \quad \delta_a \in [-2, 2].$$

Therefore,

$$\text{LATE}_{\min} = \frac{\Delta_Y - \eta_n \cdot 1 - \eta_a \cdot 2}{\eta_c} = 4,$$

$$\text{LATE}_{\max} = \frac{\Delta_Y - \eta_n(-3) - \eta_a(-2)}{\eta_c} = 8.73.$$

Then, without exclusion restrictions the LATE is *set identified*,  $\text{LATE} \in [4, 8.73]$ , and with exclusion restrictions, LATE is *point identified*,  $\text{LATE} = 6$ .

#### Example: Effects of vitamin A supplements on children's survival

This example is taken from [217]. The authors perform Bayesian inference for causal effects in a completely randomized experiment with no covariates and partial compliance. The experiment was conducted in Indonesia, where children were randomly assigned to receive vitamin A supplements. No individual assigned to the control group ( $Z_i = 0$ ) actually took the vitamin, meaning that  $D_i(0) = 1$  was never observed. This rules out always-takers and defiers. However, some individuals assigned to the treatment group ( $Z_i = 1$ ) did not take the vitamin, so  $D_i(1) = 0$  occurred for some units. Consequently,  $D_i(1) \geq D_i(0)$  holds for all individuals (there are no defiers), satisfying the monotonicity assumption. Under this setting, the population consists only of compliers and never-takers.

**TABLE 13.4**

Effects of vitamin A supplements on children's survival [217]

Type	Assignment $Z_{\text{obs},i}$	Treatment $D_{\text{obs},i}$	Survival $Y_{\text{obs},i}$	Units 23,682
Complier or never-taker	0	0	0	74
Complier or never-taker	0	0	1	11,514
Never-taker	1	0	0	34
Never-taker	1	0	1	2,385
Complier	1	1	0	12
Complier	1	1	1	9,663

Table 13.4 summarizes the data. Using this table, we can illustrate key concepts in the identification of causal effects with noncompliance. Let  $C_i$  denote the compliance type of individual  $i$ , which in this case can be either a complier ( $c$ ) or a never-taker ( $n$ ). The probability of being a complier,  $\omega := P(C_i = c)$ , is given by

$$P(C_i = c) = \underbrace{P(D_i = 1 | Z_i = 1)}_{P(\text{compliers}) + P(\text{always-takers})} - \underbrace{P(D_i = 1 | Z_i = 0)}_{P(\text{always-takers})=0}$$

$$= \frac{P(D_i = 1, Z_i = 1)}{P(Z_i = 1)}.$$

Thus, the ML estimate is

$$\hat{\omega} = \frac{12 + 9,663}{12 + 9,663 + 34 + 2,385} = 0.8.$$

Note that  $P(C_i = n) = 1 - P(C_i = c)$  in this example; thus, the ML estimate is  $1 - \hat{\omega} = 0.2$ .

The probability of survival for compliers assigned to treatment,  $\eta_{c1} := P(Y_i = 1 | C_i = c, Z_i = 1)$ , is

$$P(Y_i = 1 | C_i = c, Z_i = 1) = \frac{P(Y_i = 1, C_i = c | Z_i = 1) \times P(Z_i = 1)}{P(C_i = c | Z_i = 1) \times P(Z_i = 1)}.$$

Thus, the ML estimate is  $\hat{\eta}_{c1} = \frac{9,663}{9,663+12} = 0.999$ .

Similarly, the probability of survival for never-takers assigned to treatment,  $\eta_{n1} := P(Y_i = 1 | C_i = n, Z_i = 1)$ , is

$$P(Y_i = 1 | C_i = n, Z_i = 1) = \frac{P(Y_i = 1, C_i = n | Z_i = 1) \times P(Z_i = 1)}{P(C_i = n | Z_i = 1) \times P(Z_i = 1)}.$$

Consequently, the ML estimate is  $\hat{\eta}_{n1} = \frac{2,385}{2,385+34} = 0.986$ .

Note that  $\eta_{c0} = P(Y_i = 1 | C_i = c, Z_i = 0)$  and  $\eta_{n0} = P(Y_i = 1 | C_i = n, Z_i = 0)$  cannot be *point-identified*. However, we can obtain *set identification* because

$$\begin{aligned} \eta_0 := P(Y_i = 1 | Z_i = 0) &= P(Y_i = 1 | Z_i = 0, C_i = c) \times P(C_i = c) \\ &\quad + P(Y_i = 1 | Z_i = 0, C_i = n) \times P(C_i = n) \\ &= \frac{P(Y_i = 1, Z_i = 0)}{P(Z_i = 0)}, \end{aligned}$$

where the ML estimate is  $\hat{\eta}_0 = \frac{11,514}{11,514+74} = 0.9936$ .

The first equality holds because the probability of survival given no treatment assignment is expressed as a mixture by the law of total probability, and the second equality follows from Bayes' rule.

This implies that

$$0.9936 = \hat{\omega}\hat{\eta}_{c0} + (1 - \hat{\omega})\hat{\eta}_{n0},$$

which leads to

$$\hat{\eta}_{n0} = 4.968 - 4\hat{\eta}_{c0}.$$

Given that probabilities must lie in the unit interval, such that  $0 \leq 4.968 - 4\hat{\eta}_{c0} \leq 1$ , we obtain: (i)  $\hat{\eta}_{c0} \leq 1 \leq 1.242$  (automatically satisfied), and (ii)  $4\hat{\eta}_{c0} \geq 4.968 - 1 = 3.968$ , which implies  $\hat{\eta}_{c0} \geq 0.992$ .

Consequently, the ML estimates  $\hat{\eta}_{c0} \in [0.992, 1]$  and  $\hat{\eta}_{n0} \in [0.968, 1]$ .

These results imply that the *complier average causal effect* (CACE) is

$$\begin{aligned}\tau_{CACE} &= \mathbb{E}[Y_i(1) - Y_i(0) \mid D_i(1) = 1, D_i(0) = 0] \\ &= P[Y_i = 1 \mid C_i = c, Z_i = 1] - P[Y_i = 1 \mid C_i = c, Z_i = 0] \\ &= \eta_{c1} - \eta_{c0},\end{aligned}$$

where the first line is the formal definition, the second equality uses random assignment to express the estimand in terms of observed probabilities, and the third relies on the model parameterization.

The CACE is only set-identified, yielding the interval

$$\hat{\tau}_{CACE} \in [-0.001, 0.007],$$

meaning that the survival rate among compliers can vary between -1 and 7 per 1,000 children due to vitamin A supplementation.

Remember that point identification can be achieved by imposing the exclusion restriction, which requires that treatment assignment is unrelated to potential outcomes for never-takers and always-takers. Under this assumption,

$$P(Y_i = y \mid C_i = n, Z_i = z) = P(Y_i = y \mid C_i = n),$$

so that  $\hat{\eta}_{n1} = \hat{\eta}_{n0} = \hat{\eta}_n = 0.986$ . Substituting into the mixture equation yields  $\hat{\eta}_{c0} = 0.9955$ , and consequently,

$$\hat{\tau}_{CACE} = 0.0035,$$

indicating that survival increases by approximately 3.5 per 1,000 children due to vitamin A supplementation.<sup>3</sup>

We use conjugate families to perform Bayesian inference in this example. Specifically, we adopt the Bernoulli-Beta model (see Section 3.2) with independent non-informative Beta priors, each having parameters equal to 1. This corresponds to a uniform distribution on the interval (0, 1).

An advantage of the Bayesian formulation is that, after applying data augmentation with the compliance type  $C_i$ , all the conditional posterior distributions are Beta (See [217] for details in derivations). In particular:

$$\omega \mid \mathbf{C}, \mathbf{Z}_{\text{obs}}, \mathbf{D}_{\text{obs}}, \mathbf{Y}_{\text{obs}} \sim \text{Beta}(1 + N_c, 1 + N_n),$$

where  $N_c$  and  $N_n$  denote the number of compliers and never-takers, respectively.

$$\eta_{c1} \mid \mathbf{C}, \mathbf{Z}_{\text{obs}}, \mathbf{D}_{\text{obs}}, \mathbf{Y}_{\text{obs}} \sim \text{Beta}(1 + 9,663, 1 + 12),$$

$$\eta_{n1} \mid \mathbf{C}, \mathbf{Z}_{\text{obs}}, \mathbf{D}_{\text{obs}}, \mathbf{Y}_{\text{obs}} \sim \text{Beta}(1 + 2,385, 1 + 34),$$

---

<sup>3</sup>CACE and LATE refer to the same estimand -the average causal effect for compliers-under monotonicity and the exclusion restriction. If the exclusion restriction fails, they differ because LATE (as identified by IV) no longer equals the complier-specific causal effect (CACE).

$$\eta_{c0} \mid \mathbf{C}, \mathbf{Z}_{\text{obs}}, \mathbf{D}_{\text{obs}}, \mathbf{Y}_{\text{obs}} \sim \text{Beta}(1 + N_{c01}, 1 + N_{c00}),$$

where  $N_{c01}$  and  $N_{c00}$  are the numbers of compliers in the control group who survived and did not survive, respectively.

$$\eta_{n0} \mid \mathbf{C}, \mathbf{Z}_{\text{obs}}, \mathbf{D}_{\text{obs}}, \mathbf{Y}_{\text{obs}} \sim \text{Beta}(1 + N_{n01}, 1 + N_{n00}),$$

where  $N_{n01}$  and  $N_{n00}$  are the numbers of never-takers in the control group who survived and did not survive, respectively.

In addition, the conditional probability of being a complier is given by:

$$P(C_i = c \mid Z_{\text{obs},i}, D_{\text{obs},i}, Y_{\text{obs},i}) = \begin{cases} 0, & i \in \{Z_i = 1, D_i = 0\}, \\ 1, & i \in \{Z_i = 1, D_i = 1\}, \\ \frac{\omega g_{c0,i}}{\omega g_{c0,i} + (1 - \omega)g_{n0,i}}, & i \in \{Z_i = 0\}, \end{cases}$$

where  $g_{c0,i} = \eta_{c0}^{Y_{\text{obs},i}}(1 - \eta_{c0})^{1 - Y_{\text{obs},i}}$  and  $g_{n0,i} = \eta_{n0}^{Y_{\text{obs},i}}(1 - \eta_{n0})^{1 - Y_{\text{obs},i}}$ . Note that  $Y_{\text{obs},i} = 1$  implies  $g_{c0,i} = \eta_{c0}$  and  $g_{n0,i} = \eta_{n0}$ , while  $Y_{\text{obs},i} = 0$  implies  $g_{c0,i} = (1 - \eta_{c0})$  and  $g_{n0,i} = (1 - \eta_{n0})$ .

The following code shows the implementation, and the corresponding figure displays the posterior distribution of the CACE, the mean is 0.0024, and the 95% credible intervals is (-0.0012, 0.0071).

### R code. Effects of vitamin A supplements on children's survival

```

1 rm(list = ls()); set.seed(10101)
2 library(dplyr)
3 # Simulate data
4 Nc111 <- 9663
5 c111 <- cbind(rep(1, Nc111), rep(1, Nc111), rep(1, Nc111))
6 Nc110 <- 12
7 c110 <- cbind(rep(1, Nc110), rep(1, Nc110), rep(0, Nc110))
8 Nn101 <- 2385
9 n101 <- cbind(rep(1, Nn101), rep(0, Nn101), rep(1, Nn101))
10 Nn100 <- 34
11 n100 <- cbind(rep(1, Nn100), rep(0, Nn100), rep(0, Nn100))
12 Ncn001 <- 11514
13 cn001 <- cbind(rep(0, Ncn001), rep(0, Ncn001), rep(1, Ncn001))
14 Ncn000 <- 74
15 cn000 <- cbind(rep(0, Ncn000), rep(0, Ncn000), rep(0, Ncn000))

```

***R code. Effects of vitamin A supplements on children's survival***

---

```

1 mydata <- rbind(c111, c110, n101, n100, cn001, cn000)
2 mydata <- data.frame(Z = mydata[,1], D = mydata[,2], Y =
   mydata[,3])
3 N <- dim(mydata)[1]
4 attach(mydata)
5 # Sampling function C (type)
6 SampleType <- function(z, d, y, wc, nc0, nn0){
7   if(z == 1 & d == 0){
8     pc <- 0
9   }else{
10     if(z == 1 & d == 1){
11       pc <- 1
12     }else{
13       if(y == 1){
14         pc <- (wc * nc0) / (wc * nc0 + (1 - wc) * nn0)
15       }else{
16         pc <- (wc * (1 - nc0)) / (wc * (1 - nc0) + (1 - wc)
* (1 - nn0))
17       }
18     }
19   }
20   rbinom(1, 1, prob = pc) # 1: Complier/ 0: Never taker
21 }
22 z = 0; d = 0; y = 0; wc = 0.8; nc0 = 0.9; nn0 = 0.05
23 SampleType(z = z, d = d, y = y, wc = wc, nc0 = nc0, nn0 =
   nn0)
24 Clat <- sapply(1:N, function(i){SampleType(z = Z[i], d =
   D[i],
   ], y = Y[i], wc = wc, nc0 = nc0, nn0 = nn0)})
25 # Gibbs sampler
26 a0 <- 1; b0 <- 1 # Hyperparameters beta priors
27 burnin <- 500; S <- 2000; tot <- S + burnin
28 PosteriorDraws <- matrix(NA, tot, 5)
29 pb <- txtProgressBar(min = 0, max = tot, style = 3)
30 for(s in 1:tot){
31   dataLat <- cbind(mydata, Clat)
32   Nc011 <- sum(dataLat$Z == 0 & dataLat$Clat == 1 & dataLat$Y
   == 1)
33   Nc010 <- sum(dataLat$Z == 0 & dataLat$Clat == 1 & dataLat$Y
   == 0)
34   Nn001 <- sum(dataLat$Z == 0 & dataLat$Clat == 0 & dataLat$Y
   == 1)
35   Nn000 <- sum(dataLat$Z == 0 & dataLat$Clat == 0 & dataLat$Y
   == 0)
36   Nc <- sum(Clat == 1)
37   Nn <- sum(Clat == 0)

```

---

**R code. Effects of vitamin A supplements on children's survival**

```

1  wc <- rbeta(1, 1 + Nc, 1 + Nn)
2  nc1 <- rbeta(1, 1 + Nc111, 1 + Nc110)
3  nn1 <- rbeta(1, 1 + Nn101, 1 + Nn100)
4  nc0 <- rbeta(1, 1 + Nc011, 1 + Nc010)
5  nn0 <- rbeta(1, 1 + Nn001, 1 + Nn000)
6  Clat <- sapply(1:N, function(i){SampleType(z = Z[i], d = D
    [i], y = Y[i], wc = wc, nc0 = nc0, nn0 = nn0)})
7  PosteriorDraws[s, ] <- c(wc, nc1, nc0, nn1, nn0)
8  setTxtProgressBar(pb, s)
9 }
10 close(pb)
11 keep <- seq((burnin+1), tot)
12 LATE <- PosteriorDraws[keep, 2] - PosteriorDraws[keep, 3]
13 LATEmean <- mean(LATE)
14 LATEci <- quantile(LATE, c(0.025, 0.975))
15 # Plot posterior distribution of CATE
16 hist(LATE, breaks = 40, freq = FALSE, main = "Posterior
    Distribution of CACE",
17 xlab = "CACE", col = "lightblue", border = "white")
18 abline(v = LATEmean, col = "red", lwd = 2)
19 abline(v = LATEci, col = "darkgreen", lty = 2, lwd = 2)
20 legend("topright", legend = c("Posterior Mean", "95%
    Credible Interval"), col = c("red", "darkgreen"), lwd =
    2, lty = c(1, 2), bty = "n", cex = 0.8)
21

```

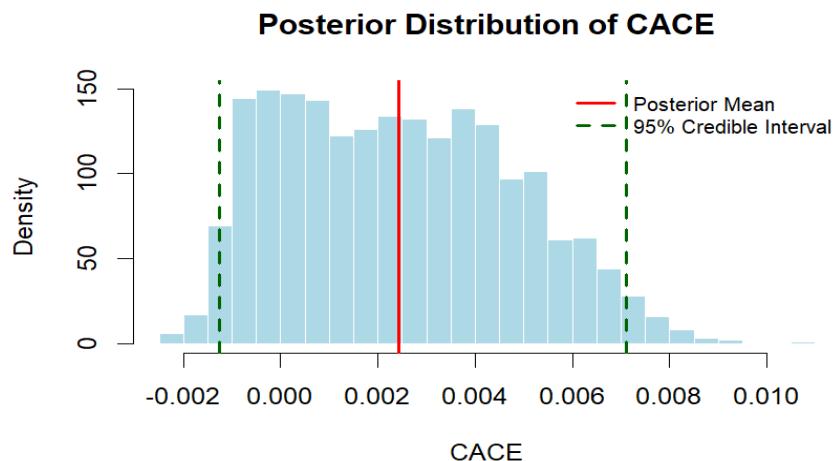
We can perform inference by imposing the exclusion restriction, setting  $\eta_n = \eta_{n0} = \eta_{n1}$ , so that the posterior distribution is

$$\eta_n | \mathbf{C}, \mathbf{Z}_{\text{obs}}, \mathbf{D}_{\text{obs}}, \mathbf{Y}_{\text{obs}} \sim \text{Beta}(1 + N_{n1}, 1 + N_{n0}),$$

where  $N_{n1}$  and  $N_{n0}$  denote the numbers of never-takers who survived and did not survive, respectively. This is the only modification to the Gibbs sampler.

The following code shows the implementation, and Figure 13.13 displays the posterior distribution of the CACE (which equals the LATE under the exclusion restriction). The posterior mean is 0.0030, and the 95% credible interval is (0.0008, 0.0054). Therefore, imposing the exclusion restriction leaves the posterior mean approximately unchanged but increases the precision and informativeness of the posterior distribution. Without the restriction, the posterior was relatively flat within the 95% credible interval, whereas under the exclusion restriction, the posterior distribution is approximately normal, resulting in a narrower 95% credible interval.

Thus, this example illustrates how the Bayesian framework provides a

**FIGURE 13.12**

Posterior distribution: Complier average causal effect of vitamin A supplements.

way to quantify uncertainty and conduct sensitivity analysis regarding the exclusion restriction. See [206] for an extension controlling for covariates.

In addition, it is worth mentioning that in partially identified models, the asymptotic equivalence between Bayesian and Frequentist inference breaks down. In particular, credible sets for partially identified parameters tend to be smaller than confidence sets asymptotically [313]. This occurs because the posterior remains asymptotically sensitive to the prior, and Bayesian credible sets lie strictly inside the true identified set. To address this disagreement, [176] propose a multi-prior robust Bayesian approach that helps reconcile Bayesian and Frequentist inference in set-identified models, which are a special case of partially identified models.

*R code. Effects of vitamin A supplements on children's survival: Imposing exclusion restrictions*

```
1 PosteriorDrawsER <- matrix(NA, tot, 4)
2 pb <- txtProgressBar(min = 0, max = tot, style = 3)
```

**R code. Effects of vitamin A supplements on children's survival: Imposing exclusion restrictions**

```

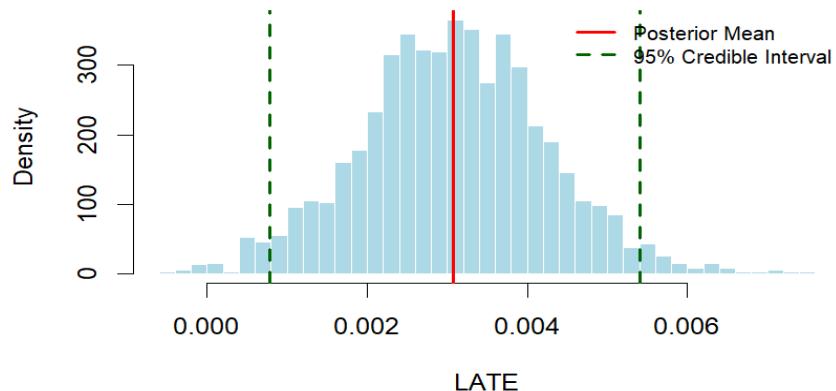
1 for(s in 1:tot){
2   dataLat <- cbind(mydata, Clat)
3   Nc011 <- sum(dataLat$Z == 0 & dataLat$Clat == 1 & dataLat$Y == 1)
4   Nc010 <- sum(dataLat$Z == 0 & dataLat$Clat == 1 & dataLat$Y == 0)
5   Nn01 <- sum(dataLat$Clat == 0 & dataLat$Y == 1)
6   Nn00 <- sum(dataLat$Clat == 0 & dataLat$Y == 0)
7   Nc <- sum(Clat == 1)
8   Nn <- sum(Clat == 0)
9   wc <- rbeta(1, 1 + Nc, 1 + Nn)
10  nc1 <- rbeta(1, 1 + Nc111, 1 + Nc110)
11  nn <- rbeta(1, 1 + Nn01, 1 + Nn00)
12  nc0 <- rbeta(1, 1 + Nc011, 1 + Nc010)
13  Clat <- sapply(1:N, function(i){SampleType(z = Z[i], d = D[i], y = Y[i], wc = wc, nc0 = nc0, nn0 = nn)})
14  PosteriorDrawsER[s, ] <- c(wc, nc1, nc0, nn)
15  setTxtProgressBar(pb, s)
16 }
17 close(pb)
18 keep <- seq((burnin+1), tot)
19 LATEER <- PosteriorDrawsER[keep, 2] - PosteriorDrawsER[keep, 3]
20 LATEERmean <- mean(LATEER)
21 LATEERci <- quantile(LATEER, c(0.025, 0.975))
22 # Plot posterior distribution of ATE
23 hist(LATEER, breaks = 40, freq = FALSE, main = "Posterior Distribution of LATE: Exclusion restrictions", xlab = "LATE", col = "lightblue", border = "white")
24 abline(v = LATEERmean, col = "red", lwd = 2)
25 abline(v = LATEERci, col = "darkgreen", lty = 2, lwd = 2)
26 legend("topright", legend = c("Posterior Mean", "95% Credible Interval"), col = c("red", "darkgreen"), lwd = 2, lty = c(1, 2), bty = "n", cex = 0.8)

```

**Example: Treatment effect of 401(k) participation on net financial assets**

In the example of the effect of *eligibility* on 401(k) on net financial assets, we calculate the ITT effect. Now, following [82], we use eligibility as an instrument for *participation* and perform inference on the (local) average treatment effect. We adopt the framework described in Section 7.3 for this example. The following code illustrates the procedure, and Figure 13.14 displays the posterior distribution of the (local) average treatment effect of participation. The 95% credible interval is (USD 5,102, USD 12,010), and the posterior mean is

### Posterior Distribution of LATE: Exclusion restrictions



**FIGURE 13.13**

Posterior distribution: Local average average atreatment effect vitamin A supplements.

USD 8,520, which is higher than the intention-to-treat effect associated with eligibility (USD 5,903).

Exercise 6 asks how to recover the ITT from the LATE and the effect of eligibility on participation. See [104] for practical methods to conduct sensitivity analysis of posterior results when relaxing the exclusion restriction in IV settings. In addition, [105] present a Bayesian inferential framework using a semi-parametric approach in which stochastic errors are modeled using a Dirichlet process (see Exercise 7), and [360] extend this approach to systems of equations.

**R code. Treatment effect: 401(k) participation on net financial assets**

```

1 rm(list = ls()); set.seed(10101)
2 library(coda); library(ggplot2)
3 mydata <- read.csv("https://raw.githubusercontent.com/
  BEsmarter-consultancy/BSTApp/reviews/master/DataApp/
  401k.csv", sep = ",", header = TRUE, quote = "")
```

**R code. Treatment effect: 401(k) participation on net financial assets**

```

1 # Attach variables
2 attach(mydata)
3 y <- net_tfa/1000 # Outcome: net financial assets
4 x <- as.vector(p401) # Endogenous regressor: participation
5 w <- as.matrix(cbind(1, age, inc, fsize, educ, marr, twoearn,
6 , db, pira, hown)) # Exogenous regressors with
7 intercept
8 z <- as.matrix(e401) # Instrument: eligibility (NO
9 intercept here)
10 X <- cbind(x, w); Z <- cbind(z, w)
11 # Dimensions
12 k <- ncol(X); kz <- ncol(Z)
13 # Priors
14 b0 <- rep(0, k); B0i <- diag(1e-5, k)
15 g0 <- rep(0, kz); G0i <- diag(1e-5, kz)
16 nu <- 3; Psi0 <- nu * 1000 * diag(2); Psi0i <- solve(Psi0)
17 # MCMC parameters
18 mcmc <- 5000; burnin <- 1000
19 tot <- mcmc + burnin; thin <- 1
20 # Auxiliary elements
21 Xtx <- t(X)%*%X; Ztz <- t(Z)%*%Z; nun <- nu + length(y)
22 # Gibbs sampling
23 PostBeta <- function(Sigma, Gamma){
24   w1 <- Sigma[1,1] - Sigma[1,2]^2/Sigma[2,2]
25   Bn <- solve(w1^(-1)*Xtx + B0i)
26   yaux <- y - (Sigma[1,2]/Sigma[2,2])*(x - Z%*%Gamma)
27   bn <- Bn%*%(B0i%*%b0 + w1^(-1)*t(X)%*%yaux)
28   Beta <- MASS::mvrnorm(1, bn, Bn)
29   return(Beta)
30 }
31 PostGamma <- function(Sigma, Beta){
32   w2 <- Sigma[2,2] - Sigma[1,2]^2/Sigma[1,1]
33   Gn <- solve(w2^(-1)*Ztz + G0i)
34   xaux <- x - (Sigma[1,2]/Sigma[1,1])*(y - X%*%Beta)
35   gn <- Gn%*%(G0i%*%g0 + w2^(-1)*t(Z)%*%xaux)
36   Gamma <- MASS::mvrnorm(1, gn, Gn)
37   return(Gamma)
38 }
39 PostSigma <- function(Beta, Gamma){
40   Uy <- y - X%*%Beta; Ux <- x - Z%*%Gamma
41   U <- cbind(Uy, Ux)
42   Psin <- solve(Psi0i + t(U)%*%U)
43   Sigmai <- rWishart::rWishart(1, df = nun, Sigma = Psin)
44   Sigma <- solve(Sigmai[,1])
45   return(Sigma)
46 }
```

**R code. Treatment effect: 401(k) participation on net financial assets**

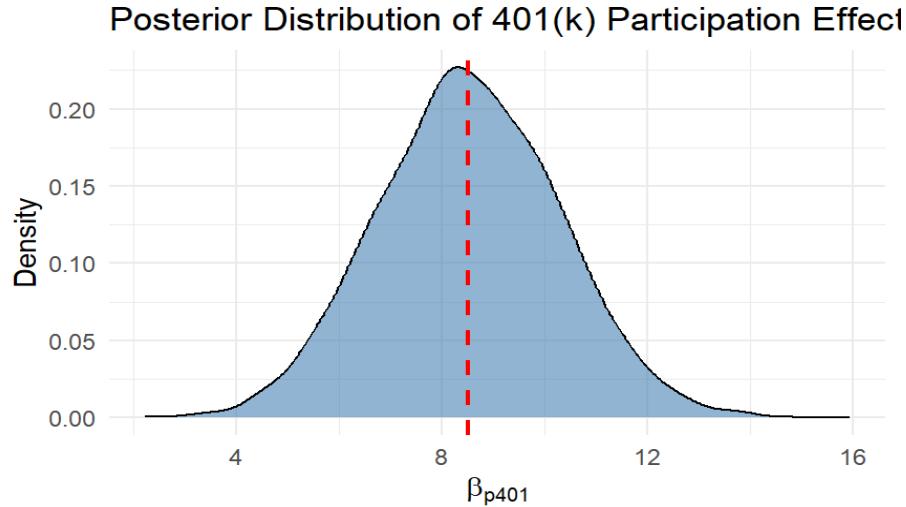
```

1 PostBetas <- matrix(0, tot, k)
2 PostGammas <- matrix(0, tot, kz)
3 PostSigmas <- matrix(0, tot, 2*(2+1)/2)
4 Beta <- rep(0, k); Gamma <- rep(0, kz)
5 pb <- txtProgressBar(min = 0, max = tot, style = 3)
6 for(s in 1:tot){
7   Sigma <- PostSigma(Beta = Beta, Gamma = Gamma)
8   Beta <- PostBeta(Sigma = Sigma, Gamma = Gamma)
9   Gamma <- PostGamma(Sigma = Sigma, Beta = Beta)
10  PostBetas[s,] <- Beta
11  PostGammas[s,] <- Gamma
12  PostSigmas[s,] <- matrixcalc::vech(Sigma)
13  setTxtProgressBar(pb, s)
14 }
15 close(pb)
16 keep <- seq((burnin+1), tot, thin)
17 Bs <- PostBetas[keep,]
18 Gs <- PostGammas[keep,]
19 Sigmas <- PostSigmas[keep,]
20 summary(coda::mcmc(Bs))
21 summary(coda::mcmc(Gs))
22 summary(coda::mcmc(Sigmas))
23 # Extract posterior draws for the treatment effect (
24   participation = p401)
25 beta_draws <- Bs[,1]
26 # Plot posterior distribution of treatment effect
27 df_beta <- data.frame(effect = as.vector(beta_draws))
28 ggplot(df_beta, aes(x = effect)) + geom_density(fill = "
  steelblue", alpha = 0.6) +
  geom_vline(xintercept = mean(beta_draws), color = "red",
  linetype = "dashed", linewidth = 1) + labs( title =
  "Posterior Distribution of 401(k) Participation Effect",
  x = expression(beta["p401"]), y = "Density") + theme_
  minimal(base_size = 14)

```

### 13.5 Difference-in-differences design (DiD)

In this section, we present how to perform inference and discuss the identification conditions when working with panel (longitudinal) data or repeated cross-sectional data in the *difference-in-differences* (DiD) framework. In DiD

**FIGURE 13.14**

Posterior distribution: Local average treatment effect 401k participation on net financial assets.

designs, we estimate causal effects by comparing changes over time between groups that receive the treatment and groups that do not (the “control” group). This approach is among the most widely used designs for identifying causal effects [23]. Its origins can be traced back to John Snow, the father of modern epidemiology, who in 1855 analyzed the spread of cholera through water in London [84].

In what follows, we first present the basic two-group, two-period ( $2 \times 2$ ) DiD setup, and then introduce the staggered DiD design.

### 13.5.1 Classic DiD: two-group and two-period ( $2 \times 2$ ) setup

Following [380], we consider a setting with two time periods,  $t = 1, 2$ , and two groups: a treated group ( $D_i = 1$ ) that receives the treatment between periods  $t = 1$  and  $t = 2$ , and a control group that never receives the treatment ( $D_i = 0$ ). We observe outcomes  $Y_{it}$  and treatment status  $D_i$  for units  $i = 1, 2, \dots, N$ , assuming a balanced panel data structure.

Let  $Y_{it}(1)$  denote the potential outcome for unit  $i$  in period  $t$  if it is untreated in the first period and treated in the second, while  $Y_{it}(0)$  denotes the potential outcome if it is never treated. In the DiD framework, the primary estimand of interest is the average treatment effect on the treated (ATT),

$$\tau_2 = \mathbb{E}[Y_{i2}(1) - Y_{i2}(0) | D_i = 1].$$

Note that we do not observe  $Y_{i2}(0) | D_i = 1$ ; that is, the outcome in the

second period for treated units if they had not received the treatment. Thus, the identification conditions in DiD are designed to express this counterfactual outcome as a function of the data.

There are two key identification assumptions:

1. *Parallel trends*:

$$\mathbb{E}[Y_{i2}(0) - Y_{i1}(0) \mid D_i = 1] = \mathbb{E}[Y_{i2}(0) - Y_{i1}(0) \mid D_i = 0],$$

which states that, in the absence of treatment, the average change in outcomes for the treated and control groups would have evolved similarly over time.

2. *No anticipation effects*:

$$\mathbb{E}[Y_{i1}(0) \mid D_i = 1] = \mathbb{E}[Y_{i1}(1) \mid D_i = 1],$$

meaning that the treatment has no causal effect prior to its implementation.

Thus, we can use the parallel trends assumption to express  $\mathbb{E}[Y_{i2}(0) \mid D_i = 1]$  as

$$\begin{aligned} \mathbb{E}[Y_{i2}(0) \mid D_i = 1] &= \mathbb{E}[Y_{i1}(0) \mid D_i = 1] + \mathbb{E}[Y_{i2}(0) - Y_{i1}(0) \mid D_i = 0] \\ &= \mathbb{E}[Y_{i1}(1) \mid D_i = 1] + \mathbb{E}[Y_{i2}(0) - Y_{i1}(0) \mid D_i = 0] \\ &= \mathbb{E}[Y_{i1} \mid D_i = 1] + \mathbb{E}[Y_{i2} - Y_{i1} \mid D_i = 0], \end{aligned}$$

where the second equality follows from the no anticipation effects assumption.

Therefore,

$$\tau_2 = \mathbb{E}[Y_{i2} - Y_{i1} \mid D_i = 1] - \mathbb{E}[Y_{i2} - Y_{i1} \mid D_i = 0], \quad (13.1)$$

that is, the ATT is the difference in outcome differences between treated and control units.

Figure 13.15 shows a graphical representation of the identification assumptions in DiD [84]. [326] show a DAG representation of DiD.

Note that we can express the observed outcome in terms of potential outcomes as

$$Y_{it} = Y_{it}(0) + [Y_{it}(1) - Y_{it}(0)]D_i.$$

Therefore,

$$Y_{i1} = Y_{i1}(0) + [Y_{i1}(1) - Y_{i1}(0)]D_i,$$

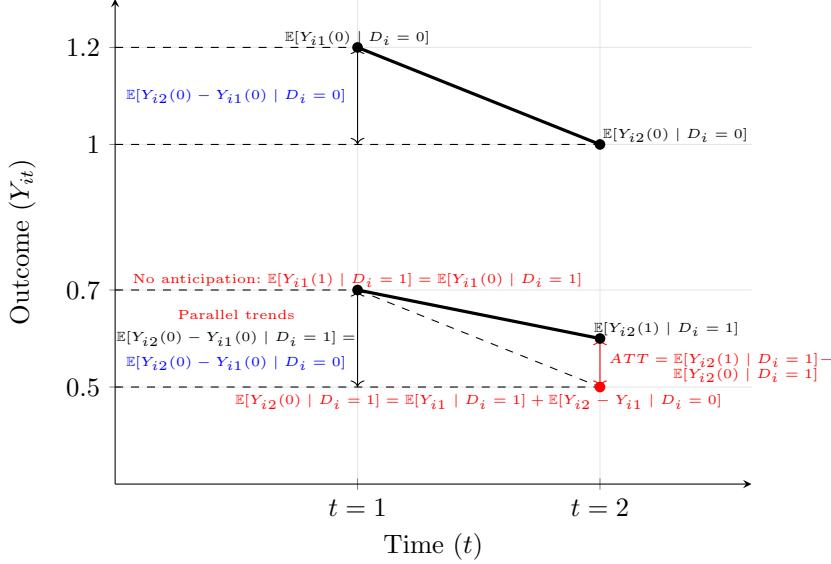
and

$$Y_{i2} = Y_{i2}(0) + [Y_{i2}(1) - Y_{i2}(0)]D_i.$$

Subtracting the first equation from the second yields

$$Y_{i2} - Y_{i1} = \underbrace{Y_{i2}(0) - Y_{i1}(0)}_{\text{common trend}} + \underbrace{[(Y_{i2}(1) - Y_{i1}(1)) - (Y_{i2}(0) - Y_{i1}(0))] D_i}_{\text{treatment effect in changes}}$$

Under the *parallel trends* assumption, the expected value of the common

**FIGURE 13.15**

Difference-in-Differences: Identification strategy

trend is the same for treated and control units. This assumption is compatible with the linear CEF function

$$Y_{it} = \alpha_i + \phi_t + \tau_2 D_i + \mu_{it},$$

which implies

$$Y_{i2} - Y_{i1} = (\phi_2 - \phi_1) + \tau_2 D_i + (\mu_{i2} - \mu_{i1}).$$

Comparing this expression with the one derived from the potential outcomes framework, it follows that regressing the time difference  $Y_{i2} - Y_{i1}$  on a constant and the treatment indicator  $D_i$  identifies  $\tau_2$ , the ATT.

Another common formulation in the linear regression setting that recovers the ATT is the two-way fixed effects (TWFE) model

$$Y_{it} = \alpha + \alpha_i + \phi_t + \tau_2 [D_i \cdot \mathbb{1}(t = 2)] + \epsilon_{it},$$

where  $\mathbb{1}(t = 2)$  is an indicator for the post-treatment period [380]. The advantage of the regression setting is that it allows for the straightforward calculation of the standard error of the ATT, enabling the construction of confidence or credible intervals under the Frequentist and Bayesian approaches, respectively.

The DiD framework can be extended to condition on covariates. In this case, we assume that the underlying identification assumptions are more plausible among units that are similar in terms of observed characteristics. Thus,

the identification assumptions are stated conditional on the exogenous variables.

Therefore, we can use covariates to assess the balance between control and treated groups in terms of levels ( $X$ ) or differences ( $\Delta X$ ) before and after treatment:

$$\frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\hat{\sigma}_1^2 + \hat{\sigma}_2^2}},$$

where  $\hat{\sigma}_l^2$  is the sample variance of  $X_l$ ,  $l = \{1, 2\}$ . As a rule of thumb, absolute standardized differences greater than 0.25 indicate potentially problematic imbalance [23]. Imbalance in covariates can lead to violations of the parallel trends assumption if these covariates would generate different outcomes under the counterfactual scenarios.

In this setting, the *conditional parallel trends* assumption becomes

$$\mathbb{E}[Y_{i2}(0) - Y_{i1}(0) | D_i = 1, \mathbf{X}_i] = \mathbb{E}[Y_{i2}(0) - Y_{i1}(0) | D_i = 0, \mathbf{X}_i], \text{ almost surely,}$$

and the *no anticipation* assumption becomes

$$\mathbb{E}[Y_{i1}(0) | D_i = 1, \mathbf{X}_i] = \mathbb{E}[Y_{i1}(1) | D_i = 1, \mathbf{X}_i].$$

In addition, the *overlap* assumption is required: there exists a treatment group whose characteristics overlap with those of the control group,

$$P(D_i = 1 | \mathbf{X}_i) < 1 - \epsilon, \epsilon > 0, \text{ almost surely, and, } P(D_i = 1) > 0.$$

Following similar arguments as before, we obtain the *conditional average treatment effect on the treated* (CATT) [23]:

$$\tau_2(\mathbf{X}_i) = \mathbb{E}[Y_{i2} - Y_{i1} | D_i = 1, \mathbf{X}_i] - \mathbb{E}[Y_{i2} - Y_{i1} | D_i = 0, \mathbf{X}_i].$$

The overall ATT can then be recovered by averaging over the distribution of  $\mathbf{X}_i$  in the treated population:

$$\begin{aligned} \tau_2 &= \mathbb{E}_{\mathbf{X}} \{ \mathbb{E}[Y_{i2} - Y_{i1} | D_i = 1, \mathbf{X}_i] - \mathbb{E}[Y_{i2} - Y_{i1} | D_i = 0, \mathbf{X}_i] | D_i = 1 \} \\ &= \mathbb{E}[Y_{i2} - Y_{i1} | D_i = 1] - \mathbb{E}_{\mathbf{X}} \{ \mathbb{E}[Y_{i2} - Y_{i1} | D_i = 0, \mathbf{X}_i] | D_i = 1 \}, \end{aligned}$$

where the second line uses the law of iterated expectations:

$$\mathbb{E}_{\mathbf{X}|D=1} (\mathbb{E}[Y_{i2} - Y_{i1} | D = 1, \mathbf{X}_i]) = \mathbb{E}[Y_{i2} - Y_{i1} | D = 1].$$

Assuming a linear CEF and a homogeneous ATT, that is, treatment effects do not vary across  $\mathbf{X}_{it}$ , we can use the first-difference specification,

$$\Delta Y_{i2} = \alpha + \beta D_i + \sum_{k=1}^K \beta_k \Delta X_{ik2} + \Delta \mu_{i2},$$

where  $\Delta Y_{i2} = Y_{i2} - Y_{i1}$  and  $\Delta X_{ik2} = X_{ik2} - X_{ik1}$ . This specification identifies

the ATT under the homogeneity assumption, in which case  $\tau_2 = \beta$ ; otherwise, it identifies a weighted average of the conditional ATTs across covariates, with the possibility of negative weights that do not sum to one (*non-convex weights*) [23].

An alternative specification is the TWFE

$$\begin{aligned} Y_{it} = & \alpha + \alpha_i + \phi_t + \beta(\mathbb{1}(t = 2) \cdot D_i) + \sum_{k=1}^K \beta_k(\mathbb{1}(t = 2) \cdot X_{ik1}) \\ & + \sum_{k=1}^K \delta_k(\mathbb{1}(t = 2) \cdot D_i \cdot X_{ik1}) + \epsilon_{it}. \end{aligned}$$

This specification can identify heterogeneous ATT under the assumption of a linear CEF.

However, alternative *semiparametric* and *nonparametric* approaches, such as *outcome regression adjustment*, *inverse probability weighting*, and *doubly robust* estimators (which combine the first two), are often preferable because they yield consistent estimates under conditional parallel trends without requiring a linear conditional expectation function (CEF) [2, 380]. These approaches typically involve two stages: first, estimating either the conditional mean of the outcome or the propensity score (probability of treatment; [376]), and second, using these estimates in a plug-in step to perform inference on the ATT. For Bayesian doubly robust frameworks for average treatment effects, see [394, 464, 282, 62] and Section 13.10.

A point to be aware of when performing DiD identification strategies is that the parallel trends assumption is generally not robust to functional form transformations of the outcome. That is, it may hold when the outcome is measured in levels but fail when it is expressed in logarithms. This is because two groups can have the same absolute increase (parallel trends in levels) but different percentage increases (non-parallel in logs).

Note that the parallel trends and no anticipation assumptions in DiD are not fundamentally testable, since they are identifying restrictions about counterfactual outcomes that are never observed. However, we can examine pre-treatment dynamics to assess their plausibility. For instance, we can check for pre-existing differences in trends (“pre-trends”) as a diagnostic for the parallel trends assumption.

Given a pre-treatment period  $t = 0$ , the no anticipation assumption implies that

$$\mathbb{E}[Y_{it}(0)] = \mathbb{E}[Y_{it}] \quad \text{for all } i \text{ and } t \leq 0.$$

Thus, one can test whether

$$\mathbb{E}[Y_{i1}(0) - Y_{i0}(0) \mid D_i = 1] = \mathbb{E}[Y_{i1}(0) - Y_{i0}(0) \mid D_i = 0],$$

which corresponds to checking for differences in pre-treatment trends between treated and control groups.

The results of such tests should be interpreted with caution. Even if pre-trends are perfectly parallel, this does not guarantee that the parallel trends assumption will hold in the post-treatment period. Moreover, there may be power issues: a true pre-existing difference in trends may go undetected if pre-treatment estimates are imprecise due to high variability.

In addition, one can conduct placebo tests by pretending that the policy started earlier and checking whether a spurious effect is detected prior to the true start. This provides indirect evidence on the plausibility of the no anticipation assumption.

Therefore, it is advisable to complement pre-trend diagnostics with sensitivity analyses, incorporating context-specific knowledge about plausible confounding factors.

### 13.5.2 Staggered (SDiD): G-group and T-period ( $G \times T$ ) setup

Many empirical settings are more complex than the two-period, two-group ( $2 \times 2$ ) design. Nevertheless, even in such settings we can view the design as an aggregation of  $2 \times 2$  comparisons between units that receive treatment and units that are not yet treated at the relevant time [23].

In a setting with multiple periods and different groups, assuming that treatment is an absorbing state—that is, once treated, units remain treated—staggered DiD (SDiD) analyses seek to estimate a weighted average of heterogeneous ATTs, where the weights may be based on economic/policy relevance.

Given  $t = 1, 2, \dots, T$ , a unit can be treated at  $t = g > 1$ ,

$$Y_{it}(g) := Y_{it}\left(\underbrace{0, \dots, 0}_{g-1}, \underbrace{1, \dots, 1}_{T-g+1}\right)$$

denotes the potential outcome path for a unit treated at  $g$ , and

$$Y_{it}(\infty) := Y_{it}\left(\underbrace{0, \dots, 0}_T\right)$$

denotes the potential outcome path for a never-treated unit. Under the absorbing treatment assumption,  $D_{it} = 1$  for all  $t \geq g$ . Let

$$G_i = \min \{t : D_{it} = 1\}$$

denote the first period in which unit  $i$  is treated.

The average treatment effect for units first treated in period  $g$  at time  $t$  is

$$ATT(g, t) = \mathbb{E}[Y_{it}(g) - Y_{it}(\infty) | G_i = g].$$

Therefore, each cohort  $g$  has its own  $T-1$  ATT estimands. Each  $ATT(g, t)$  represents the average treatment effect of initiating treatment in period  $g$ , relative to never receiving treatment.

The identifying assumptions for  $ATT(g, t)$  are:

1. *Staggered parallel trends based on not-yet-treated units:*

$$\mathbb{E}[Y_{it}(\infty) - Y_{it-1}(\infty) \mid G_i = g] = \mathbb{E}[Y_{it}(\infty) - Y_{it-1}(\infty) \mid G_i = g'], \quad \forall t \geq g, g' > t. \quad (13.2)$$

The *not-yet-treated* staggered parallel assumption establishes that the average change in potential outcomes for all treated groups would have evolved in parallel after treatment began, under the counterfactual scenario in which they had not received treatment.

Note that the *not-yet-treated* parallel trends assumption can be modified to use *never-treated* units as the comparison group (replace  $g'$  by  $\infty$  in Equation 13.2), or to impose parallel trends across *all* groups and periods (no restrictions on  $g$  and  $g'$ ). The *all-groups, all-periods* version is more demanding but can yield more precise estimates because it exploits a larger set of observations. In contrast, the *never-treated* version requires fewer assumptions about parallel trends but is typically less precise. The *not-yet-treated* version represents a middle ground between these two. Ultimately, the choice of which parallel trends assumption to adopt is context-specific [23].

2. *Staggered no anticipation:*

$$\mathbb{E}[Y_{it}(g)] = \mathbb{E}[Y_{it}(\infty)], \quad \forall i, t < g.$$

That is, treated units do not change their potential outcomes in anticipation of future treatment before treatment begins.

Thus, to identify long-term effects, the parallel trends assumption must hold for all periods after the earliest treatment period in the not-yet-treated case. If the goal is to estimate effects in a specific short-term period, parallel trends must hold in that period.

Note that under these assumptions, we can express  $ATT(g, t)$  in terms of observable outcomes as [380]

$$ATT(g, t) = \mathbb{E}[Y_{it} - Y_{i,g-1} \mid G_i = g] - \mathbb{E}[Y_{it} - Y_{i,g-1} \mid G_i = g'], \quad g' > t.$$

This formulation makes explicit the choice of the comparison group, either not-yet-treated or never-treated units in this setting. Moreover, cohorts are indexed by the first treatment date  $g$ , and effects are evaluated at calendar time  $t$ . Hence the object of interest is cohort-time specific,  $ATT(g, t)$ , and staggered designs generally imply a multiplicity of ATT parameters, where each estimand can be seen as a multi-period version of a DiD causal effect (Equation 13.1).

To study the dynamics of treatment effects, we can use *event study* plots. These involve estimating and reporting effects for a range of periods before and after treatment, allowing us to analyze the temporal pattern of the ATTs [23].

An advantage of having multiple pre-treatment periods in the design is that it allows calculating *falsification/placebo tests*. Under the *no anticipation*

assumption, any potential effect before  $G_i$  must be zero. That is, given  $b < \min\{g, g'\}$ , for all  $t$  with  $b < t < \min\{g, g'\}$ :

$$\mathbb{E}[Y_{it} - Y_{ib} | G_i = g] = \mathbb{E}[Y_{it} - Y_{ib} | G_i = g'].$$

Thus, we can test whether the differences in the expected changes of the outcome variable before treatment between different groups are statistically indistinguishable from zero. If this is the case, common practice would suggest that there is evidence in favor of parallel trends. However, we should remember that the parallel trends assumption is essentially untestable because it imposes restrictions on post-treatment periods, not on pre-treatment periods. Consequently, the same recommendations as those at the end of the previous section apply: perform sensitivity analyses regarding potential violations of the parallel trends assumption.

Assuming that parallel trends hold conditional on covariates, we obtain identification under the following assumptions [23]:

1. *Staggered parallel trends with not-yet-treated comparison (conditional on  $\mathbf{X}_i$ )*:

$$\begin{aligned} \mathbb{E}[Y_{it}(\infty) - Y_{i,t-1}(\infty) | G_i = g, \mathbf{X}_i] &= \mathbb{E}[Y_{it}(\infty) - Y_{i,t-1}(\infty) | G_i = g', \mathbf{X}_i], \\ t &\geq g, g' > t. \end{aligned}$$

2. *Staggered overlap (positivity) for not-yet-treated comparisons*:

$$P(G_i = g | \mathbf{X}_i) < 1 - \epsilon, \epsilon > 0, \text{ almost surely, and, } P(G_i = g) > 0.$$

Condition 1 means that conditional on covariates, pre-treatment changes in untreated potential outcomes for cohort  $g$  match those of units not yet treated at time  $t$ , and condition 2 means that there is positive probability of belonging to cohort  $g$  and of being not yet treated at  $t$ .

Given these assumptions and no anticipation, the  $ATT(g, t)$  can be expressed as

$$ATT(g, t) = \mathbb{E}[Y_{it} - Y_{ig-1} | G_i = g] - \mathbb{E}_{\mathbf{X}} \{\mathbb{E}[Y_{it} - Y_{ig-1} | \mathbf{X}_i, G_i = g'] | G_i = g\}.$$

Again, inference on  $ATT(g, t)$  can be conducted using regression adjustment, inverse-probability weighting, or doubly robust estimators. Although TWFE specifications are often used to estimate  $ATT(g, t)$ , [23] recommend avoiding them in staggered DiD settings because they generally do not identify  $ATT(g, t)$ . Instead, they recover a complicated, non-convex weighted average of cohort-time effects, with potentially negative weights.

#### **Example: Difference-in-Differences simulation**

Let's simulate a DiD setting where the treatment effect is 1 and the common trend is  $-1.8$ . We assume the default priors in the *MCMCpack* package. We perform inference on the ATT using

$$Y_{i2} - Y_{i1} = \underbrace{\phi_2 - \phi_1}_{\text{Common change}} + \tau_2 D_i + \epsilon_i.$$

The following code implements this example. From the posterior results, we can verify that the 95% credible interval contains the true ATT. We ask in Exercise 8 to repeat exactly the same simulation, and estimate the ATT using the specification with the interaction between treatment and the post-treatment period.

### R code. Difference-in-Differences: Simulation exercise

```

1 rm(list = ls()); set.seed(10101)
2 library(ggplot2); library(dplyr); library(fastDummies)
3 # Parameters
4 N_per_group <- 200           # units per group
5 T_periods     <- 2           # keep 2x2 for clarity
6 tau_true      <- 1           # ATT
7 sigma_eps     <- 0.5         # noise SD
8 # Panel index
9 id   <- rep(1:(2*N_per_group), each = T_periods)
10 t    <- rep(1:T_periods, times = 2*N_per_group)
11 # Group: treated (D=1) vs control (D=0)
12 D    <- rep(c(rep(0, N_per_group), rep(1, N_per_group)), each
13       = T_periods)
14 # Post indicator (t=2 is post)
15 post <- as.integer(t == 2)
16 # Unit fixed effects (random heterogeneity)
17 alpha_i <- rnorm(2*N_per_group, 0, 0.8)
18 alpha   <- alpha_i[id]
19 # Time effects (common shocks)
20 phi_t <- c(0, -1.8) # common decline from t=1 to t=2
21 phi    <- phi_t[t]
22 treat_effect <- tau_true * (D * post)
23 # Outcome:
24 eps <- rnorm(length(id), 0, sigma_eps)
25 Y    <- alpha + phi + treat_effect + eps
26 did <- data.frame(id, t, D, post, Y)
27 # Bayesian inference: Model in differences
28 dY <- did$Y[did$t==2] - did$Y[did$t==1]
29 D   <- did$D[did$t==1]
30 post_fit <- MCMCpack::MCMCregress(dY ~ 1 + D, burnin = 100,
31                                     mcmc = 1000)
32 tau_draws <- post_fit[, "D"]
33 quantile(tau_draws, c(.025,.5,.975))
34 quantile(post_fit[,1], c(.025,.5,.975))

```

## 13.6 Regression discontinuity design (RD)

In this section, we study the *regression discontinuity* (RD) design. In particular, we analyze the causal effect of a binary treatment  $D_i$  on an outcome  $Y_i$  within the potential outcomes (Rubin causal model) framework, namely  $Y_i(1) - Y_i(0)$ . As usual, we face the fundamental problem of causal inference. Identification in RD relies on a *running variable*  $Z_i$ , also known as the *forcing, index, or score variable*, which determines treatment assignment, either completely or partially, depending on whether a unit lies on either side of a known cutoff  $c$ . The running variable may be associated with the potential outcomes, but this association is assumed to be smooth at  $c$ . Under this smoothness assumption (and in the absence of precise manipulation at  $c$ ), any discontinuity in the conditional mean of  $Y_i$  as a function of  $Z_i$  at  $c$  identifies the causal effect of the treatment at the cutoff [216].

There are two common versions of RD: *sharp* and *fuzzy*. In the former, the change (jump) in the conditional probability of treatment at the cutoff equals one; in the latter, it is strictly between zero and one:

$$\Delta = \lim_{z \downarrow c} \Pr(D_i = 1 \mid Z_i = z) - \lim_{z \uparrow c} \Pr(D_i = 1 \mid Z_i = z) = \begin{cases} 1, & \text{sharp RD,} \\ \in (0, 1), & \text{fuzzy RD.} \end{cases}$$

Here  $\lim_{z \downarrow c}$  and  $\lim_{z \uparrow c}$  denote the right (approach from above) and left (approach from below) limits, respectively.

### 13.6.1 Sharp regression discontinuity design (SRD)

In the sharp case, treatment assignment is a deterministic function of the running variable:

$$D_i = \mathbb{1}\{Z_i \geq c\},$$

so the conditional probability of treatment jumps from 0 to 1 at  $Z_i = c$ . All units with  $Z_i \geq c$  are treated, whereas those with  $Z_i < c$  serve as the control group. Therefore, assignment is identical to treatment in the sharp RD (perfect compliance). Figure 13.16 displays the conditional probability of treatment.

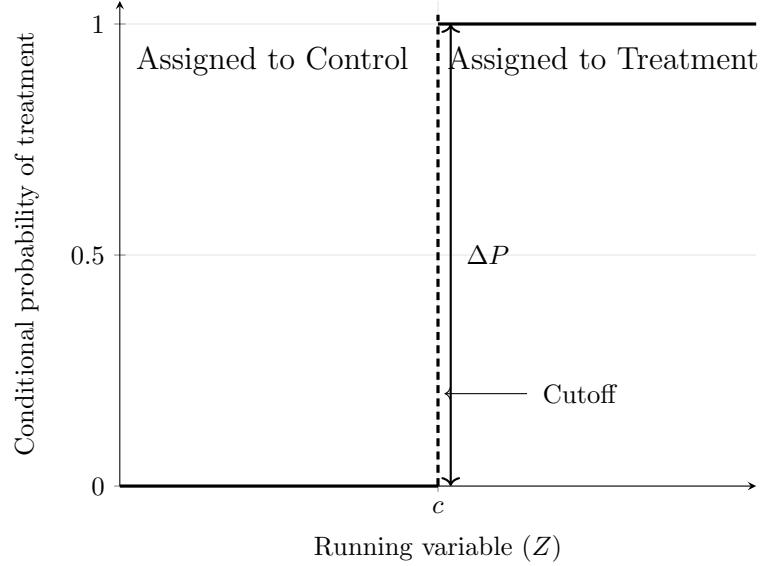
The observed outcome is given by

$$Y_i = (1 - D_i) \cdot Y_i(0) + D_i \cdot Y_i(1) = \begin{cases} Y_i(0), & Z_i < c \\ Y_i(1), & Z_i \geq c. \end{cases}$$

Thus, the causal effect that we target in SRD is

$$\tau_{SRD} = \mathbb{E}[Y_i(1) - Y_i(0) \mid Z_i = c].$$

Note that the treatment effect is evaluated at the cutoff, which is a single

**FIGURE 13.16**

Sharp regression discontinuity design: Conditional probability of treatment

point in the support of the running variable. Hence  $\tau_{SRD}$  is inherently local and need not be representative of units whose running variable lies far from  $c$ . The identification of  $\tau_{SRD}$  requires *continuity of the potential outcome regression functions*, that is,

$$\mathbb{E}[Y_i(0) | Z_i = z] \text{ and } \mathbb{E}[Y_i(1) | Z_i = z]$$

are continuous at  $z$ . Strictly speaking, it suffices that the left- and right-hand limits of  $\mathbb{E}[Y_i(d) | Z_i = z]$  coincide at  $z = c$  for  $d \in \{0, 1\}$ .

This implies

$$\begin{aligned} \mathbb{E}[Y_i(0) | Z_i = c] &= \lim_{z \uparrow c} \mathbb{E}[Y_i(0) | Z_i = z] \\ &= \lim_{z \uparrow c} \mathbb{E}[Y_i(0) | Z_i = z, D_i = 0] \\ &= \lim_{z \uparrow c} \mathbb{E}[Y_i | Z_i = z], \end{aligned}$$

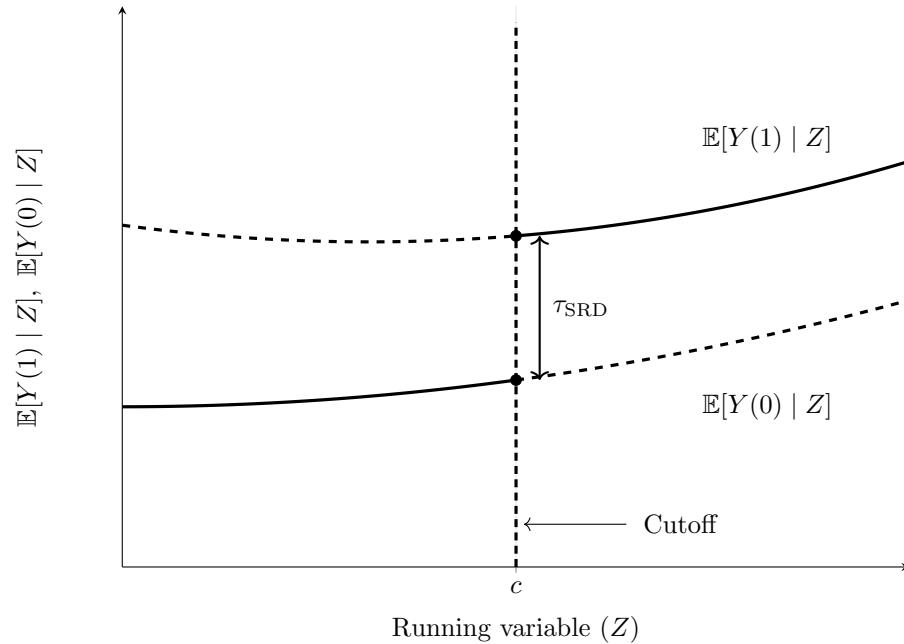
since for  $z < c$  we have  $D_i = 0$ . Similarly,

$$\begin{aligned} \mathbb{E}[Y_i(1) | Z_i = c] &= \lim_{z \downarrow c} \mathbb{E}[Y_i(1) | Z_i = z] \\ &= \lim_{z \downarrow c} \mathbb{E}[Y_i | Z_i = z], \end{aligned}$$

because for  $z > c$  we have  $D_i = 1$ . Therefore,

$$\tau_{SRD} = \lim_{z \downarrow c} \mathbb{E}[Y_i | Z_i = z] - \lim_{z \uparrow c} \mathbb{E}[Y_i | Z_i = z],$$

i.e., the target estimand is the jump in the conditional mean at the discontinuity point  $c$ . See Figure 13.17.



**FIGURE 13.17**

RD treatment effect in a sharp RD design. Solid segments are observed on each side of  $c$ ; dashed segments are counterfactual [74].

### 13.6.2 Fuzzy regression discontinuity design (FRD)

In the fuzzy case, the conditional probability of treatment has a discontinuous jump at the cutoff that is strictly between 0 and 1:

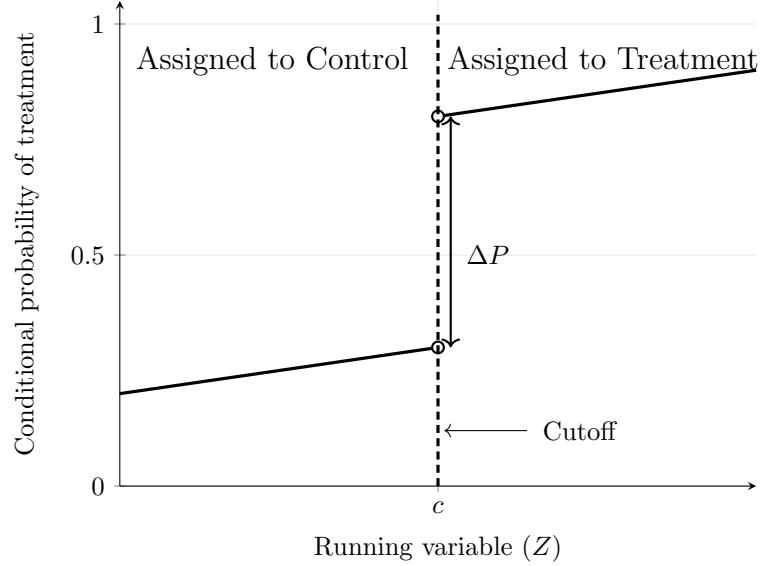
$$\lim_{z \uparrow c} \Pr(D_i = 1 | Z_i = z) \neq \lim_{z \downarrow c} \Pr(D_i = 1 | Z_i = z),$$

$$0 < \lim_{z \downarrow c} \Pr(D_i = 1 | Z_i = z) - \lim_{z \uparrow c} \Pr(D_i = 1 | Z_i = z) < 1.$$

Thus, the probability of treatment increases discontinuously at  $c$ . Note that in a fuzzy RD, compliance is imperfect: some units assigned to treatment do not take it, and some units not assigned may nonetheless receive it. Figure 13.18 displays the conditional probability of treatment.

Note that the expected value of the observed outcome given the running variable is given by

$$\mathbb{E}[Y_i | Z_i = z] = \mathbb{E}[Y_i(0) | D_i = 0, Z_i = z]P(D_i = 0 | Z_i = z) + \mathbb{E}[Y_i(1) | D_i = 1, Z_i = z]P(D_i = 1 | Z_i = z)$$

**FIGURE 13.18**

Fuzzy regression discontinuity design: Conditional probability of treatment

$$+ \mathbb{E}[Y_i(1) | D_i = 1, Z_i = z] P(D_i = 1 | Z_i = z).$$

Thus, the observed outcome is a weighted average of potential outcomes. Figure 13.19 displays the treatment effect that we are targeting.

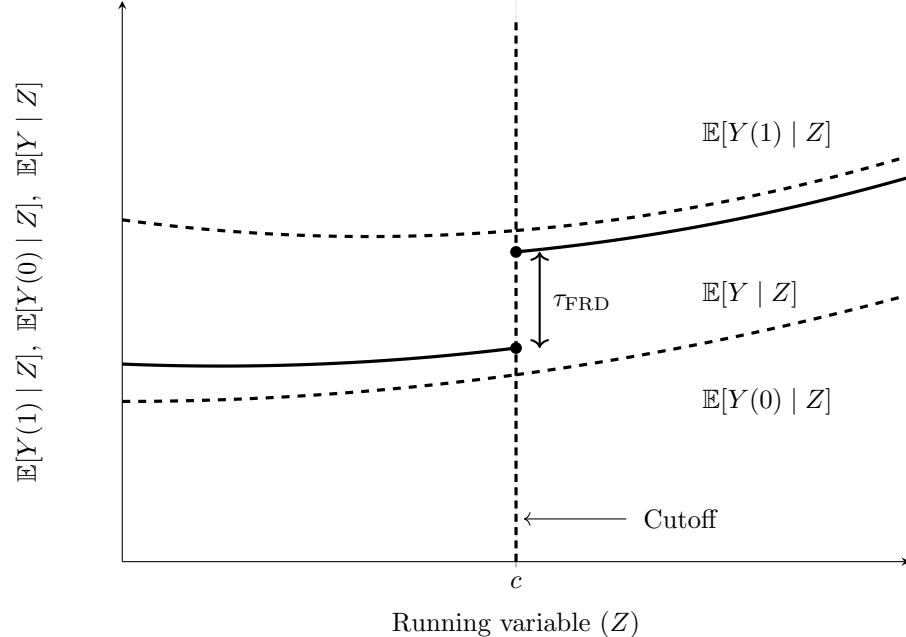
The target causal effect in the FRD is (see [188] for a proof)

$$\begin{aligned} \tau_{\text{FRD}} &= \frac{\lim_{z \downarrow c} \mathbb{E}[Y_i | Z_i = z] - \lim_{z \uparrow c} \mathbb{E}[Y_i | Z_i = z]}{\lim_{z \downarrow c} \mathbb{E}[D_i | Z_i = z] - \lim_{z \uparrow c} \mathbb{E}[D_i | Z_i = z]} \\ &= \mathbb{E}[Y_i(1) - Y_i(0) | Z_i = c], \text{ units are compliers.} \end{aligned}$$

Note that  $\tau_{\text{FRD}}$  encompasses  $\tau_{\text{SRD}}$  as  $\lim_{z \downarrow c} \mathbb{E}[D_i | Z_i = z] - \lim_{z \uparrow c} \mathbb{E}[D_i | Z_i = z] = 1$  in the SRD. In addition,  $\tau_{\text{FRD}}$  has the same structure as the local average treatment effect (LATE) in the IV literature, with the instrument given by the eligibility indicator  $T_i = \mathbb{1}\{Z_i \geq c\}$ . Because compliance is imperfect (some eligible units do not take treatment and some ineligible units receive it),  $D_i \neq T_i$  for some units, and the estimand pertains to *compliers at the cutoff*.

Let  $D_i(z)$  denote the potential treatment status if the cutoff were  $z$  (for  $z$  in a neighborhood of  $c$ ) such that  $D_i(z) = 1$  implies that the unit would be treated if the cutoff is  $z$ . The identification assumptions are:

1. *Continuity of the potential outcome regression functions*, particularly at  $c$ .
2. *A nontrivial first stage*,  $\lim_{z \downarrow c} \mathbb{E}[D_i | Z_i = z] - \lim_{z \uparrow c} \mathbb{E}[D_i | Z_i = z] \neq 0$ ,

**FIGURE 13.19**

RD treatment effect in a fuzzy RD design. Solid segments are observed on each side of  $c$ ; dashed segments are counterfactual.

and

3. *Monotonicity*:  $D_i(z)$  is non-increasing in  $z$  at  $z = c$ .

We can classify unit types (locally around the cutoff  $c$ ) using the potential treatment under a hypothetical cutoff  $z$ ,  $D_i(z)$  as follows:

*Complier status*

$$\lim_{z \downarrow Z_i} D_i(z) = 0, \text{ and } \lim_{z \uparrow Z_i} D_i(z) = 1,$$

*Never-takers*

$$\lim_{z \downarrow Z_i} D_i(z) = 0, \text{ and } \lim_{z \uparrow Z_i} D_i(z) = 0,$$

and

*Always-takers*

$$\lim_{z \downarrow Z_i} D_i(z) = 1, \text{ and } \lim_{z \uparrow Z_i} D_i(z) = 1.$$

Thus, compliers are units that would be treated if the cutoff were at or below their score  $Z_i$ , and would not be treated if the cutoff were above  $Z_i$ . Never-takers and always-takers do not change treatment status regardless of

the cutoff: the former are always untreated, whereas the latter are always treated. Monotonicity precludes *defiers*.

See [216] and [74] for discussions of graphical and formal diagnostics of the assumptions in RD.

**Example: Returns to compulsory schooling [95]**

[95] propose a Bayesian inferential framework that explicitly models unit type (complier, always-taker, and never-taker), which enables implementation of a Gibbs sampler. Their empirical application is motivated by the April 1947 reform in the United Kingdom that raised the minimum school-leaving age from 14 to 15.

In the example below, we construct a simulation calibrated to their application to illustrate Bayesian inference in a fuzzy RD. In particular, we specify priors, show the full conditional distributions, and implement a Gibbs sampler to estimate the treatment effect at the cutoff.

There are three types of individuals: compliers ( $c$ ), never-takers ( $n$ ), and always-takers ( $a$ ). These imply four potential outcomes  $Y_{ic}(0)$ ,  $Y_{ic}(1)$ ,  $Y_{in}(0)$ , and  $Y_{ia}(1)$ . [95] assume

$$Y_{ic}(0) = \beta_{0c0} + \beta_{0c1}Z_i + \mathbf{X}_i^\top \boldsymbol{\beta}_c + \mu_{0ci} = \mathbf{W}_{ic}^\top \boldsymbol{\beta}_{0c} + \mu_{0ci},$$

$$Y_{ic}(1) = \beta_{1c0} + \beta_{1c1}Z_i + \mathbf{X}_i^\top \boldsymbol{\beta}_c + \mu_{1ci} = \mathbf{W}_{ic}^\top \boldsymbol{\beta}_{1c} + \mu_{1ci},$$

$$Y_{in}(0) = \beta_{n0} + \mathbf{X}_i^\top \boldsymbol{\beta}_n + \mu_{ni} = \mathbf{W}_{in}^\top \boldsymbol{\beta}_{0n} + \mu_{ni},$$

and

$$Y_{ia}(1) = \beta_{a0} + \mathbf{X}_i^\top \boldsymbol{\beta}_a + \mu_{ai} = \mathbf{W}_{ia}^\top \boldsymbol{\beta}_{1a} + \mu_{ai},$$

where  $\mu_{dci} \sim t_v(0, \sigma_{dc}^2)$  for  $d \in \{0, 1\}$ ,  $\mu_{ni} \sim t_v(0, \sigma_{0n}^2)$ , and  $\mu_{ai} \sim t_v(0, \sigma_{1a}^2)$ ; here  $t_v$  denotes the Student- $t$  distribution with  $v$  degrees of freedom. The running variable is centered at the cutoff, so  $Z_i = 0$  at the threshold, and the FRD treatment effect is

$$\tau_{\text{FRD}} = \beta_{1c0} - \beta_{0c0}.$$

There are four cells in this setting, defined by policy epochs  $\mathbb{1}(Z_i < 0)$  and  $\mathbb{1}(Z_i \geq 0)$ , and treatment values  $D_i = 0$  and  $D_i = 1$ . Table 13.5 shows the case where, in the cell without treatment under the old policy ( $K_{00}$ ), there are compliers and never-takers, while in the cell with treatment under the new policy ( $K_{11}$ ), there are compliers and always-takers.

**TABLE 13.5**

Distribution of subject type by observed policy regime and treatment

Policy indicator	Treatment	
	$D_i = 0$	$D_i = 1$
$\mathbb{1}(Z_i < 0)$ (old policy)	c, n	a
$\mathbb{1}(Z_i \geq 0)$ (new policy)	n	c, a

The four cells by policy and treatment status are:

$$K_{00} = \{c, n\}, \quad K_{10} = \{n\}, \quad K_{01} = \{a\}, \quad K_{11} = \{c, a\},$$

and the two cells by treatment status are:

$$K_0 = \{c, n\}, \quad K_1 = \{c, a\}.$$

Under a Bayesian framework, posterior simulation simplifies because we sample each unit's latent type  $C_i \in \{c, n, a\}$ . This contrasts with the Frequentist literature, where types are not modeled explicitly. Thus, the Bayesian approach requires setting assumptions about the types (see [95] for details); for instance, types may be distributed smoothly around the cutoff with unknown distributions. In addition, given monotonicity, there are models for four potential outcomes. This differs from the Frequentist literature, which considers only two potential outcomes, given the binary treatment, and a model for the treatment.

Note that if we observe  $Z_i < 0$  and  $D_i = 0$ , only compliers and never-takers are possible. Then,

$$P(C_i = c \mid Y_i, D_i = 0, \boldsymbol{\theta}) = \frac{\eta_c t_v(Y_i \mid \mathbf{W}_{ic}^\top \boldsymbol{\beta}_{0c}, \sigma_{0c}^2)}{\eta_c t_v(Y_i \mid \mathbf{W}_{ic}^\top \boldsymbol{\beta}_{0c}, \sigma_{0c}^2) + \eta_n t_v(Y_i \mid \mathbf{W}_{in}^\top \boldsymbol{\beta}_{0n}, \sigma_{0n}^2)}, \\ i \in K_{00}$$

where  $\boldsymbol{\theta}$  denotes the set of model parameters, and  $\eta_c$  and  $\eta_n$  are the current mixture weights for a complier and a never-taker, respectively.

Conversely, if  $Z_i \geq 0$  and  $D_i = 1$ , only compliers and always-takers are possible:

$$P(C_i = c \mid Y_i, D_i = 1, \boldsymbol{\theta}) = \frac{\eta_c t_v(Y_i \mid \mathbf{W}_{ic}^\top \boldsymbol{\beta}_{1c}, \sigma_{1c}^2)}{\eta_c t_v(Y_i \mid \mathbf{W}_{ic}^\top \boldsymbol{\beta}_{1c}, \sigma_{1c}^2) + \eta_a t_v(Y_i \mid \mathbf{W}_{ia}^\top \boldsymbol{\beta}_{1a}, \sigma_{1a}^2)}, \\ i \in K_{11}$$

where  $\eta_a$  is the mixture weight for an always-taker.

Following [95], we assume a Dirichlet prior for the mixing probabilities  $(\eta_c, \eta_n, \eta_a)$  with hyperparameters  $\alpha_{0k}$ ,  $k \in \{c, n, a\}$ . The posterior distribution is then Dirichlet with parameters  $\alpha_{0k} + \sum_{i=1}^N \mathbb{1}(C_i = k)$ .

In addition, [95] use the normal-gamma mixture representation of the Student- $t$  distribution, where the mixing parameter satisfies  $\lambda_i \sim G(v/2, v/2)$ . Consequently, the posterior distribution of  $\lambda_i$  is gamma with parameters  $(v+1)/2$  and  $(v + \sigma_{dk}^{-2}(y_i - \mathbf{W}_{ik}^\top \boldsymbol{\beta}_{dk})^2)/2$ .

The prior distribution for the location parameters  $\boldsymbol{\beta}_{dk}$  is independent normal with mean  $\boldsymbol{\beta}_{dk,0}$  and variance matrix  $\mathbf{B}_{dk,0}$ . The posterior distribution of  $\boldsymbol{\beta}_{dk}$  is normal with mean

$$\boldsymbol{\beta}_{dk,n} = \mathbf{B}_{dk,n} \left( \mathbf{B}_{dk,0}^{-1} \boldsymbol{\beta}_{dk,0} + \sigma_{dk}^{-2} \sum_{i \in I_{dk}} \lambda_i \mathbf{W}_{ik} Y_i \right),$$

and covariance matrix

$$\mathbf{B}_{dk,n} = \left( \mathbf{B}_{dk,0}^{-1} + \sigma_{dk}^{-2} \sum_{i \in I_{dk}} \lambda_i \mathbf{W}_{ik} \mathbf{W}_{ik}^\top \right)^{-1},$$

where  $I_{dk} = \{i : D_i = d, C_i = k\}$  is the set of indices for observations by treatment status and unit type. Note that there are only four possible combinations determined by policy and treatment status, namely,  $I_{0c}$ ,  $I_{1c}$ ,  $I_{0n}$  and  $I_{1a}$ .

Finally, assuming that the prior distributions of the variances are independent inverse-gamma,  $\sigma_{dk}^2 \sim IG(\alpha_{dk,0}/2, \delta_{dk,0}/2)$ , the posterior distributions are also inverse-gamma,  $\sigma_{dk}^2 \sim IG(\alpha_{dk,n}/2, \delta_{dk,n}/2)$ , where

$$\alpha_{dk,n} = \alpha_{dk,0} + N_{dk},$$

and

$$\delta_{dk,n} = \delta_{dk,0} + \sum_{i \in I_{dk}} \lambda_i (Y_i - \mathbf{w}_{ik}^\top \boldsymbol{\beta}_{dk})^2.$$

Here,  $N_{dk}$  denotes the number of elements (cardinality) of  $I_{dk}$ .

The following code performs a simulation assuming that  $Z_i$  follows a discrete uniform distribution from -24 to 24, where the policy indicator equals one if  $Z_i \geq 0$ . Next,  $X_i$  follows a discrete uniform distribution between 85 and 95. We set  $\boldsymbol{\beta}_{0c} = [4.50 \ -0.20 \ 0.03]^\top$  and  $\boldsymbol{\beta}_{1c} = [5.50 \ 0.40 \ 0.03]^\top$ , so that the treatment effect is 1. In addition,  $\boldsymbol{\beta}_{0n} = [6.80 \ -0.02]^\top$  and  $\boldsymbol{\beta}_{1a} = [5.50 \ -0.04]^\top$ , with  $\sigma_{0c}^2 = \sigma_{1c}^2 = 0.10$ ,  $\sigma_{0n}^2 = 0.15$ , and  $\sigma_{1a}^2 = 0.20$ . The type probabilities are  $\eta_c = 0.70$ ,  $\eta_n = 0.15$ , and  $\eta_a = 0.15$ . We set the sample size to 3,000 and  $v = 5$ . In addition, we set non-informative priors, a burn-in of 1,000, and 5,000 MCMC iterations. The following code illustrates the implementation, and Figure 13.20 presents the posterior distribution of the LATE. The 95% credible interval encompasses the population value, and the posterior mean lies close to this value.

### R code. Simulation: Fuzzy regression discontinuity design

```

1 rm(list = ls()); set.seed(10101); library(ggplot2)
2 # Simulation setup
3 N <- 3000
4 Zi <- sample(seq(-24, 24, 1), N, replace = TRUE) # Running
5 Z <- Zi - mean(Zi)
6 T <- as.integer(Z >= 0) # Policy (0/1)
7 X <- sample(seq(85, 95, 1), N, replace = TRUE) # Regressor
8 Wc <- cbind(1, Z, X) # Compliers
9 W <- cbind(1, X) # No compliers

```

**R code. Simulation: Fuzzy regression discontinuity design**

```

1 B0c <- c(4.5, -0.2, 0.03); B1c <- c(5.5, 0.4, 0.03)
2 B0n <- c(6.8, -0.02); B1a <- c(5.5, -0.04)
3 s2c <- 0.1; s2n <- 0.15; s2a <- 0.2
4 Etac <- 0.7; Etan <- 0.15; Etaa <- 0.15
5 v <- 5
6 # Potential outcomes
7 mu0c <- sqrt(s2c) * rt(N, v); Y0c <- as.numeric(Wc %*% B0c
+ mu0c)
8 mu1c <- sqrt(s2c) * rt(N, v); Y1c <- as.numeric(Wc %*% B1c
+ mu1c)
9 mu0n <- sqrt(s2n) * rt(N, v); Y0n <- as.numeric(W %*% B0n
+ mu0n)
10 mu1a <- sqrt(s2a) * rt(N, v); Y1a <- as.numeric(W %*% B1a
+ mu1a)
11 # Latent types: first n, then a, remaining are c
12 id <- seq_len(N)
13 n_n <- as.integer(round(Etan * N))
14 n_a <- as.integer(round(Etaa * N))
15 id0n <- sample(id, n_n, replace = FALSE) # never-takers
16 id_rem <- setdiff(id, id0n)
17 id1a <- sample(id_rem, n_a, replace = FALSE) # always-
    takers
18 idc <- setdiff(id, c(id0n, id1a)) # compliers
19 type <- rep(NA_character_, N)
20 type[id0n] <- "n"
21 type[id1a] <- "a"
22 type[idc] <- "c"
23 # Realized treatment under monotonicity:
24 # n -> D=0, a -> D=1, c -> D=T
25 D <- integer(N)
26 D[type == "n"] <- 0; D[type == "a"] <- 1
27 D[type == "c"] <- T[type == "c"]
28 # 2x2 table of D vs T
29 tab_DT <- table(T = factor(T, levels = 0:1), D = factor(D,
    levels = 0:1))
30 print(tab_DT)
31 Y <- ifelse(type == "n", Y0n, ifelse(type == "a", Y1a,
    ifelse(D == 0, Y0c, Y1c)))
32 ## Gibbs sampler ##
33 I00 <- which(T == 0 & D == 0); I11 <- which(T == 1 & D == 1)
34 I01 <- which(T == 0 & D == 1); I10 <- which(T == 1 & D == 0)
35 # Hyperparameters
36 a0c <- 1; a0n <- 1; a0a <- 1 # Hyperpar Dirichlet
37 v <- 5 # Student-t
38 a0 <- 0.01; d0 <- 0.01 # Inv-Gamma
39 # Prob complier | D=0
40 dens_t_locscale <- function(y, mu, sig2, v) {
41   z <- (y - mu) / sqrt(sig2); dt(z, df = v) / sqrt(sig2)
42 }
```

**R code. Simulation: Fuzzy regression discontinuity design**

```

1 Pc0 <- function(b0c, sigma20c, etac, b0n, sigma20n, etan, i)
2   {
3     p0c <- etac * dens_t_loyscale(Y[i], mu = sum(Wc[i,]*b0c),
4       sig2 = sigma20c, v = v)
5     p0n <- etan * dens_t_loyscale(Y[i], mu = sum(W[i, ]*b0n),
6       sig2 = sigma20n, v = v)
7     pc0 <- p0c / (p0c + p0n + 1e-12)
8     return(pc0)
9   }
10 Pc1 <- function(b1c, sigma21c, etac, b1a, sigma21a, etaa, i)
11   {
12     p1c <- etac * dens_t_loyscale(Y[i], mu = sum(Wc[i,]*b1c),
13       sig2 = sigma21c, v = v)
14     p1a <- etaa * dens_t_loyscale(Y[i], mu = sum(W[i, ]*b1a),
15       sig2 = sigma21a, v = v)
16     pc1 <- p1c / (p1c + p1a + 1e-12)
17     return(pc1)
18   }
19 rtype <- function(alphas){
20   types <- MCMCpack::rdirichlet(1, alphas)
21   return(types)
22 }
23 PostLambda_i <- function(sig2, Beta, yi, hi, v){
24   shape <- (v + 1)/2
25   rate <- (v + (yi - sum(hi * Beta))^2 / sig2)/2
26   rgamma(1, shape = shape, rate = rate)
27 }
28 PostBeta <- function(sig2, lambda, y, H){
29   k <- dim(H)[2]; B0i <- solve(1000*diag(k))
30   b0 <- rep(0, k)
31   Bn <- solve(B0i + sig2^(-1)*t(H) %*% diag(lambda) %*% H)
32   bn <- Bn %*% (B0i %*% b0 + sig2^(-1)*t(H) %*% diag(lambda) %*% y)
33   Beta <- MASS::mvrnorm(1, bn, Bn)
34   return(Beta)
35 }
36 PostSig2 <- function(Beta, lambda, y, H){
37   Ndk <- length(y); an <- a0 + Ndk
38   dn <- d0 + t(y - H %*% Beta) %*% diag(lambda) %*% (y - H %*% Beta)
39   sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
40   return(sig2)
41 }
42 # MCMC parameter
43 burnin <- 1000; S <- 5000; tot <- S + burnin
44 ETAs <- matrix(NA, tot, 3)
45 BETASOC <- matrix(NA, tot, 3); BETAS1C <- matrix(NA, tot, 3)
46 BETASA <- matrix(NA, tot, 2); BETASN <- matrix(NA, tot, 2)
47 SIGMAS <- matrix(NA, tot, 4)
48 b0c <- rep(0, 3); b1c <- rep(0, 3); b1a <- rep(0, 2); b0n <-
49   rep(0, 2)

```

**R code. Simulation: Fuzzy regression discontinuity design**

```

1 sigma21c <- 0.1; sigma20c <- 0.1; sigma20n <- 0.1; sigma21a
2 <- 0.1
3 EtacNew <- 0.5; EtanNew <- 0.25; EtaaNew <- 0.25
4 pb <- txtProgressBar(min = 0, max = tot, style = 3)
5 for(s in 1:tot){
6   ProbC0 <- sapply(I00, function(i)Pc0(b0c, sigma20c,
7     EtacNew, b0n, sigma20n, EtanNew, i = i))
8   ProbC1 <- sapply(I11, function(i)Pc1(b1c, sigma21c,
9     EtacNew, b1a, sigma21a, EtaaNew, i = i))
10  Typec0 <- sapply(ProbC0, function(p) {sample(c("c", "n"),
11    1, prob = c(p, 1-p))})
12  Typec1 <- sapply(ProbC1, function(p) {sample(c("c", "a"),
13    1, prob = c(p, 1-p))})
14  typeNew <- rep(NA_character_, N)
15  typeNew[I10] <- "n"; typeNew[I01] <- "a"
16  typeNew[I00] <- Typec0; typeNew[I11] <- Typec1
17  Nc <- sum(typeNew == "c"); Na <- sum(typeNew == "a")
18  Nn <- sum(typeNew == "n");
19  anc <- a0c + Nc; ann <- a0n + Nn; ana <- a0a + Na
20  alphas <- c(anc, ann, ana)
21  EtasNew <- rtype(alphas = alphas)
22  EtacNew <- EtasNew[1]; EtanNew <- EtasNew[2]
23  EtaaNew <- EtasNew[3]
24  ETAs[s, ] <- EtasNew; Lambda <- numeric(N)
25  idx_n <- which(typeNew == "n")
26  idx_a <- which(typeNew == "a")
27  idx_c0 <- which(typeNew == "c" & D == 0); idx_c1 <- which(
28    typeNew == "c" & D == 1)
29  Lambda[idx_n] <- sapply(idx_n, function(i) PostLambda_i(
30    sigma20n, b0n, Y[i], W[i, ], v))
31  Lambda[idx_a] <- sapply(idx_a, function(i) PostLambda_i(
32    sigma21a, b1a, Y[i], W[i, ], v))
33  Lambda[idx_c0] <- sapply(idx_c0, function(i) PostLambda_i(
34    sigma20c, b0c, Y[i], Wc[i, ], v))
35  Lambda[idx_c1] <- sapply(idx_c1, function(i) PostLambda_i(
36    sigma21c, b1c, Y[i], Wc[i, ], v))
37  b1a <- PostBeta(sig2 = sigma21a, lambda = Lambda[idx_a], y
38    = Y[idx_a], H = W[idx_a,])
39  b0n <- PostBeta(sig2 = sigma20n, lambda = Lambda[idx_n], y
40    = Y[idx_n], H = W[idx_n,])
41  b0c <- PostBeta(sig2 = sigma20c, lambda = Lambda[idx_c0],
42    y = Y[idx_c0], H = Wc[idx_c0,])
43  b1c <- PostBeta(sig2 = sigma21c, lambda = Lambda[idx_c1],
44    y = Y[idx_c1], H = Wc[idx_c1,])
45  BETASA[s, ] <- b1a; BETASN[s, ] <- b0n
46  BETASOC[s, ] <- b0c; BETAS1C[s, ] <- b1c
47  sigma21a <- PostSig2(Beta = b1a, lambda = Lambda[idx_a], y
48    = Y[idx_a], H = W[idx_a,])

```

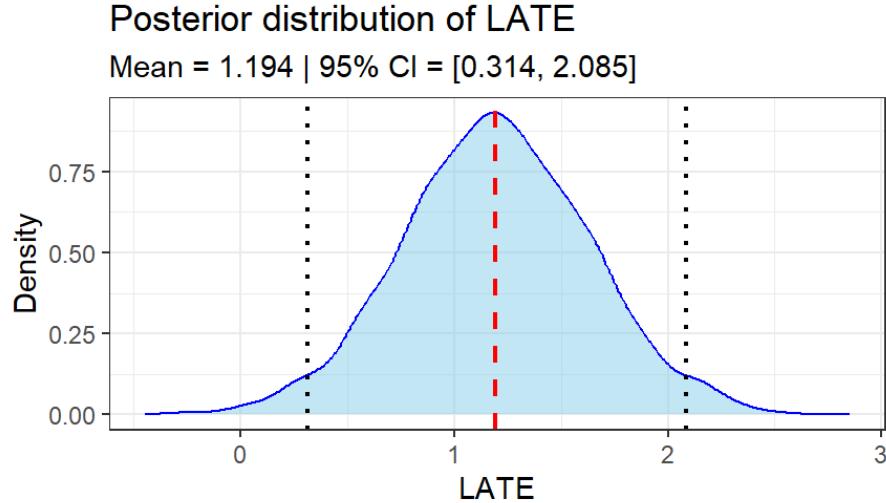
### R code. Simulation: Fuzzy regression discontinuity design

```

1 sigma20n <- PostSig2(Beta = b0n, lambda = Lambda[idx_n], y
2   = Y[idx_n], H = W[idx_n,])
3 sigma20c <- PostSig2(Beta = b0c, lambda = Lambda[idx_c0],
4   y = Y[idx_c0], H = Wc[idx_c0,])
5 sigma21c <- PostSig2(Beta = b1c, lambda = Lambda[idx_c1],
6   y = Y[idx_c1], H = Wc[idx_c1,])
7 SIGMAS[s, ] <- c(sigma21a, sigma20n, sigma20c, sigma21c)
8 setTxtProgressBar(pb, s)
9 }
10 close(pb)
11 keep <- seq((burnin+1), tot)
12 LATE <- coda::mcmc(BETAS1C[keep, 1] - BETASOC[keep, 1])
13 # Extract samples as numeric
14 late_draws <- as.numeric(LATE)
15 # Posterior mean and 95% CI
16 LATEmean <- mean(late_draws)
17 LATEci <- quantile(late_draws, c(0.025, 0.975))
18 # Plot posterior density
19 df <- data.frame(LATE = late_draws)
20 ggplot(df, aes(x = LATE)) + geom_density(fill = "skyblue",
21   alpha = 0.5, color = "blue") + geom_vline(xintercept =
22     LATEmean, color = "red", linetype = "dashed", linewidth
23     = 1) + geom_vline(xintercept = LATEci, color = "black",
24     linetype = "dotted", linewidth = 1) + labs(title =
25     "Posterior distribution of LATE",
26     subtitle = paste0("Mean = ", round(LATEmean,3), " | 95% CI =
27     [", round(LATEci[1],3), ", ", round(LATEci[2],3), "]"),
28     x = "LATE", y = "Density") +
29 theme_bw(base_size = 14)

```

A difference in this example is that, unlike most of the Frequentist literature, the entire dataset on the forcing variable is used. [94] propose a non-parametric Bayesian framework that emphasizes the data near the cutoff using a *soft-window* approach. Other Bayesian approaches to inference in RD include [58], who use Gaussian process regression to model the conditional expectation function, and [369], who develop a non-parametric Bayesian framework for spatial RD.

**FIGURE 13.20**

Posterior distribution of LATE: Simulation of fuzzy regression discontinuity design

### 13.7 Sample selection

In Figure 13.8 in Section 13.3, we depict a situation of collider bias that induces selection bias. Specifically, conditioning on a particular subset of the population ( $D_i = 1$ ), where  $D_i$  is influenced by both the treatment and confounders, opens a back-door path and creates a spurious association between their causes. Placing this situation within the well-known sample selection framework [198], the observed outcome can be represented as

$$Y_i = \begin{cases} Y_i(1) & \text{if } D_i = 1, \\ \text{NA} & \text{if } D_i = 0, \end{cases}$$

that is, we only observe  $Y_i = Y_i(1)$  for  $i = 1, 2, \dots, N$ , while  $Y_i(0)$  remains unobserved (“missing”).

In this setting, inference can be performed based on the likelihood of the observed outcomes together with the *selection (missingness) mechanism* ( $(Y_i(1), D_i)$ ), integrating out the unobserved  $Y_i(0)$  from the joint distribution of  $\{Y_i(1), Y_i(0), D_i\}$ . However, one must consider whether the missingness mechanism is *ignorable* or not. According to [275], the missingness mechanism is ignorable in Bayesian inference if the following two conditions hold:

1. the likelihood can be factorized as

$$p(y_i, d_i | \boldsymbol{\theta}, \boldsymbol{\gamma}) = p(y_i | \boldsymbol{\theta}) p(d_i | \boldsymbol{\gamma}),$$

and

2. the parameters  $\boldsymbol{\theta}$  and  $\boldsymbol{\gamma}$  are *a priori* independent,

$$\pi(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \pi(\boldsymbol{\theta}) \pi(\boldsymbol{\gamma}).$$

Under these conditions, posterior inference can be based on

$$\pi(\boldsymbol{\theta} | \mathbf{y}) \propto \pi(\boldsymbol{\theta}) p(\mathbf{y} | \boldsymbol{\theta}).$$

Thus, if the missingness mechanism is *Missing At Random* (MAR) and the parameters are *a priori* independent, the missingness mechanism is ignorable for Bayesian inference (see Chapter 6 of [275]).

In contrast, when the missingness mechanism is non-ignorable (i.e., Not Missing At Random, NMAR), the probability of observing  $Y_i$  depends on the unobserved values themselves. In this case, Bayesian inference must be performed using the full joint likelihood,

$$p(y_i, d_i | \boldsymbol{\theta}, \boldsymbol{\gamma}),$$

which requires specifying and estimating both the outcome model and the selection model simultaneously (see Chapter 15 of [275]). In particular, the classical sample selection model of [198] establishes,

$$Y_i = \begin{cases} \mathbf{X}_i^\top \boldsymbol{\beta} + \mu_i & \text{if } D_i = 1, \\ \text{NA} & \text{if } D_i = 0, \end{cases}$$

$$D_i^* = \begin{cases} 1 & \text{if } D_i^* = \mathbf{Z}_i^\top \boldsymbol{\gamma} + \nu_i > 0, \\ 0 & \text{if } D_i^* = \mathbf{Z}_i^\top \boldsymbol{\gamma} + \nu_i \leq 0, \end{cases}$$

where

$$\begin{bmatrix} \mu_i \\ \nu_i \end{bmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_\mu^2 & \sigma_{\mu\nu} \\ \sigma_{\mu\nu} & 1 \end{bmatrix}\right). \quad (13.3)$$

The restriction  $\text{Var}(\nu_i) = 1$  is imposed for identification. Without this normalization, the latent index  $D_i^* = \mathbf{Z}_i^\top \boldsymbol{\gamma} + \nu_i$  is only identified up to scale, since

$$P(D_i = 1) = P(\mathbf{Z}_i^\top \boldsymbol{\gamma} + \nu_i > 0) = P\left(\frac{\nu_i}{c} > -\frac{\mathbf{Z}_i^\top \boldsymbol{\gamma}}{c}\right), \quad \forall c > 0.$$

Thus, setting  $\sigma_\nu^2 = 1$  yields point identification of the parameters.

Note that since  $D_i = 1 \iff \nu_i > -\mathbf{Z}_i^\top \boldsymbol{\gamma}$ ,

$$\mathbb{E}[\mu_i | D_i = 1, \mathbf{Z}_i] = \mathbb{E}[\mathbb{E}(\mu_i | \nu_i) | \nu_i > -\mathbf{Z}_i^\top \boldsymbol{\gamma}] = \sigma_{\mu\nu} \mathbb{E}[\nu_i | \nu_i > -\mathbf{Z}_i^\top \boldsymbol{\gamma}].$$

This is because  $\mu_i | \nu_i \sim N(\sigma_{\mu\nu} \nu_i, \sigma_\mu^2 - \sigma_{\mu\nu}^2)$ .

If  $\nu \sim \mathcal{N}(0, 1)$ , then for  $a = \mathbf{Z}_i^\top \boldsymbol{\gamma}$ ,

$$\mathbb{E}[\nu | \nu > -a] = \frac{\phi(a)}{\Phi(a)} \equiv \lambda(a),$$

where  $\phi$  and  $\Phi$  are the standard normal pdf and cdf. Hence

$$\mathbb{E}[\mu_i | D_i = 1, \mathbf{Z}_i] = \sigma_{\mu\nu} \lambda(\mathbf{Z}_i^\top \boldsymbol{\gamma}) = \rho \sigma_\mu \lambda(\mathbf{Z}_i^\top \boldsymbol{\gamma}),$$

where  $\rho = \sigma_{\mu\nu} / \sigma_\mu$ .

Then,

$$\begin{aligned} \mathbb{E}[Y_i | D_i = 1, \mathbf{X}_i, \mathbf{Z}_i] &= \mathbf{X}_i^\top \boldsymbol{\beta} + \mathbb{E}[\mu_i | D_i = 1, \mathbf{Z}_i] \\ &= \mathbf{X}_i^\top \boldsymbol{\beta} + \sigma_{\mu\nu} \lambda(\mathbf{Z}_i^\top \boldsymbol{\gamma}) \\ &= \mathbf{X}_i^\top \boldsymbol{\beta} + \rho \sigma_\mu \lambda(\mathbf{Z}_i^\top \boldsymbol{\gamma}). \end{aligned}$$

This is the Heckman selection-bias correction: the inverse Mills ratio  $\lambda(\mathbf{Z}_i^\top \boldsymbol{\gamma})$  enters as an additional regressor in the selected sample. Point identification can be achieved through functional form (since the inverse Mills ratio is non-linear) together with the normality assumption. This implies that, in principle, one may have  $\mathbf{X}_i = \mathbf{Z}_i$ . However, such identification is weak, and it is preferable to include at least one variable in  $\mathbf{Z}_i$  that is excluded from  $\mathbf{X}_i$ , i.e.,  $\mathbf{X}_i \neq \mathbf{Z}_i$ , as this strengthens identification and also improves the precision of inference.

To perform Bayesian inference in this model, we use the augmented likelihood. Thus,

$$\begin{aligned} \pi(\boldsymbol{\delta}, \sigma_\mu^2, \sigma_{\mu\nu}, D_i^* | \mathbf{y}, \mathbf{d}) &\propto \prod_{i \in I_0} \pi(D_i^* | \boldsymbol{\delta}, \sigma_\mu^2, \sigma_{\mu\nu}) \mathbb{1}(d_i = 0) \mathbb{1}(D_i^* \leq 0) \\ &\quad \times \prod_{i \in I_1} p(y_i, D_i^* | \boldsymbol{\delta}, \sigma_\mu^2, \sigma_{\mu\nu}) \mathbb{1}(d_i = 1) \mathbb{1}(D_i^* > 0) \\ &\quad \times \pi(\boldsymbol{\delta}, \sigma_\mu^2, \sigma_{\mu\nu}), \end{aligned}$$

where  $\boldsymbol{\delta} = [\boldsymbol{\beta}^\top \boldsymbol{\gamma}^\top]^\top$ ,  $I_0 = \{i : d_i = 0\}$  and  $I_1 = \{i : d_i = 1\}$ .

Following Chapter 11 in [184], we set

$$\boldsymbol{\eta}_i = \begin{cases} [0, D_i^*]^\top, & i \in I_0 \\ [y_i, D_i^*]^\top, & i \in I_1 \end{cases}, \quad \mathbf{W}_i = \begin{bmatrix} \mathbf{X}_i^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_i^\top \end{bmatrix}, \quad \mathbf{J} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

Assuming  $\pi(\boldsymbol{\delta}) \sim N(\boldsymbol{\delta}_0, \mathbf{D}_0)$ , the conditional posterior distribution of  $\boldsymbol{\delta}$  is normal with mean

$$\boldsymbol{\delta}_n = \mathbf{D}_n \left[ \sum_{i \in I_0} \mathbf{W}_i^\top \mathbf{J} \boldsymbol{\eta}_i + \sum_{i \in I_1} \mathbf{W}_i^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\eta}_i + \mathbf{D}_0^{-1} \boldsymbol{\delta}_0 \right],$$

and variance matrix

$$\mathbf{D}_n = \left[ \sum_{i \in I_0} \mathbf{W}_i^\top \mathbf{J} \mathbf{W}_i + \sum_{i \in I_1} \mathbf{W}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{W}_i + \mathbf{D}_0^{-1} \right]^{-1},$$

where

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_\mu^2 & \sigma_{\mu\nu} \\ \sigma_{\mu\nu} & 1 \end{bmatrix}.$$

Let  $\omega = \sigma_\mu^2 - \sigma_{\mu\nu}$  denote the conditional variance of  $\mu_i \mid \nu_i$ . Assuming  $\omega^{-1} \sim G(\alpha_0/2, \delta_0/2)$  and noting that

$$p(y_i, D_i^* \mid \boldsymbol{\delta}, \omega, \sigma_{\mu\nu}) = p(y_i \mid D_i^*, \boldsymbol{\delta}, \omega, \sigma_{\mu\nu}) \times p(D_i^* \mid \boldsymbol{\delta}, \omega, \sigma_{\mu\nu}),$$

the conditional posterior distribution of  $\omega^{-1}$  is Gamma,

$$\omega^{-1} \mid \boldsymbol{\delta}, \sigma_{\mu\nu}, \mathbf{y}, \mathbf{d} \sim G(\alpha_n/2, \delta_n/2),$$

with

$$\alpha_n = \alpha_0 + N_1, \quad \delta_n = \delta_0 + \sum_{i \in I_1} \left[ y_i - \mathbf{X}_i^\top \boldsymbol{\beta} - \sigma_{\mu\nu}(D_i^* - \mathbf{Z}_i^\top \boldsymbol{\gamma}) \right]^2,$$

where  $N_1$  is the number of observations with  $D_i = 1$ .

In addition, assuming  $\sigma_{\mu\nu} \sim N(s_0, S_0)$ , the conditional posterior distribution is normal with mean

$$s_n = S_n \left( \omega^{-1} \sum_{i=1}^N (D_i^* - \mathbf{Z}_i^\top \boldsymbol{\gamma})(y_i - \mathbf{X}_i^\top \boldsymbol{\beta}) + S_0^{-1} s_0 \right),$$

and variance

$$S_n = \left[ \omega^{-1} \sum_{i=1}^N (D_i^* - \mathbf{Z}_i^\top \boldsymbol{\gamma})^2 + S_0^{-1} \right]^{-1}.$$

Finally, since

$$p(y_i, D_i^* \mid \boldsymbol{\delta}, \omega, \sigma_{\mu\nu}) = p(D_i^* \mid y_i, \boldsymbol{\delta}, \omega, \sigma_{\mu\nu}) \times p(y_i \mid \boldsymbol{\delta}, \omega, \sigma_{\mu\nu}),$$

the conditional posterior distribution of  $D_i^*$  is

$$D_i^* \sim \begin{cases} TN_{(-\infty, 0]}(\mathbf{Z}_i^\top \boldsymbol{\gamma}, 1), & i \in I_0, \\ TN_{(0, \infty)}\left(\mathbf{Z}_i^\top \boldsymbol{\gamma} + \frac{\sigma_{\mu\nu}}{\sigma_{\mu\nu}^2 + \omega} (y_i - \mathbf{X}_i^\top \boldsymbol{\beta}), \frac{\omega}{\sigma_{\mu\nu}^2 + \omega}\right), & i \in I_1, \end{cases}$$

where  $TN_A(\mu, \sigma^2)$  denotes a normal distribution with mean  $\mu$  and variance  $\sigma^2$ , truncated to the set  $A$ .

**Example: Simulation of the sample selection model**

We simulate the model

$$Y_i = 12 + X_{i1} + X_{i2} + \mu_i,$$

where  $X_{i1} \sim N(0, 4)$ ,  $X_{i2} \sim \text{Bin}(1, 0.5)$ , and  $\mu_i \sim N(0, 1.2)$  for  $i = 1, 2, \dots, 1,000$ . In addition,

$$D_i^* = 1 + Z_{i1} - X_{i2} + \nu_i,$$

where  $Z_{i1} \sim N(0, 4)$ .

The covariance between  $\mu_i$  and  $\nu_i$  is set to 0.6.

The hyperparameters are  $\delta_0 = [0 \ 0 \ 0 \ 0 \ 0]^T$ ,  $D_0 = 1,000\mathbf{I}_6$ ,  $\alpha_0 = \delta_0 = 0.001$ ,  $s_0 = 0$ , and  $S_0 = 1,000$ . We perform 1,500 iterations with a burn-in of 500. The following code illustrates the Gibbs sampler. Finally, we compare the results with an implementation that does not account for the sample selection issue.

Figure 13.21 shows the posterior distributions of the second parameter in the outcome equation, whose population value is 1 (black dashed line). The posterior distribution of the model that accounts for selection encompasses the population value, whereas the model that ignores selection does not.

### R code. Simulation: Sample selection model

```

1 rm(list = ls()); set.seed(10101)
2 library(doParallel); library(snow)
3 N <- 1000
4 w1 <- rbinom(N, 1, 0.5) # rnorm(N, 0, sigExo);
5 delta <- c(1, 1, -1); beta <- c(2,1,1)
6 zx <- MASS::mvrnorm(N, mu = rep(0, 2), matrix(c(1,0.7,0.7,1),
7 ,2,2))
7 z1 <- zx[,1]
8 Z <- cbind(1, z1, w1)
9 sig12 <- 0.8; sig11 <- 1.2
10 SIGMA <- matrix(c(sig11, sig12, sig12, 1), 2, 2)
11 E <- MASS::mvrnorm(N, mu = rep(0, 2), SIGMA)
12 cl <- Z%*%delta + E[,2]; c <- cl > 0
13 x1 <- zx[,2]; X <- cbind(1, x1, w1)
14 y <- X%*%beta + E[,1]
15 y[c==0] <- NA
16 # Hyperparameters
17 b0 <- rep(0, 6); B0 <- 1000*diag(6); B0i <- solve(B0)
18 a0 <- 0.001; d0 <- 0.001
19 s0 <- 0; S0 <- 1000; Soi <- 1/S0

```

**R code. Simulation: Sample selection model**

```

1 # Location
2 idc1 <- which(c==1)
3 ncl <- length(idc1)
4 PostThetaNew <- function(Sigma, clat){
5   J <- matrix(c(0,0,0,1),2,2)
6   WW <- matrix(0, 6, 6)
7   Wy <- matrix(0, 6, 1)
8   for(i in 1:N){
9     if(i %in% idc1){
10       yclat <- c(y[i], clat[i])
11       Auxi <- solve(Sigma)
12     }else{
13       yclat <- c(0, clat[i])
14       Auxi <- J
15     }
16     Wi <- as.matrix(Matrix:::bdiag(X[i,], Z[i,]))
17     WWi <- Wi%*%Auxi%*%t(Wi)
18     Wyi <- Wi%*%Auxi%*%yclat
19     WW <- WW + WWi
20     Wy <- Wy + Wyi
21   }
22   Bn <- solve(B0i + WW)
23   bn <- Bn%*%(B0i%*%b0 + Wy)
24   Beta <- MASS::mvrnorm(1, bn, Bn)
25   return(Beta)
26 }
27 PostOmega11 <- function(theta, sig12, clat){
28   an <- a0 + ncl
29   mui <- y[idc1] - X[idc1, ]%*%theta[1:3] - sig12*(clat[idc1]
30   - Z[idc1, ]%*%theta[4:6])
31   dn <- d0 + t(mui)%*%mui
32   omega11 <- LaplacesDemon::rinvgamma(1, an/2, dn/2)
33   return(omega11)
34 }
35 PostSig12 <- function(omega11, theta, clat){
36   Sn <- (omega11^(-1)*sum((clat[idc1] - Z[idc1, ]%*%theta
37   [4:6])^2) + S0i)^(-1)
38   sn <- Sn*(omega11^(-1)*sum((clat[idc1] - Z[idc1, ]%*%theta
39   [4:6])* (y[idc1] - X[idc1, ]%*%theta[1:3])) + s0*S0i)
40   sig12 <- rnorm(1, sn, sd = Sn^0.5)
41   return(sig12)
42 }
43 PostClat <- function(theta, sig12, omega11, i){
44   if(i %in% idc1){
45     mu <- Z[i, ]%*%theta[4:6] + (sig12/(omega11+sig12^2))*(y[
46     i] - X[i, ]%*%theta[1:3])
47     sig2 <- omega11/(omega11+sig12^2)
48     clat <- EnvStats::rnormTrunc(1, mean = mu, sd = sig2
49     ^0.5, min = 0, max = Inf)

```

### R code. Simulation: Sample selection model

```

1   }else{
2     mu <- Z[i,]*theta[4:6]
3     clat <- EnvStats::rnormTrunc(1, mean = mu, sd = 1, min =
4       -Inf, max = 0)
5   }
6   return(clat)
7 }
# Sampler
8 S <- 1500; burnin <- 500; thin <- 2
9 keep <- seq(burnin+thin, S, thin)
10 PostThetasDraws <- matrix(NA, S, 6)
11 PostSigmaDraws <- matrix(NA, S, 2)
12 # Initial parameters
13 Thetap <- rep(0, 6); Sigmap <- diag(2)
14 Sig12p <- 0; Omega11p <- 1
15 for(i in 1:N){
16 if(c[i] == 0){
17   LatPost <- EnvStats::rnormTrunc(1, mean = 0, sd = 1, min =
18     -Inf, max = 0)
19   }else{
20     LatPost <- EnvStats::rnormTrunc(1, mean = 0, sd = 1, min
21     = 0, max = Inf)
22   }
23 #### Parallel code #####
24 cn <- detectCores()
25 ClusterHope <- makeCluster(cn, type = "SOCK")
26 registerDoParallel(ClusterHope)
27 clusterExport(ClusterHope, list("Z", "X", "c", "y", "N",
28   "idc1", "nc1", "PostClat", "Thetap", "Sig12p", "Omega11p"))
29 pb <- txtProgressBar(min = 0, max = S, style = 3)
30 for(rep in 1:S){
31   LatPost <- t(parSapply(ClusterHope, 1:N, function(i){
32     PostClat(theta = Thetap, sig12 = Sig12p, omega11 =
33       Omega11p, i)}))
34   Thetap <- PostThetaNew(Sigma = Sigmap, clat = LatPost)
35   Omega11p <- PostOmega11(theta = Thetap, sig12 = Sig12p,
36     clat = LatPost)
37   Sig12p <- PostSig12(omega11 = Omega11p, theta = Thetap,
38     clat = LatPost)
39   Sigmap <- matrix(c(Omega11p+Sig12p^2, Sig12p, Sig12p, 1
40     ,2,2)
41   PostThetasDraws[rep,] <- Thetap
42   PostSigmaDraws[rep, ] <- c(Omega11p+Sig12p^2, Sig12p)
43   clusterExport(ClusterHope, list("Thetap", "Sig12p", "
44     Omega11p"))
45   setTxtProgressBar(pb, rep)
46 }

```

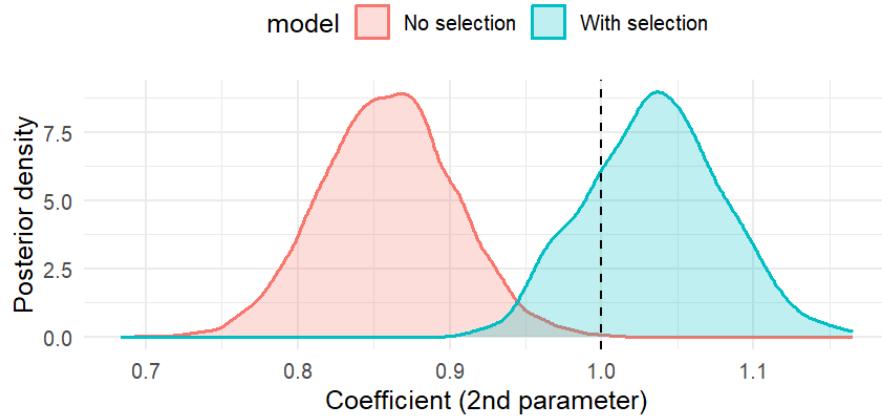
**R code. Simulation: Sample selection model**

```

1 stopCluster(ClusterHope)
2 close(pb)
3 thetaHat <- coda::mcmc(PostThetasDraws[keep ,])
4 PostDrawsSigma <- coda::mcmc(PostSigmaDraws[keep ,])
5 summary(thetaHat)
6 summary(PostDrawsSigma)
7 RegN0sel <- MCMCpack::MCMCregress(y ~ x - 1)
8 summary(RegN0sel)

```

Posterior of 2nd Coeff: No-Selection vs Selection Models  
Dashed line = population value (1)



**FIGURE 13.21**

Posterior distribution second parameter in outcome: Selection vs no selection models

The original Heckman selection model was not intended to calculate the ATE, since  $Y_i(0)$  is not observed and, therefore, a treatment effect is not defined. Its main purpose is to correct the expected value of  $Y_i \mid D_i = 1$  by accounting for sample selection. Subsequently, [199] and [201] extended the sample selection framework to incorporate  $Y_i(0)$ , framing the Roy model [384] within the treatment effect literature and establishing the identification conditions for different treatment parameters. In particular, [201] introduced

the marginal treatment effect (MTE),

$$\tau_{MTE}(x, u_D) = \mathbb{E}[Y_i(1) - Y_i(0) \mid \mathbf{X}_i = \mathbf{x}, U_{Di} = u_D],$$

which is the expected effect of treatment for individuals with observed characteristics  $\mathbf{X}_i = \mathbf{x}$  and unobservable factors from the treatment assignment  $U_{Di} = u_D$  [201].

The MTE is a unifying concept that connects the treatment effect, selection, and matching literatures. Moreover, these authors show that many treatment effect parameters, such as the ATE, ATT, and LATE, can be expressed as weighted averages of the MTE (see Table IA in [201]).

In the Bayesian literature, several authors have estimated different versions of selection models. In particular, [440] and [122] analyze the sample selection model using flexible specifications for the error terms, while [249] and [90] estimate the Roy model. Furthermore, [93] extend the basic framework to a semiparametric setting that accounts for endogeneity, and [200] introduce latent factors to address the fundamental problem of causal inference in the Roy model. This approach allows researchers to learn about the otherwise unidentified covariance between  $Y_i(1)$  and  $Y_i(0)$ , which in turn enables the estimation of distributional versions of treatment effects that require the joint distribution of  $Y_i(1)$  and  $Y_i(0)$ .

### 13.8 Bayesian exponentially tilted empirical likelihood

Bayesian parametric approaches are often criticized because they require distributional assumptions that may be arbitrary or remain unchecked. For example, in this chapter we model a continuous outcome as normally distributed. This choice is defensible: among all distributions with a fixed mean and variance, the normal imposes the least prior structure (it maximizes entropy; see Exercise 2), and the same reasoning extends to regression via Gaussian errors with fixed conditional variance [467]. Nevertheless, in some settings it may be preferable to use partial information methods that rely only on moment conditions rather than full distributional assumptions. The trade-off is familiar: unless the parametric model is correctly specified, these semiparametric approaches typically reduce efficiency relative to a well-specified parametric model.

The point of departure is a set of moment conditions

$$\mathbb{E}[\mathbf{g}(\mathbf{W}, \boldsymbol{\theta})] = \mathbf{0}_d,$$

where the expectation is with respect to the population distribution,  $\mathbf{W}_{1:N} := [\mathbf{W}_1 \ \mathbf{W}_2 \ \dots \ \mathbf{W}_N]$  is a random sample from  $\mathbf{W} \subset \mathbb{R}^{d_w}$ ,  $\mathbf{g} : \mathbb{R}^{d_w} \times \Theta \rightarrow \mathbb{R}^d$  is a vector of known functions, and  $\boldsymbol{\theta} = [\theta_1 \ \theta_2 \ \dots \ \theta_p]^\top \in \Theta \subset \mathbb{R}^p$ . If  $d > p$

the model is *over-identified*; if  $d = p$  it is *exactly identified* (and if  $d < p$  it is *under-identified*).

**Example: Linear regression**

Let

$$y_i = \mathbf{X}_i^\top \boldsymbol{\beta} + \mu_i, \quad \mathbb{E}[\mu_i | \mathbf{X}_i] = 0,$$

with  $\mathbf{X}_i \in \mathbb{R}^p$  and  $\boldsymbol{\beta} \in \mathbb{R}^p$ . Then the unconditional moment conditions are

$$\mathbb{E}[\mathbf{X}_i \mu_i] = \mathbb{E}[\mathbf{X}_i (y_i - \mathbf{X}_i^\top \boldsymbol{\beta})] = \mathbf{0}_p.$$

**Example: Instrumental variables**

If there is endogeneity,  $\mathbb{E}[\mu_i | \mathbf{X}_i] \neq 0$ , but there exist instruments  $\mathbf{Z}_i \in \mathbb{R}^d$  that are exogenous,  $\mathbb{E}[\mu_i | \mathbf{Z}_i] = 0$ , and relevant,  $\text{rank}(\mathbb{E}[\mathbf{Z}_i \mathbf{X}_i^\top]) = p$ , then

$$\mathbb{E}[\mathbf{Z}_i \mu_i] = \mathbb{E}[\mathbf{Z}_i (y_i - \mathbf{X}_i^\top \boldsymbol{\beta})] = \mathbf{0}_d,$$

with  $d \geq p$ .

Moment conditions can be used in a Bayesian framework via *Bayesian Empirical Likelihood* (BEL) [258] and *Bayesian Exponentially Tilted Empirical Likelihood* (BETEL) [397]. We focus on BETEL because, while BEL inherits the attractive properties of empirical likelihood under correct specification, it can lose them under model misspecification. In contrast, Exponentially Tilted Empirical Likelihood (ETEL) remains well behaved under misspecification and retains root- $N$  consistency and asymptotic normality [398].

Thus, the posterior distribution is

$$\pi(\boldsymbol{\theta} | \mathbf{W}_{1:N}) \propto \pi(\boldsymbol{\theta}) L_{\text{ETEL}}(\boldsymbol{\theta}), \quad L_{\text{ETEL}}(\boldsymbol{\theta}) = \prod_{i=1}^N p_i^*(\boldsymbol{\theta}),$$

where  $\pi(\boldsymbol{\theta})$  is the prior and  $L_{\text{ETEL}}$  is the exponentially tilted empirical likelihood. The weights  $(p_1^*(\boldsymbol{\theta}), \dots, p_N^*(\boldsymbol{\theta}))$  are obtained from the maximum-entropy problem

$$\max_{\{p_i\}_{i=1}^N} \left\{ -\sum_{i=1}^N p_i \log p_i \right\} \quad \text{subject to} \quad \sum_{i=1}^N p_i = 1, \quad p_i \geq 0, \quad \sum_{i=1}^N p_i \mathbf{g}(\mathbf{W}_i, \boldsymbol{\theta}) = \mathbf{0}_d.$$

Equivalently (dual/saddlepoint form; see [397, 398, 98]),

$$p_i^*(\boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\lambda}(\boldsymbol{\theta})^\top \mathbf{g}(\mathbf{W}_i, \boldsymbol{\theta}))}{\sum_{j=1}^N \exp(\boldsymbol{\lambda}(\boldsymbol{\theta})^\top \mathbf{g}(\mathbf{W}_j, \boldsymbol{\theta}))}, \quad \text{where} \quad \sum_{i=1}^N p_i^*(\boldsymbol{\theta}) \mathbf{g}(\mathbf{W}_i, \boldsymbol{\theta}) = \mathbf{0}_d.$$

$\boldsymbol{\lambda}(\boldsymbol{\theta})$  can be characterized as

$$\boldsymbol{\lambda}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\lambda} \in \mathbb{R}^d} \log \left( \frac{1}{N} \sum_{i=1}^N \exp(\boldsymbol{\lambda}^\top \mathbf{g}(\mathbf{W}_i, \boldsymbol{\theta})) \right), \quad (13.4)$$

whose gradient condition is precisely the moment constraint above. Therefore, the BETEL posterior is

$$\pi(\boldsymbol{\theta} | \mathbf{w}_{1:N}) \propto \pi(\boldsymbol{\theta}) \prod_{i=1}^N \frac{\exp(\boldsymbol{\lambda}(\boldsymbol{\theta})^\top \mathbf{g}(\mathbf{W}_i, \boldsymbol{\theta}))}{\sum_{j=1}^N \exp(\boldsymbol{\lambda}(\boldsymbol{\theta})^\top \mathbf{g}(\mathbf{W}_j, \boldsymbol{\theta}))}.$$

Posterior inference of the BETEL can be performed using a Metropolis-Hastings algorithm where the proposal distribution is  $q(\boldsymbol{\theta} | \mathbf{W}_{1:N})$ . See Algorithm A35 [98].

---

**Algorithm A35** Bayesian Exponentially Tilted Empirical Likelihood: Metropolis-Hastings algorithm

---

- 1: Set  $\boldsymbol{\theta}^0$  in the support of  $\pi(\boldsymbol{\theta} | \mathbf{W}_{1:N})$
- 2: **for**  $s = 1, \dots, S$  **do**
- 3:     Propose  $\boldsymbol{\theta}^c$  from  $q(\boldsymbol{\theta} | \mathbf{W}_{1:N})$ , and solve for  $p_i^*(\boldsymbol{\theta}^c)$ ,  $i = 1, 2, \dots, N$  in Equation 13.4
- 4:     Calculate the acceptance probability

$$\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c | \mathbf{W}_{1:N}) = \min \left( 1, \frac{\pi(\boldsymbol{\theta}^c | \mathbf{W}_{1:N})}{\pi(\boldsymbol{\theta}^{(s-1)} | \mathbf{W}_{1:N})} \frac{q(\boldsymbol{\theta}^{(s-1)} | \mathbf{W}_{1:N})}{q(\boldsymbol{\theta}^c | \mathbf{W}_{1:N})} \right)$$

- 5:     Draw  $U$  from  $U(0, 1)$
  - 6:      $\boldsymbol{\theta}^{(s)} = \begin{cases} \boldsymbol{\theta}^c & \text{if } U \leq \alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c | \mathbf{W}_{1:N}) \\ \boldsymbol{\theta}^{(s-1)} & \text{otherwise} \end{cases}$
  - 7: **end for**
- 

**Example: Classical measurement error in regressor**

Let's set the unobserved (latent) regressor  $X_i^*$  such that

$$X_i = X_i^* + \nu_i,$$

where  $X_i$  is the observed regressor, and  $\nu_i$  is a *classical measurement error* such that  $\mathbb{E}[\nu_i] = 0$  and  $\nu_i \perp \{X_i^*, \mu_i\}$ , where  $\mu_i$  is the stochastic error of the structural model

$$Y_i = \beta_0 + \beta_1 X_i^* + \mu_i,$$

where  $\mathbb{E}[\mu_i] = 0$  and  $\mu_i \perp X_i^*$ .

If we perform the regression using the observed regressor, then

$$Y_i = \beta_0 + \beta_1(X_i - \nu_i) + \mu_i = \beta_0 + \beta_1 X_i + \underbrace{\mu_i - \beta_1 \nu_i}_{\epsilon_i},$$

where the new error term is  $\epsilon_i = \mu_i - \beta_1 \nu_i$ . We will show that

$$\mathbb{E}[\epsilon_i | X_i] \neq 0.$$

By the law of iterated expectations,

$$\mathbb{E}[\nu_i X_i] = \mathbb{E}[X_i \mathbb{E}[\nu_i | X_i]].$$

This implies:

- If  $\mathbb{E}[\nu_i | X_i] = 0$  almost surely, then  $\mathbb{E}[\nu_i X_i] = 0$ .
- If  $\mathbb{E}[\nu_i X_i] \neq 0$ , then  $\mathbb{E}[\nu_i | X_i]$  cannot be equal to zero almost surely.

Now compute

$$\mathbb{E}[\nu_i X_i] = \mathbb{E}[\nu_i(X_i^* + \nu_i)] = \underbrace{\mathbb{E}[\nu_i X_i^*]}_{=0} + \mathbb{E}[\nu_i^2] = \sigma_\nu^2 \neq 0.$$

Hence it must be that  $\mathbb{E}[\nu_i | X_i] \neq 0$ . Therefore,

$$\mathbb{E}[\epsilon_i | X_i] = \underbrace{\mathbb{E}[\mu_i | X_i]}_{=0} - \beta_1 \underbrace{\mathbb{E}[\nu_i | X_i]}_{\neq 0} \neq 0,$$

that is, the regressor is not exogenous, and consequently,

$$\mathbb{E}[X_i \underbrace{(Y_i - \beta_0 - \beta_1 X_i)}_{\epsilon_i}] \neq 0,$$

thus, we need a relevant (strong) ( $\mathbb{E}[Z_i X_i^*] \neq 0$ ) and exogeneous ( $\mathbb{E}[Z_i | \mu_i] = \mathbb{E}[Z_i | \nu_i] = 0$ ) instrument to identify the causal effect. Therefore,

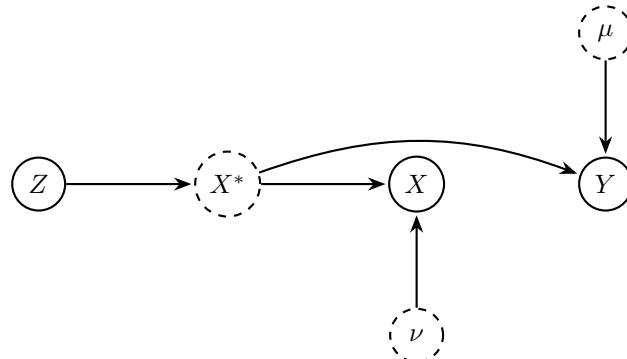
$$\mathbb{E}\left[\begin{bmatrix} 1 \\ Z_i \end{bmatrix} (Y_i - \beta_0 - \beta_1 X_i)\right] = \mathbf{0}.$$

The DAG 13.22 illustrates the situation of measurement error, and how the instrument helps to identify the causal effect. The instrument solves the endogeneity problem because it exploits variation in the regressor  $X_i$  that is correlated with the true latent regressor  $X_i^*$  but uncorrelated with the measurement error  $\nu_i$ . See Chapter 9 in [204] for a nice review of the effects of measurement error in causal inference.

We simulate the latent process  $X_i^* = 0.5Z_i + e_i$ , with the observed regressor defined as  $X_i = X_i^* + \nu_i$ , and the outcome equation specified as  $Y_i = 1 + 1.2X_i^* + \mu_i$ . The variables  $Z_i$ ,  $e_i$ , and  $\nu_i$  are standard normal, while  $\mu_i$  follows a mixture of two normal distributions with means 0.5 and -0.5, and standard deviations 0.5 and 1.2. The sample size is 2,000, the burn-in is 1,000, and the number of MCMC draws retained after burn-in is 10,000.

We perform Bayesian exponentially tilted empirical likelihood (BETEL) using the package *betel*, and adopt the default hyperparameter values provided there.<sup>4</sup> This package implements Bayesian estimation and marginal likelihood computation for moment condition models, employing a Student-t prior distribution and a Student-t proposal distribution following [98]. We also compare

<sup>4</sup> Available at <https://apps.olin.wustl.edu/faculty/chib/rpackages/betel/>.

**FIGURE 13.22**

DAG with measurement error and instrument  $Z$ . Observed  $X$  depends on the latent  $X^*$  and error  $\nu$ . Because the composed error in the estimated equation includes  $\nu$ ,  $X$  is endogenous when regressed on  $Y$ . Instrument  $Z$  isolates exogenous variation in  $X^*$ .

the results with a model that assumes  $X_i$  is exogenous, and with an instrumental Gibbs sampler that assumes  $Y_i$  is normally distributed. We also use the default hyperparameter of the used packages. The following code shows the implementation.

Figure 13.23 displays the posterior distributions. We see that ignoring measurement error produces a biased posterior distribution. In particular, the absolute value of the causal effect is smaller than the population value; this is called *attenuation bias*. By contrast, methods that account for measurement error and use valid instruments yield well-centered posterior distributions.

**R code. Simulation: Classical measurement error in regressor**

```

1 rm(list = ls()); set.seed(10101)
2 library(betel); library(ucminf)
3 # Simulate data
4 N <- 2000; d <- 2; k <- 2
5 gamma <- 0.5; beta <- c(1, 1.2)
6 # Mixture
7 mum1 <- 1/2; mum2 <- -1/2
8 mu1 <- rnorm(N, mum1, 0.5); mu2 <- rnorm(N, mum2, 1.2)

```

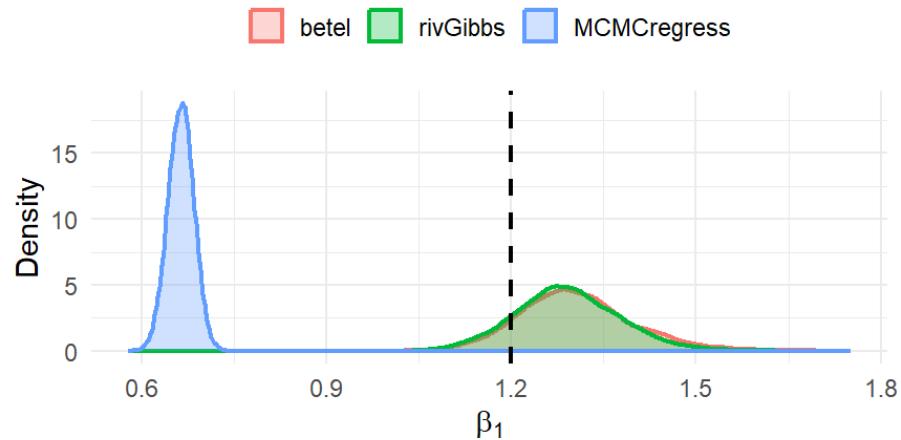
**R code. Simulation: Classical measurement error in regressor**

```

1 mu <- sapply(1:N, function(i){sample(c(mu1[i], mu2[i]), 1,
2   prob = c(0.5, 0.5))})
3 e <- rnorm(N)
4 z <- rnorm(N) # Instrument
5 xlat <- gamma*z + e # Unobserved regressor
6 nu <- rnorm(N) # Measurement error
7 x <- xlat + nu # Observed regressor
8 Xlat <- cbind(1, xlat)
9 y <- Xlat%*%beta + mu
10 dat <- cbind(1, x, z) # Data
11 # Function g_i by row in BETEL
12 gfunc <- function(psi = psi, y = y, dat = dat) {
13   X <- dat[,1:2]
14   e <- y - X %*% psi
15   E <- e %*% rep(1,d)
16   Z <- dat[,c(1,3)]
17   G <- E * Z;
18   return(G)
19 }
20 nt <- round(N * 0.1, 0); # training sample size for prior
21 psi0 <- lm(y[1:nt] ~ x[1:nt])$coefficients # Starting value of
22           psi = (theta, v), v is the slack parameter in CSS
23           (2018)
24 names(psi0) <- c("alpha","beta")
25 psi0_ <- as.matrix(psi0) # Prior mean of psi
26 Psi0_ <- 5*rep(1,k) # Prior dispersions of psi
27 lam0 <- .5*rnorm(d) # Starting value of lambda
28 nu <- 2.5 # df of the prior student-t
29 nuprop <- 15 # df of the student-t proposal
30 n0 <- 1000 # burn-in
31 m <- 10000 # iterations beyond burn-in
32 # MCMC ESTIMATION BY THE CSS (2018) method
33 psim <- betel::bayesetel(gfunc = gfunc, y = y[-(1:nt)], dat
34   = dat[-(1:nt),], psi0 = psi0, lam0 = lam0, psi0_ = psi0_,
35   , Psi0_ = Psi0_, nu = nu, nuprop = nuprop,
36   controlpsi = list(maxiterpsi = 50, mingrpsi = 1.0e-8), #
37   list of parameters in maximizing likelihood over psi
38   controllam = list(maxiterlam = 50, # list of parameters in
39   minimizing dual over lambda
40   mingrlam = 1.0e-7),
41   n0 = n0, m = m)
42 MCMCexg <- MCMCpack::MCMCregress(y ~ x, burnin = n0, mcmc =
43   m)
44 Data <- list(y = c(y), x = x, z = matrix(z, N, 1), w =
45   matrix(rep(1, N), N, 1))
46 Mcmc <- list(R = m)
47 MCMCivr <- bayesm::rivGibbs(Data, Mcmc = Mcmc)

```

### Posterior Densities with Population Value

**FIGURE 13.23**

Posterior distribution slope parameter given measurement error: betel, ivGibbs and MCMCregress

#### **Example: Omission of relevant correlated regressor**

Consider the true model

$$Y_i = \beta_0 + \beta_1 X_i + \gamma W_i + \mu_i, \quad \mathbb{E}[\mu_i | X_i, W_i] = 0.$$

Suppose that the relevant regressor  $W_i$  is omitted from the specification, with  $\gamma \neq 0$  and  $\text{Cov}(X_i, W_i) \neq 0$ .

Thus, if we regress  $Y_i$  only on  $X_i$ , the composite error is

$$\epsilon_i = \gamma W_i + \mu_i.$$

Therefore,

$$\mathbb{E}[\epsilon_i | X_i] = \gamma \mathbb{E}[W_i | X_i] + \mathbb{E}[\mu_i | X_i] = \gamma \mathbb{E}[W_i | X_i] \neq 0,$$

since  $\mathbb{E}[W_i | X_i]$  is nonconstant whenever  $X_i$  and  $W_i$  are correlated, that is, the expected value of  $W_i$  is a function of  $X_i$ .

Consequently,

$$\mathbb{E}[\epsilon_i X_i] = \gamma \mathbb{E}[W_i X_i] + \mathbb{E}[\mu_i X_i].$$

Because  $\mathbb{E}[\mu_i X_i] = 0$ , we obtain

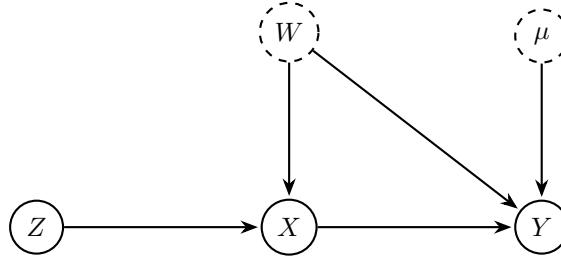
$$\mathbb{E}[\epsilon_i X_i] = \gamma (\text{Cov}(W_i, X_i) + \mathbb{E}[W_i] \mathbb{E}[X_i]) \neq 0.$$

In this situation, we can use an instrument that is relevant ( $\mathbb{E}[Z_i X_i] \neq 0$ )

and exogenous ( $\mathbb{E}[Z_i \mid \mu_i] = \mathbb{E}[Z_i \mid W_i] = 0$ ) to address the endogeneity problem, and identify the causal effect.<sup>5</sup> Therefore,

$$\mathbb{E} \left[ \begin{bmatrix} 1 \\ Z_i \end{bmatrix} (Y_i - \beta_0 - \beta_1 X_i) \right] = \mathbf{0}.$$

Figure 13.24 illustrates the case of an omitted relevant regressor that is correlated with  $X_i$ , and how an instrument can be used to identify the causal effect. The instrument resolves the endogeneity problem by exploiting variation in  $X_i$  that is driven by  $Z_i$ , which is orthogonal to the problematic error term. In other words, it replaces “bad” correlation with “clean” variation.



**FIGURE 13.24**

DAG with omitted relevant regressor  $W$  (dashed, unobserved) that is correlated with  $X$  and affects  $Y$ , inducing endogeneity of  $X$  in a regression of  $Y$  on  $X$ . A valid instrument  $Z$  affects  $X$  (relevance) but has no direct path to  $Y$  and is independent of  $W$  and  $u$  (exclusion/independence), enabling identification of the causal effect  $X \rightarrow Y$ .

We ask in Exercise 10 to perform a simulation exercise to assess the ability of BETEL to identify the causal effect when an instrument is used to address the omission of relevant regressors.

#### Example: Simultaneous causality

This example is based on [214], where Professor Imbens illustrates the problem of analyzing the causal effect on traded quantity of a tax of  $100 \times r\%$  in a market. He defines the average causal effect on the logarithm of traded quantity as

$$\tau = \mathbb{E}[q_i(r) - q_i(0)],$$

where  $q_i(r) = \log Q_i(r)$  and  $Q_i(r)$  denotes the potential traded quantity if the tax were  $r\%$ .

This situation is more challenging than in the standard treatment effects literature, because we do not observe any unit facing the tax. Instead, we only observe all units facing no tax, that is,  $Q_i^{\text{Obs}} = Q_i(0)$ . This setting requires the

---

<sup>5</sup>Another potential solution for the omission of relevant variables is the use of *proxy variables*; see [455] for details.

use of a *structural model* to define the counterfactual scenarios of the potential outcomes.

The starting point for inference on the treatment effect of this new tax is the price determination mechanism, that is, the assignment mechanism in the potential outcome framework. We specify a *structural demand function* that defines the potential demand given the price (treatment) and other exogenous variables:

$$q_i^d(p) = \beta_1 + \beta_2 p + \beta_3 y_i + \beta_4 pc_i + \beta_5 ps_i + u_{i1},$$

where  $q^d$  is demand, and  $p$ ,  $y$ ,  $pc$ , and  $ps$  are the logarithms of price, income, the price of a complementary good, and the price of a substitute good, respectively. All coefficients are interpreted as demand elasticities. The term  $u_{i1}$  represents unobserved demand factors such that  $\mathbb{E}[u_{i1}] = 0$ . Therefore,

$$\mathbb{E}[q_i^d(p) | p, y_i, pc_i, ps_i] = \beta_1 + \beta_2 p + \beta_3 y_i + \beta_4 pc_i + \beta_5 ps_i.$$

This expectation does not represent the conditional expectation of the observed quantity in markets where the observed price equals  $p$ . Rather, it is the expectation of potential demand functions given  $p$  and other exogenous controls, irrespective of the realized market price.

Similarly, the assignment mechanism requires specifying the *structural supply function*:

$$q_i^s(p) = \alpha_1 + \alpha_2 p + \alpha_3 er_i + u_{i2},$$

which represents the quantity that sellers are willing to supply given the price and the (exogenous) exchange rate  $er_i$ . The unobserved supply factors satisfy  $\mathbb{E}[u_{i2}] = 0$ , so that

$$\mathbb{E}[q_i^s(p) | p, er_i] = \alpha_1 + \alpha_2 p + \alpha_3 er_i.$$

This expectation represents the average of all potential supply functions given the price and exogenous controls, again irrespective of the realized market price.

Thus, the assignment mechanism is given by the market equilibrium where the observed price is such that the observed quantity is equal to the demand and supply potential outcomes at the observed price,

$$q_i^{Obs} = q_i^d(p_i^{Obs}) = q_i^s(p_i^{Obs}).$$

Using this market equilibrium condition, we get

$$p_i^{Obs} = \pi_1 + \pi_2 er_i + \pi_3 y_i + \pi_4 pc_i + \pi_5 ps_i + v_{i1},$$

where  $\pi_1 = \frac{\alpha_1 - \beta_1}{\beta_2 - \alpha_2}$ ,  $\pi_2 = \frac{\alpha_3}{\beta_2 - \alpha_2}$ ,  $\pi_3 = \frac{-\beta_3}{\beta_2 - \alpha_2}$ ,  $\pi_4 = \frac{-\beta_4}{\beta_2 - \alpha_2}$ ,  $\pi_5 = \frac{-\beta_5}{\beta_2 - \alpha_2}$ , and  $v_{i1} = \frac{u_{i2} - u_{i1}}{\beta_2 - \alpha_2}$  given  $\beta_2 \neq \alpha_2$ , that is, the equations should be independent. This condition is given by economic theory due to  $\beta_2 < 0$  and  $\alpha_2 > 0$ , the effect of price on demand and supply should be negative and positive, respectively.

The equation of price into the demand equation gives

$$q_i^{Obs} = \tau_1 + \tau_2 er_i + \tau_3 y_i + \tau_4 pc_i + \tau_5 ps_i + v_{i2},$$

where  $\tau_1 = \beta_1 + \beta_2 \pi_1$ ,  $\tau_2 = \beta_2 \pi_2$ ,  $\tau_3 = \beta_2 \pi_3 + \beta_3$ ,  $\tau_4 = \beta_2 \pi_4 + \beta_4$ ,  $\tau_5 = \beta_2 \pi_5 + \beta_5$ , and  $v_{i2} = \beta_2 v_{i1} + u_{i1}$ .

The expressions for  $p_i^{Obs}$  and  $q_i^{Obs}$  are called the *reduced-form* representations. In Section 7.1, we presented the order condition, which is necessary, and the rank condition, which is both necessary and sufficient, to identify the structural parameters from the reduced-form parameters. A key point to note is that only the prior distribution of the reduced-form parameters is updated by the sample information, while the updating of the structural parameters occurs solely through the reduced-form parameters, that is,

$$\pi(\boldsymbol{\beta}, \boldsymbol{\alpha} | \boldsymbol{\gamma}, \boldsymbol{\pi}, \mathbf{W}_{1:N}) \propto \pi(\boldsymbol{\beta}, \boldsymbol{\alpha} | \boldsymbol{\gamma}, \boldsymbol{\pi}).$$

See Section 9.3 in [465] for details about identification in Bayesian inference. This implies that in *under-identified* models, the posterior distribution of the structural parameters is not concentrated at a single point, but rather remains spread out (e.g., uniformly) over a range of values.<sup>6</sup>

To analyze the causal effects of the new tax, we can find the new equilibrium,

$$q_i^d(P_i(r) \times (1+r)) = q_i^s(P_i(r)),$$

where  $P_i(r)$  is the price level ( $p_i(r) = \log P_i(r)$ ) that sellers get, and  $(P_i(r) \times (1+r))$  is the price that buyers pay.

Let's define the (log) price that buyers pay  $p_b = p + \log(1+r)$  while sellers receive  $p$ . The demand function becomes

$$q_i^d(p + \log(1+r)) = \beta_1 + \beta_2(p + \log(1+r)) + \beta_3 y_i + \beta_4 pc_i + \beta_5 ps_i + u_{i1},$$

and the supply function remains

$$q_i^s(p) = \alpha_1 + \alpha_2 p + \alpha_3 er_i + u_{i2}.$$

Thus, the equilibrium price is given by

$$p_i^*(r) = \frac{(\beta_1 - \alpha_1) + \beta_3 y_i + \beta_4 pc_i + \beta_5 ps_i - \alpha_3 er_i + (u_{i1} - u_{i2}) + \beta_2 \log(1+r)}{\alpha_2 - \beta_2}.$$

The equilibrium quantity is

$$q_i^*(r) = \beta_1 + \beta_2(p_i^*(r) + \log(1+r)) + \beta_3 y_i + \beta_4 pc_i + \beta_5 ps_i + u_{i1}.$$

Thus, the expected equilibrium price and quantity are

$$\mathbb{E}[p_i^*(r) | y_i, pc_i, ps_i, er_i] = \frac{(\beta_1 - \alpha_1) + \beta_3 y_i + \beta_4 pc_i + \beta_5 ps_i - \alpha_3 er_i + \beta_2 \log(1+r)}{\alpha_2 - \beta_2},$$

---

<sup>6</sup>This identification issue is not peculiar to simultaneous equation models, it arises in other econometric/statistical models [465]. See for instance the example of the effects of vitamin A, Figures 13.12 and 13.13.

$$\mathbb{E}[q_i^*(r) | y_i, pc_i, ps_i, er_i] = \alpha_1 + \alpha_2 \mathbb{E}[p_i^*(r) | \cdot] + \alpha_3 er_i.$$

We can see that given  $\beta_2 < 0 < \alpha_2$ , we have:

$$\begin{aligned} \frac{dp^*(r)}{dr} &= \frac{\beta_2}{(\alpha_2 - \beta_2)(1+r)} < 0, & \frac{dp_b^*(r)}{dr} &= \frac{\alpha_2}{(\alpha_2 - \beta_2)(1+r)} > 0, \\ \frac{dq^*(r)}{dr} &= \alpha_2 \cdot \frac{dp^*(r)}{dr} < 0. \end{aligned}$$

Thus, the price received by sellers decreases with the tax, the price paid by buyers increases, and the equilibrium quantity falls.

Thus, the average causal effect on the logarithm of traded quantity is

$$\tau = \mathbb{E}[q_i(r) - q_i(0)] = \frac{\alpha_2 \beta_2}{\alpha_2 - \beta_2} \log(1+r).$$

Note that the treatment effect depends on the price elasticities of supply and demand. Therefore, we need to identify the demand and supply functions using instruments, since the observed quantities and prices cannot be used directly due to the issue of simultaneous causality. In particular, given the assumption of exogeneity of the other control variables, the only part of  $p_i^{Obs}$  that can correlate with  $u_{i1}$  or  $u_{i2}$  is the error component

$$\frac{u_{i2} - u_{i1}}{\beta_2 - \alpha_2}.$$

Hence

$$\mathbb{E}[u_{i1} p_i^{Obs}] = \frac{1}{\beta_2 - \alpha_2} \mathbb{E}[u_{i1}(u_{i2} - u_{i1})] = \frac{\text{Cov}(u_{i1}, u_{i2}) - \text{Var}(u_{i1})}{\beta_2 - \alpha_2} \neq 0,$$

and

$$\mathbb{E}[u_{i2} p_i^{Obs}] = \frac{1}{\beta_2 - \alpha_2} \mathbb{E}[u_{i2}(u_{i2} - u_{i1})] = \frac{\text{Var}(u_{i2}) - \text{Cov}(u_{i1}, u_{i2})}{\beta_2 - \alpha_2} \neq 0,$$

Due to the usual slopes  $\beta_2 < 0 < \alpha_2$  so that  $\alpha_2 - \beta_2 > 0$ .

We can use supply shifters to identify the demand, and demand shifters to identify the supply. In this setting, the exchange rate serves as an instrument to identify the demand, while income together with the prices of complementary and substitute goods serve as instruments to identify the supply. Thus, we can use the moment conditions,

$$\mathbb{E} \left[ \begin{bmatrix} 1 \\ y_i \\ pc_i \\ ps_i \\ er_i \end{bmatrix} (q_i^{Obs} - \beta_1 - \beta_2 p_i^{Obs} - \beta_3 y_i - \beta_4 pc_i - \beta_5 ps_i) \right] = \mathbf{0},$$

and

$$\mathbb{E} \left[ \begin{bmatrix} 1 \\ y_i \\ pc_i \\ ps_i \\ er_i \end{bmatrix} (q_i^{Obs} - \alpha_1 - \alpha_2 p_i^{Obs} - \alpha_3 er_i) \right] = \mathbf{0}$$

to identify the demand and supply functions. Note that the demand equation is *exactly identified*, whereas the supply equation is *over identified*.

We perform a simulation exercise to analyze the hypothetical causal effects of a new tax rate of 10% simulating the demand and supply equations using the following structural parameters  $\beta = [5 \ -0.5 \ 0.8 \ -0.4 \ 0.7]^\top$ ,  $\alpha = [-2 \ 0.5 \ -0.4]^\top$ ,  $u_1 \sim N(0, 0.5^2)$ , and  $u_2 \sim N(0, 0.5^2)$  assuming a sample size equal to 5,000. Additionally, assume that  $y \sim N(10, 1)$ ,  $pc \sim N(5, 1)$ ,  $ps \sim N(5, 1)$ , and  $er \sim N(15, 1)$ .

The population causal effect of the tax on quantity is

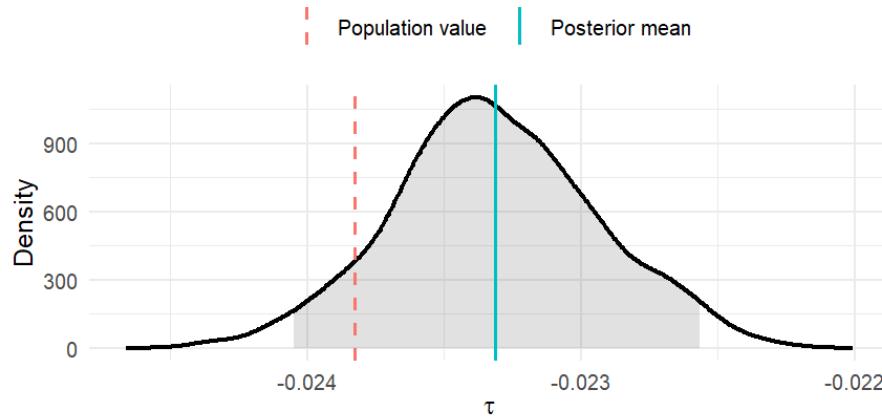
$$\tau = \frac{(-0.5) \times 0.5}{0.5 - (-0.5)} \log(1 + 0.1) \approx -0.0238,$$

which implies that a 10% tax reduces traded quantity by approximately 2.4%.

In Exercise 11, you are asked to program a BETEL algorithm from scratch to perform inference in this example. Figure 13.25 displays the posterior distribution of the causal effect. The 95% credible interval contains the population value, and the posterior mean lies close to it.

### Posterior of the causal effect

95% Crl: [-0.024, -0.023]; mean = -0.023



**FIGURE 13.25**

Posterior distribution of the causal effect of a tax: BETEL

[98] propose a unified framework based on marginal likelihoods and Bayes factors for comparing different moment-restricted models and for discarding misspecified restrictions. They demonstrate the model selection consistency of the marginal likelihood, showing that it favors the specification with the fewest parameters and the largest number of valid moment restrictions. When models are misspecified, the marginal likelihood procedure selects the model that is closest to the (unknown) true data-generating process in terms of Kullback–Leibler divergence. See [259, 277] for comprehensive reviews of recent advances in empirical likelihood and exponentially tilted empirical likelihood methods, including their Bayesian variants.

---

### 13.9 General Bayes posteriors

It is well known that, under model misspecification, that is, when the assumed likelihood function does not coincide with the true data-generating likelihood, the posterior concentrates its mass near those points in the support of the prior that minimize the Kullback–Leibler divergence with respect to the true model [245]. However, this updating process may yield credible sets that fail to achieve the desired coverage, Bayes factors that can be misleading, and predictions that may remain approximately valid but require careful checking. In addition, the misspecified posterior distribution may exhibit suboptimal risk performance. In this setting, a modified posterior directly linked to the risk function of interest, known as the *Gibbs posterior*, can still perform very well [229],

$$\hat{\pi}(\boldsymbol{\theta} | \mathbf{y}) = \frac{\exp\{-Nw R_N(\boldsymbol{\theta}, \mathbf{y})\} \pi(\boldsymbol{\theta})}{\int_{\Theta} \exp\{-Nw R_N(\boldsymbol{\theta}, \mathbf{y})\} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}}, \quad (13.5)$$

where  $w > 0$  is the *learning rate* (or the inverse temperature in simulated annealing), which balances the information in the data with that in the prior,  $N$  is the sample size, and

$$R_N(\boldsymbol{\theta}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N l(\boldsymbol{\theta}, \mathbf{y}_i)$$

is the empirical risk associated with the loss function  $l(\boldsymbol{\theta}, \mathbf{y}_i)$ .

Note that Equation 13.5 reduces to the ordinary Bayesian posterior when the loss is chosen as the negative log-likelihood,  $l(\boldsymbol{\theta}, y_i) = -\log p(y_i | \boldsymbol{\theta})$ , under i.i.d. sampling and with  $w = 1$ . In this case, we are asserting knowledge of the data-generating process  $p(y_i | \boldsymbol{\theta})$ . More generally, however,  $R_N(\boldsymbol{\theta}, \mathbf{y})$  can be based on any loss function that satisfies the required conditions for coherent belief updating: non-negativity, existence of a well-defined expectation, identifiability, and additivity across observations. Importantly, correctness of

the parametric model is not required, since

$$\frac{1}{N} \sum_{i=1}^N l(\boldsymbol{\theta}, \mathbf{y}_i) \xrightarrow{p} \int_{\mathcal{Y}} l(\boldsymbol{\theta}, \mathbf{y}) p(\mathbf{y} \mid \boldsymbol{\theta}) d\mathbf{y},$$

by the law of large numbers as  $N \rightarrow \infty$ .

Thus, the Gibbs posterior provides inference on the parameter values that minimize the chosen risk function, with minimal modeling assumptions. In contrast, the standard Bayesian posterior allows inference on essentially any feature of the data-generating distribution, but at the cost of strong model assumptions [427].

[46] show that Equation 13.5 is a valid, coherent mechanism to update prior beliefs. In particular, there must exist a mapping  $h$  such that

$$\pi(\boldsymbol{\theta} \mid \mathbf{y}) = h(l(\boldsymbol{\theta}, \mathbf{y}), \pi(\boldsymbol{\theta})),$$

where  $h$  satisfies the coherence property

$$h[l(\boldsymbol{\theta}, y_2), h(l(\boldsymbol{\theta}, y_1), \pi(\boldsymbol{\theta}))] = h(l(\boldsymbol{\theta}, y_2) + l(\boldsymbol{\theta}, y_1), \pi(\boldsymbol{\theta})).$$

This ensures that the updated posterior  $\pi(\boldsymbol{\theta} \mid y_1, y_2)$  is the same whether we update with  $(y_1, y_2)$  jointly or sequentially.

Equation 13.5 is the solution of minimizing the loss function  $L(\nu; \pi, \mathbf{y})$  on the space of probability measures on  $\boldsymbol{\theta}$  – space,

$$\hat{\pi} = \arg \min_{\nu} L(\nu; \pi, \mathbf{y}),$$

such that  $\hat{\pi}$  is the representation of beliefs about  $\boldsymbol{\theta}$  given prior beliefs ( $\pi$ ) and data ( $\mathbf{y}$ ). Given that the prior beliefs and the data are two independent pieces of information, it makes sense that the loss function is additive in these two arguments,

$$L(\nu; \pi, \mathbf{y}) = h_1(\nu, \mathbf{y}) + w^{-1} h_2(\nu, \pi),$$

where  $h_1(\cdot)$  and  $h_2(\cdot)$  are loss functions in their arguments. The coherence requirement implies that

$$h_2(\nu, \pi) = d_{KL}(\nu, \pi) = \int \log \frac{\nu(d\boldsymbol{\theta})}{\pi(d\boldsymbol{\theta})} \nu(d\boldsymbol{\theta}),$$

that is,  $h_2(\nu, \pi)$  must be the Kullback-Leibler divergence [46].

On the other hand,

$$h_1(\nu, \mathbf{y}) = \int l(\boldsymbol{\theta}, \mathbf{y}) \nu(d\boldsymbol{\theta})$$

is the expected loss of the action with respect to the data.

Therefore, the minimizer of

$$L(\nu; \pi, \mathbf{y}) = \int l(\boldsymbol{\theta}, \mathbf{y}) \nu(d\boldsymbol{\theta}) + d_{KL}(\nu, \pi),$$

is given by Equation 13.5 [469, 229, 46].

A very important parameter in the Gibbs posterior [229], or *general Bayes posterior* [46], is the learning rate, as it determines the asymptotic sampling properties of  $\hat{\pi}(\boldsymbol{\theta} | \mathbf{y})$  used to perform inference on  $\boldsymbol{\theta}$ . For instance, [46] propose different strategies to set this parameter, and [427] propose a Monte Carlo algorithm that selects the learning rate so that the resulting credible region attains the nominal Frequentist coverage probability.

[85] introduce the *Laplace-type estimator* (LTE) or *quasi-posterior* distribution, which can be interpreted as a special case of the *general Bayes* update of [46], taking as loss a scaled sample criterion based on the moment conditions (e.g., the GMM quadratic form).

In particular, given the moment conditions

$$\mathbb{E}[\mathbf{g}_i(\mathbf{w}_i, \boldsymbol{\theta})] = \mathbf{0}_d \quad \text{if and only if } \boldsymbol{\theta} = \boldsymbol{\theta}_0,$$

where the expectation is taken with respect to the population distribution,  $\mathbf{w}_{1:N} := [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_N]$  is a random sample from  $\mathbf{W} \subset \mathbb{R}^{d_w}$ ,  $\mathbf{g} : \mathbb{R}^{d_w} \times \Theta \rightarrow \mathbb{R}^d$  is a vector of known functions, and  $\boldsymbol{\theta} = [\theta_1 \ \theta_2 \ \dots \ \theta_p]^\top \in \Theta \subset \mathbb{R}^p$  with  $d \geq p$ , the risk function can be defined as

$$R_N(\boldsymbol{\theta}, \mathbf{w}) = \frac{1}{2} \left( \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{g}_i(\mathbf{w}_i, \boldsymbol{\theta})}_{\mathbf{g}_N(\boldsymbol{\theta})} \right)^\top \mathbf{W}_N \left( \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{g}_i(\mathbf{w}_i, \boldsymbol{\theta})}_{\mathbf{g}_N(\boldsymbol{\theta})} \right),$$

where  $\mathbf{W}_N$  is a positive semi-definite weighting matrix such that

$$\mathbf{W}_N \rightarrow \left( \text{Var} \left[ \sqrt{N} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i(\mathbf{w}_i, \boldsymbol{\theta}_0) \right) \right] \right)^{-1} \quad \text{as } N \rightarrow \infty.$$

Then, the *quasi-posterior* in [85] is similar to Equation 13.5 with  $w = 1$ .

Algorithm A36 shows the Metropolis-Hastings to perform inference using the general Bayes posterior.

Under suitable regularity conditions, [85] show that the LTE posterior mean is first-order equivalent to the efficient GMM estimator and that posterior quantiles yield confidence sets with asymptotically correct Frequentist coverage. Relatedly, [315] demonstrates that, under model misspecification, replacing the original likelihood with a curvature-adjusted (“sandwich”) log-likelihood can improve Frequentist risk, that is, lower the average loss over repeated samples when making inference on the parameters. Moreover, [79]

**Algorithm A36** General Bayes posterior: Metropolis-Hastings algorithm

---

```

1: Set  $\boldsymbol{\theta}^{(0)}$  in the support of  $\hat{\pi}(\boldsymbol{\theta} \mid \mathbf{y})$ 
2: for  $s = 1, \dots, S$  do
3:   Draw  $\boldsymbol{\theta}^c$  from  $q(\boldsymbol{\theta}^c \mid \boldsymbol{\theta}^{(s-1)})$ 
4:   Calculate  $\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c) = \min \left\{ \frac{q(\boldsymbol{\theta}^{(s-1)} | \boldsymbol{\theta}^c) \exp\{-Nw R_N(\boldsymbol{\theta}^c, \mathbf{y})\} \pi(\boldsymbol{\theta}^c)}{q(\boldsymbol{\theta}^c | \boldsymbol{\theta}^{(s-1)}) \exp\{-Nw R_N(\boldsymbol{\theta}^{(s-1)}, \mathbf{y})\} \pi(\boldsymbol{\theta}^{(s-1)})}, 1 \right\}$ 
5:   Draw  $U$  from  $U(0, 1)$ 
6:    $\boldsymbol{\theta}^{(s)} = \begin{cases} \boldsymbol{\theta}^c & \text{if } U \leq \alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^c) \\ \boldsymbol{\theta}^{(s-1)} & \text{otherwise} \end{cases}$ 
7: end for

```

---

develop confidence sets based on quasi-posteriors that achieve exact asymptotic Frequentist coverage for identified sets of parameters in complex nonlinear structural models, regardless of whether the parameters are point identified. Finally, [12] study decision rules for weak GMM and provide support for quasi-Bayesian procedures built from the GMM quadratic form, in line with the spirit of the LTE in settings with weak identification.

**Example: Instrumental variable quantile regression (IVQR)**

[83] propose an instrumental variable model of quantile treatment effects to characterize the heterogeneous impact of treatments across different points of the outcome distribution. Their model requires conditions that restrict the evolution of ranks across treatment states. Using these conditions, they address the endogeneity problem and recover quantile treatment effects using instrumental variables for the entire population, not only for compliers.

Assume a binary treatment variable  $D_i \in \{0, 1\}$ , with potential outcomes  $\{Y_i(0), Y_i(1)\}$ . The estimand of interest is the quantile treatment effect, which summarizes the differences in the quantiles of potential outcomes across treatment states,

$$q(D_i = 1, \mathbf{X}_i = \mathbf{x}_i, \tau) - q(D_i = 0, \mathbf{X}_i = \mathbf{x}_i, \tau),$$

where  $q(D_i = d, \mathbf{X}_i = \mathbf{x}_i, \tau)$  denotes the  $\tau$ -th quantile treatment response function.

The potential outcomes are related to the quantile treatment response via

$$Y_i(d) = q(D_i = d, \mathbf{X}_i = \mathbf{x}_i, U(d)),$$

where  $U(d_i) \sim U(0, 1)$  is the *rank variable*. This variable captures unobserved heterogeneity that explains differences in outcomes given observed characteristics  $\mathbf{x}_i$  and treatment  $d_i$ .

For instance, consider a retirement savings model, where the potential outcome is net financial assets under different retirement plan statuses  $d$ , and  $q(D = d, \mathbf{X} = \mathbf{x}, \tau)$  is the net financial asset function describing how an individual with retirement status  $d$  and “financial ability”  $\tau$  is rewarded in

the financial market. Because the function depends on  $\tau$ , treatment effects are heterogeneous.

The identification conditions for the IVQR model are stated in [83] as follows:

1. *Potential outcomes:*

$$Y_i(d) = q(D_i = d_i, \mathbf{X}_i = \mathbf{x}_i, U(d_i)),$$

where  $q(D_i = d_i, \mathbf{X}_i = \mathbf{x}_i, \tau)$  is strictly increasing in  $\tau$  and  $U(d_i) \sim U(0, 1)$ .

2. *Independence:* Conditional on  $\mathbf{X}_i = \mathbf{x}_i$ , the rank variables  $\{U(d_i)\}$  are independent of the instruments  $\mathbf{Z}_i$ .

3. *Selection:* The treatment assignment is given by  $D_i = \delta(\mathbf{Z}_i, \mathbf{X}_i, \mathbf{V}_i)$  for some unknown function  $\delta$  and unobserved heterogeneity  $\mathbf{V}_i$ .

4. *Rank invariance:* Conditional on  $\mathbf{X}_i = \mathbf{x}_i, \mathbf{Z}_i = \mathbf{z}_i$ ,

- either  $\{U(d_i)\}$  coincide ( $U(d_i) = U$  for all  $d$ ), or
- $\{U(d_i)\}$  are identically distributed conditional on  $\mathbf{V}_i$ .

5. *Observables:* The observed data consist of  $Y_i = q(D_i, \mathbf{X}_i, U(D_i))$ ,  $D_i$ ,  $\mathbf{X}_i$ , and  $\mathbf{Z}_i$ ,  $i = 1, 2, \dots, N$ .

The most important identification restriction is *rank invariance*, which implies that individuals with higher unobserved rank remain higher-ranked regardless of treatment status. This condition accommodates more general selection mechanisms than the monotonicity assumption used in the LATE framework, while being less restrictive than full independence assumptions between instruments ( $\mathbf{Z}$ ) and unobserved variables in the selection equation ( $\mathbf{V}$ ) in other common models (see [82, 83] for details).<sup>7</sup>

The main testable implication of the identification restrictions is that

$$P(Y_i \leq q(D_i, \mathbf{X}_i, \tau) \mid \mathbf{X}_i, \mathbf{Z}_i) = P(Y_i < q(D_i, \mathbf{X}_i, \tau) \mid \mathbf{X}_i, \mathbf{Z}_i) = \tau,$$

for all  $\tau$  almost surely, and  $U(D_i) \sim U(0, 1)$  conditional on  $(\mathbf{X}_i, \mathbf{Z}_i)$ .

This conditional moment restriction implies the unconditional moment conditions that form the basis for estimation and inference in the IVQR model [85, 82],

$$\mathbf{g}_N(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\tau - \mathbb{1}(Y_i \leq \alpha_\tau D_i + \mathbf{X}_i^\top \boldsymbol{\beta}_\tau)) \begin{bmatrix} \mathbf{X}_i \\ \mathbf{Z}_i \end{bmatrix}.$$

Thus, the empirical risk function is given by

$$R_N(\boldsymbol{\theta}, \mathbf{w}) = \frac{1}{2} \mathbf{g}_N(\boldsymbol{\theta})^\top \mathbf{W}_N \mathbf{g}_N(\boldsymbol{\theta}),$$

where

$$\mathbf{W}_N = \frac{1}{\tau(1-\tau)} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{Z}_i \mathbf{Z}_i^\top \right)^{-1}.$$

The following code implements Algorithm A36 to perform inference on the

---

<sup>7</sup>[83]'s proposal allows for both discrete and continuous  $D$  and  $\mathbf{Z}$ .

quantile treatment effect of participation in the 401(k) retirement program on net financial assets, using eligibility as an instrument (see the 401(k) treatment effects example in Section 13.4 for details of the data).

### R code. Quantile treatment effects: 401(k) participation on net financial assets

```

1 rm(list = ls()); set.seed(10101)
2 # Load data
3 df <- read.csv("https://raw.githubusercontent.com/BESmarter-
  consultancy/BSTApp/refs/heads/master/DataApp/401k.csv",
4 sep = ",", header = TRUE, quote = "")
5 # Attach variables
6 attach(df)
7 y <- net_tfa/1000    # Outcome: net financial assets
8 x <- as.vector(p401) # Endogenous regressor: participation
9 w <- as.matrix(cbind(age, inc, fsize, educ, marr, twoearn,
  db, pira, hown))   # Exogenous regressors
10 z <- as.matrix(e401) # Instrument: eligibility (NO
  intercept here)
11 library(quantreg)
12 tau <- 0.5
13 QuanReg <- rq(y ~ p401 + age + inc + fsize + educ + marr +
  twoearn + db + pira + hown, tau = tau, data = df)
14 summary(QuanReg)
15 Reg <- MCMCpack::MCMCquantreg(y ~ x + w, data = df, tau =
  tau)
16 BayesExo <- summary(Reg)
17 LossFunct <- function(par, tau, y, z, x, w){
18   n <- length(y)
19   X <- cbind(1, x, w)
20   Z <- cbind(1, z, w)
21   Ind <- as.numeric(y <= X%*%par)
22   gn <- colMeans((tau - Ind) * Z)
23   Wni <- lapply(1:n, function(i) {Z[i,] %*% t(Z[i,])})
24   Wn <- 1 / (tau * (1 - tau)) * solve(Reduce("+", Wni)/n)
25   Ln <- - 0.5 * n * t(gn) %*% Wn %*% gn
26   return(Ln)
27 }
28 par0 <- colMeans(Reg)
29 LossFunct(par = par0, tau = tau, y = y, z = z, x = x, w = w)
30 # ----- MH using Ln -----
31 k <- length(par0); b0 <- rep(0, k); B0 <- 1000*diag(k)
32 S <- 10000; burnin <- 10000; thin <- 1; tot <- S + burnin
33 BETA <- matrix(NA, tot, k); accept <- logical(tot)
34 SIGMA <- diag(BayesExo[["statistics"]][,2]); tune <- 2.4 /
  sqrt(k)

```

***R code. Quantile treatment effects: 401(k) participation on net financial assets***

---

```

1 BETA[1,] <- par0
2 LL <- LossFunct(par = BETA[1,], tau = tau, y = y, z = z, x =
   x, w = w)
3 pb <- txtProgressBar(min=0, max=tot, style=3)
4 for(s in 2:tot){
5   cand <- BETA[s-1,] + MASS::mvrnorm(1, rep(0, k), tune *
   SIGMA)
6   LLc <- LossFunct(par = cand, tau = tau, y = y, z = z, x =
   x, w = w)
7   priorRat <- mvtnorm::dmvnorm(cand, b0, B0, log=TRUE) -
   mvtnorm::dmvnorm(BETA[s-1,], b0, B0, log=TRUE)
8   loga <- (LLc - LL) + priorRat
9   if (is.finite(loga) && log(runif(1)) <= loga) {
10     BETA[s,] <- cand; LL <- LLc; accept[s] <- TRUE
11   } else {
12     BETA[s,] <- BETA[s-1,]; accept[s] <- FALSE
13   }
14   if (s <= burnin && s %% 100 == 0) {
15     acc <- mean(accept[(s-99):s])
16     if (acc > 0.35) tune <- tune * 1.25
17     if (acc < 0.15) tune <- tune / 1.25
18   }
19   setTxtProgressBar(pb, s)
20 }
21 close(pb)
22 cat("\nAcceptance rate:", mean(accept), "\n")
23 keep <- seq(burnin, tot, by = thin)
24 cat("\nAcceptance rate:", mean(accept[keep]), "\n")
25 post <- BETA[keep, , drop=FALSE] # posterior draws in
   scaled space
26 colnames(post) <- c("Int", "p401", "age", "inc", "fsiz", "
   educ", "marr", "twoearn", "db", "pira", "hown")
27 # Posterior summaries
28 post_mean <- colMeans(post)
29 post_ci <- apply(post, 2, quantile, c(0.025, 0.975))
30 round(post_mean, 3); round(post_ci, 3)

```

---

Setting  $\tau = 0.5$  (the median), the quantile treatment effect is USD 6,719, with a 95% credible interval of (USD 6,612, USD 6,806). For  $\tau = 0.9$ , the treatment effect is USD 19,318, with a 95% credible interval of (USD 18,634, USD 20,085). Note that, in the instrumental variable example, the posterior mean of the LATE was USD 8,520, which is higher than the treatment effect at the median ( $\tau = 0.5$ ). This suggests that the LATE overstates the causal effect of 401(k) participation due to the large impact at higher quantiles of the outcome distribution.

### 13.10 Doubly robust Bayesian inferential framework (DRB)

[62] propose a doubly robust Bayesian (DRB) approach to perform inference on the average treatment effect (ATE) under unconfoundedness. They adjust the prior for the conditional mean outcome function using an estimator of the propensity score, the conditional probability of treatment assignment [376], and then re-center the posterior distribution of the treatment effect using the propensity score estimator together with a pilot estimator of the outcome regression. The authors show that under unconfoundedness and overlap, the corrected posterior converges in distribution to a normal law centered at a semiparametrically efficient estimator, and that the resulting credible intervals achieve asymptotically exact Frequentist coverage.

Let  $Y_i \in \{0, 1\}$  be the outcome and  $D_i \in \{0, 1\}$  the treatment,<sup>8</sup> thus

$$Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0),$$

where  $Y_i(1)$  and  $Y_i(0)$  are the potential outcomes.

Let  $\mathbf{X}_i$  denote pre-treatment covariates with distribution  $F_0$  and density  $f_0$ , and define the propensity score

$$\pi_0(\mathbf{x}) := P(D_i = 1 | \mathbf{X}_i = \mathbf{x})$$

and the conditional mean outcome (success probability)

$$m_0(d, \mathbf{x}) := P(Y_i = 1 | D_i = d, \mathbf{X}_i = \mathbf{x}).$$

For an i.i.d. sample  $W_i = (Y_i, D_i, \mathbf{X}_i^\top)^\top$ , the joint density can be written as

$$p_{f,\pi,m}(y, d, \mathbf{x}) = f(\mathbf{x}) [\pi(\mathbf{x})]^d [1 - \pi(\mathbf{x})]^{1-d} [m(d, \mathbf{x})]^y [1 - m(d, \mathbf{x})]^{1-y}.$$

The parameter of interest is the average treatment effect (ATE),

$$\tau := \mathbb{E}_0[Y_i(1) - Y_i(0)],$$

where  $\mathbb{E}_0$  is the expectation with respect to the population distribution.

[62] impose support restrictions via link functions. For the outcome regression and propensity score use logits:

$$m_\eta(d, \mathbf{x}) = \frac{1}{1 + \exp\{-\eta^m(d, \mathbf{x})\}}, \quad \pi_\eta(\mathbf{x}) = \frac{1}{1 + \exp\{-\eta^\pi(\mathbf{x})\}},$$

so that

$$\eta^m(d, \mathbf{x}) = \log\left(\frac{m_\eta(d, \mathbf{x})}{1 - m_\eta(d, \mathbf{x})}\right), \quad \eta^\pi(\mathbf{x}) = \log\left(\frac{\pi_\eta(\mathbf{x})}{1 - \pi_\eta(\mathbf{x})}\right).$$

---

<sup>8</sup>[62] extend their framework to continuous, counting and multinomial outcomes.

For the covariate density, they use a log-density parametrization

$$f_\eta(\mathbf{x}) = \exp\{\eta^f(\mathbf{x})\}, \quad \eta^f(\mathbf{x}) = \log f_\eta(\mathbf{x}).$$

Given a prior on  $(\eta^m, \eta^\pi, \eta^f)$ , the induced ATE is

$$\tau_\eta := \mathbb{E}_\eta[m_\eta(1, \mathbf{X}) - m_\eta(0, \mathbf{X})].$$

The basis for inference on  $\tau_\eta$  is the *efficient influence function* (EIF):

$$\tilde{\tau}_\eta(W) = \{m_\eta(1, \mathbf{X}) - m_\eta(0, \mathbf{X})\} + \gamma_\eta(D, \mathbf{X}) \{Y - m_\eta(D, \mathbf{X})\} - \tau_\eta,$$

where the Riesz representer for the ATE functional is

$$\gamma_\eta(D, \mathbf{X}) = \frac{D}{\pi_\eta(\mathbf{X})} - \frac{1-D}{1-\pi_\eta(\mathbf{X})}.$$

An influence function measures the first-order effect of an infinitesimal contamination of the underlying distribution on the estimand. The *efficient* influence function attains the semiparametric efficiency bound (its variance is the minimal asymptotic variance in the model). At the truth  $\eta_0$ , the EIF is mean-zero:  $\mathbb{E}_0[\tilde{\tau}_{\eta_0}(W)] = 0$ . Moreover, it is *doubly robust*:

$$\mathbb{E}_0[\tilde{\tau}_\eta(W)] = 0 \quad \text{if either } m_\eta \text{ or } \pi_\eta \text{ is correctly specified.}$$

To conduct Bayesian inference, [62] place independent priors on the outcome regression and the covariate distribution (equivalently, on  $\eta^m$  and  $\eta^f$ ). For the covariates they use a Dirichlet Process (DP) prior on  $F$  (hence on  $f$ ; see Section 11.1.2), which, under standard conditions, induces posterior draws that approximate the Bayesian bootstrap for sample supported functionals [279] (see Section 6.10).

For the outcome regression they place a Gaussian Process (GP) prior on  $\eta^m$  (see Section 12.4), with a correction term that incorporates a preliminary estimator of the Riesz representer. The propensity score is estimated in a separate first step and is used in the correction/recentering; thus, a prior for  $\eta^\pi$  is not required for their procedure.

Algorithm A37 outlines the doubly robust Bayesian procedure of [62]. The pilot estimate of the Riesz representer,  $\hat{\gamma}_\eta$ , is obtained by plug-in, using a propensity score  $\hat{\pi}(\mathbf{x})$  estimated via logistic Lasso.

The pilot estimate of the outcome regression,  $\hat{m}_\eta(d, \mathbf{x}) = \Pr(Y = 1 \mid D = d, \mathbf{X} = \mathbf{x})$ , is taken as the posterior mean from a GP *classification* model fit with a Laplace approximation.

Concretely, place a GP prior on the latent logit  $\eta^m(d, \mathbf{x})$  with a squared exponential kernel, and map to success probabilities via the logistic link  $m(d, \mathbf{x}) = \text{logit}^{-1}\{\eta^m(d, \mathbf{x})\}$ . Because the Bernoulli likelihood makes the required integrals over  $\eta^m$  analytically intractable, the posterior is approximated using Laplace's method (Expectation Propagation is a common alternative). See Sections 3.3–3.5 of [364] for details.

The *adjusted* GP prior specifies the latent logit as

$$\eta^m(d, \mathbf{x}_i) = W^m(d, \mathbf{x}_i) + \lambda \hat{\gamma}_\eta(d, \mathbf{x}_i), \quad W^m \sim \text{GP}(0, K), \quad \lambda \sim N(0, \sigma_n^2),$$

where  $K$  is the base kernel. Following [62], take

$$\sigma_n = \frac{\log n}{\sqrt{n} \Gamma_n}, \quad \Gamma_n := \frac{1}{n} \sum_{i=1}^n |\hat{\gamma}_\eta(d_i, \mathbf{x}_i)|,$$

where  $n$  is the effective sample size (after trimming).

Since  $W^m$  and  $\lambda$  are Gaussian, this is equivalent to a GP with an augmented kernel,

$$\eta^m \sim \text{GP}(0, K_c), \quad K_c((d, \mathbf{x}), (d', \mathbf{x}')) = K((d, \mathbf{x}), (d', \mathbf{x}')) + \sigma_n^2 \hat{\gamma}_\eta(d, \mathbf{x}) \hat{\gamma}_\eta(d', \mathbf{x}'),$$

which, for the vector of latent logits evaluated at the training inputs, corresponds to the covariance matrix  $\mathbf{K} + \sigma_n^2 \hat{\gamma} \hat{\gamma}^\top$ .

---

**Algorithm A37** Double Robust Bayesian Algorithm

---

- 1: Have initial estimates of  $\hat{\gamma}_\eta$  based on logistic Lasso
  - 2: Trim the propensity score using standard values, for instance, 10%, or the optimal trimming rule proposed by [108]
  - 3: Estimate the Riesz representer using the estimate of the propensity score
  - 4: Calculate the prior correction term  $\sigma_n = \frac{\log(n)}{n^{1/2} \Gamma_n}$ ,  $\Gamma_n = \frac{1}{n} \sum_{i=1}^n |\hat{\gamma}_\eta(d_i, \mathbf{x}_i)|$ ,  $n$  is the effective sample size (after trimming)
  - 5: Calculate the prior adjustment  $\sigma_n \times \hat{\gamma}_\eta$
  - 6: Calculate  $\hat{m}(d, \mathbf{x}) = \frac{1}{1+\exp\{-\eta^m(d, \mathbf{x})\}}$  as the posterior mean of a Gaussian Process (GP) based on  $d$  and  $\mathbf{x}$  (unadjusted mean)
  - 7: Set an initial GP prior for the adjusted mean  $m_\eta^c(d, \mathbf{x}) = \frac{1}{1+\exp\{-\eta^m\}}$ , where  $\eta^m \sim \text{GP}(0, K_c)$ ,  $K_c((d, \mathbf{x}), (d', \mathbf{x}')) = K((d, \mathbf{x}), (d', \mathbf{x}')) + \sigma_n^2 \hat{\gamma}_\eta(d, \mathbf{x}) \hat{\gamma}_\eta(d', \mathbf{x}')$
  - 8: **for**  $s = 1, \dots, S$  **do**
  - 9:     Generate the  $s$ -th draw of the adjusted GP posterior,  $m_\eta^{c,s}(d, \mathbf{x})$
  - 10:    Draw Bayesian bootstrap weights  $M_{ni}^s = \frac{e_i^s}{\sum_{j=1}^n e_j^s}$ ,  $e_i^s \sim \text{Exp}(1)$ ,  $i = 1, 2, \dots, n$
  - 11:    Calculate the corrected posterior draw for the ATE  $\tilde{\tau}_\eta^s = \tau_\eta^s - \hat{b}_\eta^s$ , where  $\tau_\eta^s = \sum_{i=1}^n M_{ni}^s (m_\eta^{c,s}(1, \mathbf{x}_i) - m_\eta^{c,s}(0, \mathbf{x}_i))$ , and  $\hat{b}_\eta^s = \frac{1}{n} \sum_{i=1}^n \tau[m_\eta^s(\mathbf{w}_i) - \hat{m}(\mathbf{w}_i)]$ ,  $\tau[m] = m(1, \mathbf{x}) - m(0, \mathbf{x}) + \hat{\gamma}(d, \mathbf{x})(y - m(d, \mathbf{x}))$
  - 12: **end for**
- 

Algorithm A37 produces posterior draws of the ATE. The  $(1-\alpha)\%$  credible interval is given by the empirical quantiles at levels  $\alpha/2$  and  $1 - \alpha/2$ .

It is worth noting that [62] do not perform sample splitting, whereas [81] recommend it in the context of *Double/Debiased Machine Learning for treatment and structural parameters*. While sample splitting (or cross-fitting) helps

reduce overfitting, its role in [81] is more fundamental: it ensures that the nuisance functions are estimated on data independent of the observations used for inference, which is critical for their theoretical guarantees of  $\sqrt{N}$ -consistency and valid asymptotic inference. Nevertheless, [62] report that they did not find evidence of overfitting in their setting and verified that their results are robust to the use of sample splitting.

Finally, observe that the doubly robust approach does not propagate the uncertainty arising from the estimation of the propensity score, which is treated as a plug-in. This omission is justified because the EIF is *Neyman orthogonal*: estimation error in the propensity score enters only at second order and is therefore asymptotically negligible for uncertainty quantification. [278] propose a quasi-Bayesian framework in which a *control function*, explicitly modeling the source of endogeneity, is estimated in a first stage using Frequentist methods and then used in a second stage where inference is conducted using Bayesian methods. The posterior mean converges asymptotically to the true parameter, and a bootstrap procedure is employed to account for uncertainty propagation from the first stage, thereby achieving valid asymptotic coverage.

#### **Example: Doubly robust Bayesian inference**

Let us simulate the process

$$\eta_i^\pi = -0.3 + 0.8X_{i1} - 0.5X_{i2} + 0.7X_{i3}, \quad \eta_i^m = -0.2 + 0.6D_i + 0.5X_{i1} - 0.4X_{i2} + 0.3X_{i3},$$

where  $X_1 \sim N(0, 1)$ ,  $X_2 \sim \text{Bernoulli}(0.5)$ , and  $X_3 \sim U(-1, 1)$ , and the sample size is 2,000.

The population ATE from this simulation is approximately 0.138. The posterior mean and 95% credible interval of the ATE using the DRB are 0.155 and (0.107, 0.206), respectively.

The following code shows how to implement the DRB inferential framework in this setting.

#### **R code. Doubly robust Bayesian inference: Simulation exercise**

```

1 ##### Doubly Robust Bayesian: Simulation #####
2 rm(list = ls()); set.seed(123); library(glmnet)
3 ## simulate covariates
4 n <- 2000
5 X1 <- rnorm(n) # continuous
6 X2 <- rbinom(n, 1, 0.5) # binary
7 X3 <- runif(n, -1, 1) # continuous
8 Z <- cbind(X1, X2, X3)

```

### R code. Doubly robust Bayesian inference: Simulation exercise

```

1 # True propensity and outcome
2 logit <- function(t) 1/(1+exp(-t))
3 true_pi <- function(X){ logit(-0.3 + 0.8*X[,1] -0.5*X[,2] +
0.7*X[,3]) }
4 true_eta <- function(d,X){ -0.2 + 0.6*d + 0.5*X[,1] -0.4*X
[,2] + 0.3*X[,3] }
5 true_p <- function(d,X){ logit(true_eta(d,X)) }
6 ATE <- mean(true_p(1,Z) - true_p(0,Z)) # Population ATE
7 # Data
8 pi <- true_pi(Z); D <- rbinom(n,1,pi)
9 P <- true_p(D,Z); Y <- rbinom(n,1,P)
10 dat <- data.frame(Y, D, X1, X2 = factor(X2), X3)
11 ##### Doubly Robust Bayesian: Implementation #####
12 p <- dim(Z)[2]; covars <- c("X1","X2","X3")
13 Z.df <- as.data.frame(Z) # covariates as a dataframe
14 # Propensity score model: lasso
15 cvfit <- cv.glmnet(Z, D, family = "binomial", alpha=1,
nfolds = 10) # lasso,min
16 ps.coef <- coef(cvfit, s = "lambda.min")[0:p+1]
17 ps_hat <- predict(cvfit, newx = Z, s = "lambda.min",type =
"response")
18 ps_hat <- pmin(pmax(ps_hat, 1e-6), 1 - 1e-6)
19 # Riesz representer
20 riesz_fun <- function(d, ps) d/ps - (1 - d)/(1 - ps)
21 # GP model
22 make_gp <- function(X_train, y_bin, use_correction = FALSE,
gamma_scaled = NULL,
23 approx_method = gplite::approx_laplace(),
24 init_lscale = 0.3) {
25 stopifnot(length(y_bin) == nrow(X_train))
26 # Initializing the covariance function
27 cf_se <- gplite::cf_sexp(vars = colnames(X_train), lscale
= init_lscale, magn = 1, normalize = TRUE)
28 cfs <- list(cf_se)
29 if (use_correction) {
30 stopifnot(!is.null(gamma_scaled))
31 # Add linear kernel on the last column (gamma_scaled)
32 cf_lin <- gplite::cf_lin(vars = "gamma", magn = 1,
normalize = FALSE)
33 cfs <- list(cf_se, cf_lin)
34 }
35 # Initializes a GP model
36 gp <- gplite::gp_init(cfs = cfs, lik = gplite::lik_
bernoulli(), approx = approx_method)
37 # Optimizing hyperparameters
38 gp <- gplite::gp_optim(gp, x = X_train, y = y_bin, maxiter
= 1000, restarts = 3, tol = 1e-05)
39 return(gp)
40 }
```

**R code. Doubly robust Bayesian inference:  
Simulation exercise**

---

```

1 posterior_draws_prob <- function(gp, X_test, ndraws = 5000)
2   {
3     out <- gplite::gp_draw(gp, xnew = X_test, draws = ndraws,
4       transform = TRUE, target = FALSE, jitter = 1e-6)
5     return(t(out))
6   }
7 # Main function
8 run_trim <- function(t_trim = 0.10, N_post = 5000, h_r = 1/
9   2, cor_size = 1) {
10   # Trim by estimated PS
11   keep <- which(ps_hat >= t_trim & ps_hat <= (1 - t_trim))
12   n_eff <- length(keep)
13   if (n_eff < 50) stop("Too few observations after trimming.
14   ")
15   Yk <- Y[keep]
16   Dk <- D[keep]
17   Zk <- Z[keep, , drop = FALSE]
18   psk <- ps_hat[keep]
19   # Riesz representer and sigma_n choice
20   gamma_k <- riesz_fun(Dk, psk)
21   sigma_n <- cor_size * (log(n_eff) / ( (n_eff)^h_r * mean
22   (abs(gamma_k)) ))
23   gamma_scaled <- sigma_n * gamma_k
24   # Training inputs for GP: [Z, D] and, if corrected, an
25   # extra column 'gamma'
26   X_train_nc <- data.frame(Zk)
27   X_train_nc$D <- Dk
28   colnames(X_train_nc) <- c(covars, "D")
29   X_train_c <- X_train_nc
30   X_train_c$gamma <- gamma_scaled
31   # Test inputs: for each i, (Z_i, d=0) and (Z_i, d=1)
32   X_t0 <- X_train_nc; X_t0$D <- 0
33   X_t1 <- X_train_nc; X_t1$D <- 1
34   X_t0c <- X_t0; X_t1c <- X_t1
35   X_t0c$gamma <- sigma_n * riesz_fun(0, psk)
36   X_t1c$gamma <- sigma_n * riesz_fun(1, psk)
37   # GP WITHOUT prior correction #
38   gp_nc <- make_gp(X_train = X_train_nc, y_bin = Yk, use_
39   correction = FALSE)
40   # joint draws for [m(Z,0); m(Z,1)]
41   M_nc_0 <- posterior_draws_prob(gp_nc, X_t0, ndraws = N_
42   post) # N_post x n_eff
43   M_nc_1 <- posterior_draws_prob(gp_nc, X_t1, ndraws = N_
44   post)
45   # observed arm probabilities
46   M_nc_obs <- ifelse(matrix(Dk, nrow = N_post, ncol = n_eff,
47   byrow = TRUE) == 1, M_nc_1, M_nc_0)

```

---

### R code. Doubly robust Bayesian inference: Simulation exercise

```

1  # GP WITH prior correction #
2  gp_c <- make_gp(X_train = X_train_c, y_bin = Yk, use_
+    correction = TRUE, gamma_scaled = gamma_scaled)
3  M_c_0 <- posterior_draws_prob(gp_c, X_t0c, ndraws = N_post
+    )
4  M_c_1 <- posterior_draws_prob(gp_c, X_t1c, ndraws = N_post
+    )
5  M_c_obs <- ifelse(matrix(Dk, nrow = N_post, ncol = n_eff,
+    byrow = TRUE) == 1, M_c_1, M_c_0)
6  # Bayesian bootstrap weights #
7  W <- matrix(rexp(N_post * n_eff, rate = 1), nrow = N_post)
8  W <- W / rowSums(W)
9  # ATEs
10 # (1) Unadjusted Bayes
11 Ate_nc_draws <- rowSums((M_c_1 - M_c_0) * W)
12 # (2) Prior-adjusted Bayes
13 Ate_c_draws <- rowSums((M_c_1 - M_c_0) * W)
14 # (3) DR Bayes (posterior recentering)
15 # Pre-centering using UNCorrected GP posterior means
16 mu_nc_diff <- colMeans(M_c_1 - M_c_0)
17 mu_nc_obs <- colMeans(M_c_obs)
18 ATE_dr_pre <- mean(mu_nc_diff + gamma_k * (Yk - mu_nc_obs)
+    )
19 # Recenter draws using corrected GP draws
20 DR_rec_1 <- (matrix(gamma_k, nrow = N_post, ncol = n_eff,
+    byrow = TRUE)) * (matrix(Yk, nrow = N_post, ncol = n_eff
+    , byrow = TRUE) - M_c_obs)
21 # The first component is \tau_{\eta}'s in Algorithm 1
22 # The second component is \hat{m} a scalar. This has
# positive sign because
23 # minus x minus, the bias is with minus, and then, this
# component enters with minus in the bias term
24 # The third component is \gamma_k x (y-m(d,x)), m(d,x) is
# with prior correction
25 # The fourth component is m(1,x)-m(0,x) also with
# correction in the prior
26 Ate_drb_draws <- rowSums((M_c_1 - M_c_0) * W) + ATE_dr_pre
- rowSums(DR_rec_1) / n_eff - rowSums(M_c_1 - M_c_0) /
n_eff
27 qfun <- function(x) {
28   qs <- quantile(x, c(0.025, 0.975))
29   c(mean = mean(x), median = median(x), lo = qs[1], hi
= qs[2], len = diff(qs))
30 }
31 vals <- list(Bayes = qfun(Ate_nc_draws), 'PA Bayes' = qfun
(Ate_c_draws), 'DR Bayes' = qfun(Ate_drb_draws))
32 out_df <- as.data.frame(do.call(rbind, vals), check.names
= FALSE)
33 list(t_trim = t_trim, results = out_df)
34 }
```

***R code. Doubly robust Bayesian inference:  
Simulation exercise***

---

```
1 N_post <- 5000; h_r <- 1/2
2 cor_size <- 1; trim <- 0.05
3 res_list <- run_trim(t_trim = trim, N_post = N_post, h_r = h
   _r, cor_size = cor_size)
4 print(res_list$results)
```

---

### 13.11 Summary

In this chapter, we review some of the most common approaches to causal inference. We begin by outlining the identification restrictions underlying each method and then describe how these restrictions are incorporated into a Bayesian framework for inference. A natural point of departure is the use of Directed Acyclic Graphs (DAGs), which provide a graphical representation of the underlying structural or causal model. The next step is to identify an exogenous source of variation in the assignment rule, such as instruments or institutional arrangements, whose variability enables the identification of the causal effect of the treatment or relevant regressor on the outcome.

Given such sources of variation, one can use parametric, semiparametric, or nonparametric models to specify the conditional means of the potential outcomes given treatment assignment, thereby constructing counterfactual scenarios and enabling inference on causal estimands, such as the average treatment effect. Particular attention must be paid to the correct specification of the treatment assignment mechanism (propensity score) and the outcome regression, as both play a crucial role in credible causal inference.

This highlights the growing importance of modern machine learning methods for estimating high-dimensional or complex nuisance functions (see Chapter 12 and Section 13.10). At the same time, however, it is essential to address and correct for the biases that machine learning methods may introduce in order to obtain valid causal inference.

Finally, the Bayesian framework provides a coherent way to integrate these identification strategies with prior information and to obtain full posterior distributions of causal estimands, based on both the marginal distributions of the potential outcomes and their joint distribution, thereby offering a unified approach to causal inference.

---

### 13.12 Exercises

1. Show that the Average Treatment Effect (ATE) in the simple linear regression framework

$$Y_i = \beta_0 + \tau D_i + \mu_i,$$

assuming non-informative prior distributions, so that the posterior mean of the location parameter coincides with the maximum likelihood estimator, is equal to

$$\bar{y}_1 - \bar{y}_0.$$

2. Some readers may question the assumption that potential outcomes are normally distributed. However, it is important to note that the normal distribution is the *maximum entropy* continuous distribution given a specified mean  $\mu$  and finite variance  $\sigma^2$ . In other words, among all distributions with the same mean and variance, the normal distribution represents the one with the greatest level of uncertainty or unpredictability [106].

Show that the normal distribution is the *maximum entropy* continuous distribution given a specified mean  $\mu$  and finite variance  $\sigma^2$  by considering the formal definition of entropy:

$$H(f) = - \int_{-\infty}^{\infty} f(y) \log f(y) dy,$$

where  $f(y)$  is a probability density function.

3. Use the package *dagitty* to construct the DAG in Figure 13.5, verify that it is acyclic, and check whether the causal effect of  $D$  on  $Y$  is identifiable by controlling for  $\mathbf{X}$ .
4. Use the package *dagitty* to construct the DAG in Figure 13.8, verify that it is acyclic, and check that the causal effect of  $D$  on  $Y$  is identifiable by controlling for  $\mathbf{X}$  but not for  $C$ .
5. Use the package *dagitty* to construct the DAG in Figure 13.11, taking into account that  $\mathbf{U}$  is unobserved (latent). Verify that it is acyclic, check that  $Z$  is a valid instrument, and determine whether the causal effect of  $D$  on  $Y$  is identifiable.
6. **401(k) participation on net financial assets continues I**  
Apply the framework from this example to compute the intention-to-treat effect, the local average treatment effect, and the effect of eligibility on participation.

**7. 401(k) participation on net financial assets continues II**

Perform inference in this example assuming that the stochastic errors follow a Dirichlet process, using the function *rivDP* from the *bayesm* package to analyze 401(k) participation and its effect on net financial assets under the same specification as in the main text, and plot the LATE.

**8. Difference-in-Differences simulation continues I**

Perform the simulation of the DiD example, and perform inference using the specification:

$$Y_{it} = \alpha + \alpha_i + \phi_t + \tau_2 [D_{it} \cdot \mathbb{1}(t = 2)] + \epsilon_{it},$$

**9. Difference-in-Differences simulation continues II**

Note that another strategy to perform inference on the ATT is to estimate the saturated model

$$Y_{it} = \sum_{t,l} \mu_{tl} [D_{il} \cdot \mathbb{1}(t = t)] + \epsilon_{it}, \quad t = 1, 2, \quad l = 1, 0,$$

and then use the posterior draws to compute

$$\tau_2 = (\mu_{21} - \mu_{11}) - (\mu_{20} - \mu_{10}).$$

Explain why inference on  $\tau_2$  using this approach in the simulation setting shows that the posterior mean is similar to that from the previous approaches, but the level of uncertainty is higher.

10. Perform a simulation exercise to assess the ability of BETEL to identify the causal effect when an instrument is used to address the omission of relevant regressors. Illustrate the consequences of varying the dependence between the omitted regressor and the observed regressor.

**11. Simultaneous causality continues**

Use the demand-supply simulation and the moment conditions to infer the causal effect of a 10% tax, implementing a BETEL algorithm from scratch.

12. Perform a simulation of the instrumental variable quantile regression, and then perform inference using the general Bayes posterior framework. Use the Laplace asymmetric distribution to simulate the stochastic error, such that

$$\int_{-\infty}^0 f_\tau(\mu_i) d\mu_i = \tau,$$

which implies that the  $\tau$ -th quantile of  $\mu_i$  is 0. In addition, let

$$q_\tau(Y_i | X_i, D_i) = \beta_0 + \beta_1 X_i + \beta_2 D_i$$

denote the  $\tau$ -th quantile regression function of  $Y_i$  given  $X_i$  and  $D_i$ , with  $0 < \tau < 1$ . Specifically, consider the model

$$Y_i = \beta_{0\tau} + \beta_{1\tau}X_i + \beta_{2\tau}D_i + \mu_i,$$

together with

$$D_i = \gamma Z_i + \delta \mu_i, \quad Z_i \sim N(0, 1),$$

see Section 6.9.

**13. Example: Doubly robust Bayesian inference continues**

Perform a simulation exercise to assess the consequences of misspecifying the propensity score function and the outcome regression in doubly robust Bayesian inference.

# 14

---

## *Approximate Bayesian methods*

---

*Approximate Bayesian methods* are a family of techniques designed to handle situations where the likelihood function lacks an analytical expression, is highly complex, or the problem is high-dimensional, whether due to a large parameter space or a massive dataset [296]. In the former case, traditional Markov Chain Monte Carlo (MCMC) and importance sampling algorithms fail to provide a solution. In the latter, these algorithms struggle to produce accurate estimates within a reasonable timeframe, unless users modify them (see Chapter 12).

However, there is no free lunch, *Approximate Bayesian methods* address these challenges at the cost of providing an approximation to the posterior distribution rather than the *exact* posterior. Nonetheless, asymptotic results show that the approximation improves as the sample size increases.

In this chapter, I first present *simulation-based approaches*, which are designed to address situations where the likelihood is highly complex and may lack an analytical solution. In the second part, I introduce *optimization approaches*, which are intended to handle high-dimensional problems. Specifically, I discuss approximate Bayesian computation (ABC) and Bayesian synthetic likelihood (BSL), the two most common *simulation-based approaches*. Then, I present integrated nested Laplace approximations (INLA) and *variational Bayes* (VB), the two most common *optimization approaches* for high-dimensional problems.

---

### 14.1 Simulation-based approaches

Taking into account the fundamental equation for performing parameter inference in the Bayesian framework,

$$\pi(\boldsymbol{\theta} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta}),$$

we see in Section 4.1 that MCMC algorithms, such as the Gibbs sampler (Section 4.1.1) and Metropolis-Hastings (Section 4.1.2), require evaluation of the likelihood function  $p(\mathbf{y} | \boldsymbol{\theta})$  in the posterior conditional distribution or the acceptance probability, respectively. This is also the case for importance sampling when calculating the importance weights (Section 4.2).

Thus, what happens when the likelihood function does not have an analytical expression? This situation arises in many models involving unobserved heterogeneity (i.e., unobserved taste preferences), models defined by quantile functions (e.g., the g-and-k distribution), or dynamic equilibrium models (e.g., repeated game models).

*Simulation-based algorithms* provide a Bayesian solution when we face this situation, namely, when the likelihood function lacks an analytical expression or is highly complex. The only requirement is that we must be able to simulate synthetic data from the model conditional on the parameters. Therefore, these algorithms obtain an approximation to the posterior draws by simulating from the prior distribution  $\pi(\boldsymbol{\theta})$  and then using these draws to simulate from the likelihood  $p(\mathbf{y} | \boldsymbol{\theta})$ .

#### 14.1.1 Approximate Bayesian computation

*Approximate Bayesian Computation* (ABC) is designed to handle inferential situations where the likelihood function  $p(\mathbf{y} | \boldsymbol{\theta})$  is intractable or highly complex,  $\mathbf{y} \in \mathbb{R}^N$ . It was introduced in population genetics by [429, 342] and later generalized by [34]. The basic intuitive origin of ABC appears to have been introduced by [387]. A growing body of literature explores its applications in biology, cosmology, finance, economics, and other fields.

The requirement in ABC is the ability to simulate from the parametric model. The process begins by drawing samples from the prior distribution  $\pi(\boldsymbol{\theta})$  multiple times,  $\boldsymbol{\theta} \in \mathbb{R}^K$ , and then simulating data from the model given each  $\boldsymbol{\theta}^{(s)}, s = 1, 2, \dots, S$ . The resulting synthetic data,  $\mathbf{z}^{(s)} \in \mathbb{R}^N$  is used to compute summary statistics  $\boldsymbol{\eta}(\mathbf{z}^{(s)}) \in \mathbb{R}^L, L \geq K$ . These summary statistics are crucial for the performance of ABC and should be selected based on a thorough understanding of the model.

Next, we compare the synthetic summary statistics with the observed summary statistics  $\boldsymbol{\eta}(\mathbf{y})$  using a distance metric  $d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^{(s)})\}$ , typically the Euclidean distance. We retain the prior draws that generate synthetic summary statistics closest to the observed ones, that is,  $d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^{(s)})\} \leq \epsilon$ , forming an approximation of the posterior distribution  $\pi_\epsilon(\boldsymbol{\theta}, \mathbf{z} | \boldsymbol{\eta}(\mathbf{y}))$ .

The simplest algorithm is the accept/reject approximate Bayesian computation (ABC-AR) (see Algorithm A38).

---

**Algorithm A38** Accept/reject ABC

---

- 1: **for**  $s = 1, \dots, S$  **do**
  - 2:     Draw  $\boldsymbol{\theta}^s$  from  $\pi(\boldsymbol{\theta})$ ,
  - 3:     Simulate  $\mathbf{z}^s = (z_1^s, z_2^s, \dots, z_N^s)^\top$  from the model,  $p(\cdot | \boldsymbol{\theta}^s)$
  - 4:     Calculate  $d_{(s)} = d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^s)\}$
  - 5: **end for**
  - 6: Order the distances  $d_{(1)} \leq \dots \leq d_{(S)}$
  - 7: Select all  $\boldsymbol{\theta}^s$  such that  $d_{(i)} \leq \epsilon$ , where  $\epsilon > 0$  is the tolerance level.
-

Note that the posterior distribution is conditional on the summary statistics  $\boldsymbol{\eta}(\mathbf{y})$  and the tolerance parameter  $\epsilon$ . This implies that we obtain an approximation to the target distribution  $\pi(\boldsymbol{\theta} | \mathbf{y})$ , that is  $\pi_\epsilon(\boldsymbol{\theta} | \boldsymbol{\eta}(\mathbf{y}))$ , because  $\boldsymbol{\eta}(\mathbf{y})$  is not a sufficient statistic in most cases, and  $\epsilon > 0$ , these conditions introduce bias [51]. However, ABC performs well compared to full-likelihood approaches in low-dimensional parameter spaces [34].

Furthermore, [144] show in Theorems 1 and 2 that Bayesian consistency and asymptotic normality hold, provided that  $\epsilon \rightarrow 0$  fast enough as  $N \rightarrow +\infty$ . In particular, the requirement is that the proportion of accepted draws converges to 0 at a rate faster than  $N^{-K/2}$ . Additionally, Theorem 2 in [144] shows that  $100(1 - \alpha)\%$  Bayesian credible regions using ABC have frequentist coverage of  $100(1 - \alpha)\%$ .

We should note from these asymptotic results that ABC suffers from the *curse of dimensionality*. Specifically, given a sample size of 1,000 and two parameters, the proportion of accepted draws should be 0.1%, meaning we would require one million prior draws to obtain 1,000 posterior draws. On the other hand, if the number of parameters is three, we would require 31.62 million prior draws. This limitation of ABC has attracted attention; see Chapter 8 of [413] for some potential solutions.

It is a common practice in ABC to perform a regression adjustment after retaining the draws [34, 265, 413]. This adjustment reduces bias in posterior draws by performing a simple linear regression between the selected draws and the discrepancy between the observed and simulated summary statistics,  $\theta_k^{(s)} = \alpha_k + (\boldsymbol{\eta}(\mathbf{y}) - \boldsymbol{\eta}(\mathbf{z}^{(s)}))^\top \boldsymbol{\beta}_k + \mu_k^{(s)}$ ,  $k = 1, 2, \dots, K$ . Then, the posterior draws are adjusted using the slope estimate  $\theta_k^{\text{adj},(s)} = \theta_k^{(s)} - (\boldsymbol{\eta}(\mathbf{y}) - \boldsymbol{\eta}(\mathbf{z}^{(s)}))^\top \hat{\boldsymbol{\beta}}_k$ . Other regression adjustment strategies are also used, such as local linear regression, ridge regression, and neural networks. See the *abc* package in **R**.

The favorable asymptotic sampling properties of ABC rely on correct model specification. [145] demonstrate that when the assumed model is misspecified, the asymptotic behavior of ABC can deteriorate. In particular, the posterior shape becomes asymptotically non-Gaussian, and the behavior of the posterior mean remains generally unknown. Additionally, regression adjustment approaches can produce posteriors that differ significantly from their simpler accept/reject counterparts.

Given these concerns, testing model specification in ABC is essential. This can be done using simulated goodness-of-fit statistics [43, 273], predictive p-values [43], discrepancy diagnostics [145], and asymptotic tests [358] to evaluate model adequacy.

The accept/reject ABC algorithm is inefficient, as all draws are independent; thus, there is no learning from previous draws. This intensifies the computational burden. Therefore, [292, 448] introduced Markov Chain Monte Carlo ABC (ABC-MCMC) algorithms, and [414, 127, 116, 263] proposed sequential Monte Carlo approaches (ABC-SMC). However, results comparing

ABC-MCMC and ABC-SMC with ABC-AR are controversial regarding computational efficiency [43]. In addition, ABC-AR is very simple and easily allows parallel computing [143]. Nevertheless, ABC-SMC is now the recommended approach, as it does not require tuning the algorithm's tolerance [296], and there are open-source implementations that facilitate its use.

New developments in ABC have focused on using empirical measures calculated from the observed ( $\hat{\mu}_n$ ) and synthetic ( $\hat{\mu}_{\theta}^{(s)}$ ) data to replace summary statistics. Thus,  $d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^{(s)})\}$  is replaced by  $\mathcal{D}\left\{\hat{\mu}_n, \hat{\mu}_{\theta}^{(s)}\right\}$ , where the latter is a discrepancy measure, such as the Kullback-Leibler divergence [228]. However, [126] found in their simulation exercises that the best-performing summary statistics approach performs at least as well as the best discrepancy-measure approaches. The key point is to select informative summary statistics.

#### **Example: g-and-k distribution for financial returns**

The g-and-k distribution is a highly flexible distribution capable of capturing skewness and heavy tails through its parameters. This makes it particularly useful for modeling real-world data that deviate from normality, especially in fields like finance, where outliers are common. However, this distribution lacks a closed-form expression for its density function.

The g-and-k distribution is defined by its quantile function [128]. Specifically, it is specified through its inverse cumulative distribution function,

$$Q(p \mid \theta) = F^{-1}(p \mid \theta),$$

where  $F = P(U \leq u)$ , and  $Q(p \mid \theta)$  represents the  $p$ -quantile [365]. The quantile function of the g-and-k distribution is given by

$$Q^{gk}\{z(p) \mid a, b, c, g, k\} = a + b \left[ 1 + c \frac{1 - \exp\{-gz(p)\}}{1 + \exp\{-gz(p)\}} \right] \{1 + z(p)^2\}^k z(p),$$

where  $z(p)$  is the standard normal quantile function, and  $c = 0.8$  is a commonly suggested value.

In the g-and-k distribution,  $a$  is the location parameter, and  $b$  is the scale parameter, controlling the dispersion. The parameters  $g$  and  $k$  determine the levels of skewness and kurtosis, respectively, while  $c$  modifies the impact of skewness.

[128] propose a moving average of order one using a g-and-k distribution to model exchange rate log returns. In particular,

$$z_t = \epsilon_t + \theta_1 \epsilon_{t-1}, \quad t = 1, \dots, 524,$$

where  $\epsilon_t \sim N(0, 1)$ .

The values of  $z_t$  are then divided by  $(1 + \theta_1^2)^{1/2}$  to ensure that they marginally follow a standard normal distribution. Thus, simulating g-and-k data requires only substituting  $z_t$  into the quantile function.

We model exchange rate log daily returns from USD/EUR one year before and after the WHO declared the COVID-19 pandemic on 11 March 2020. We use the dataset *ExchangeRate.csv* from our GitHub repository. Our ABC implementation uses twelve summary statistics: the seven octiles, the interquartile range, robust measures of skewness and kurtosis, and the autocorrelations of order one and two (see [128] and code below). We adopt the prior distributions proposed by [358],

$$\begin{aligned}\theta_1 &\sim U(-1, 1), \quad a \sim U(0, 5) \quad b \sim U(0, 5) \\ g &\sim U(-5, 5), \quad k \sim U(-0.5, 5).\end{aligned}$$

We use the *EasyABC* package in **R** to implement the ABC accept/reject (ABC-AR) algorithm using 150,000 prior draws with an acceptance rate of 0.67%. We also apply the ABC Markov chain Monte Carlo (ABC-MCMC) method [292] and the sequential Monte Carlo ABC (ABC-SMC) method [263] to compare the results across different ABC algorithms.<sup>1</sup> We generate 100,000 samples and retain 1% in ABC-MCMC, and 30,000 samples, keeping 3.4%, with a stopping criterion of 5% in ABC-SMC. These settings imply that the three algorithms require approximately the same computational time. As ABC is a simulation-based method, the computational burden is mostly driven by the speed of simulating the process; thus, the *EasyABC* package has parallel computing algorithms to speed up the processes. Users can refer to the cited references for algorithmic details, and the *EasyABC* package for parallel computing implementation.

In Exercise 1, we ask to program the ABC accept/reject algorithm (Algorithm A38) from scratch and compare the results with those obtained using the ABC-AR implementation in the *EasyABC* package.

The following code presents the results, and Figure 14.1 compares the posterior distributions of  $\theta_1$ ,  $g$ , and  $k$  using the three methods. In this figure, we observe that ABC-MCMC (red) and ABC-SMC (green) exhibit similar performance, and both approaches provide more information than ABC-AR (blue). There is marginal positive evidence that the moving average coefficient is positive, and the three algorithms yield similar means, although ABC-AR exhibits lower precision. Since the posterior distribution of  $g$  is centered around zero, the distribution appears to be symmetric around its median. Meanwhile, a positive  $k$  indicates that the distribution has heavier tails than a normal distribution, implying a higher likelihood of extreme values (outliers) in the exchange rate USD/EURO.

---

<sup>1</sup>Note that this setting does not satisfy the asymptotic requirements for Bayesian consistency. However, it serves as a pedagogical exercise.

**R code. Exchange rate log returns: Approximate Bayesian computation**

```

1 ##### ABC Exchange rate og returns: USD/EURO
2 rm(list = ls()); set.seed(010101)
3 library(EasyABC)
4 dfExcRate <- read.csv(file = "https://raw.githubusercontent.com/BEsmarter-consultancy/BSTApp/refs/heads/master/DataApp/ExchangeRate.csv", sep = ",", header = T)
5 attach(dfExcRate); n <- length(USDEUR)
6 # Summary statistics
7 SumSt <- function(y) {
8   Oct <- quantile(y, c(0.125, 0.25, 0.375, 0.5, 0.625, 0.75,
9   0.875))
10  eta1 <- Oct[6] - Oct[2]
11  eta2 <- (Oct[6] + Oct[2] - 2 * Oct[4]) / eta1
12  eta3 <- (Oct[7] - Oct[5] + Oct[3] - Oct[1]) / eta1
13  autocor <- acf(y, lag = 2, plot = FALSE)
14  autocor[["acf"]][2:3]
15  Etay <- c(Oct, eta1, eta2, eta3, autocor[["acf"]][2:3])
16  return(Etay)
17 }
18 # g-and-k distribution
19 RGKnewSum <- function(par) {
20   z <- NULL
21   theta <- par[1]; a <- par[2]; b <- par[3]; g <- par[4]; k
22   <- par[5]
23   e <- rnorm(n + 1)
24   for(t in 2:(n + 1)){
25     zt <- e[t] + theta * e[t-1]
26     z <- c(z, zt)
27   }
28   zs <- z / (1 + theta^2)^0.5
29   x <- a + b * (1 + 0.8 * (1 - exp(-g * zs)) / (1 + exp(-g *
30   zs))) * (1 + zs^2)^k * zs
31   Etaz <- SumSt(x)
32   return(Etaz)
33 }
34 toy_prior <- list(c("unif",-1,1), c("unif",0,5), c("unif",
35 0,5), c("unif", -5,5), c("unif", -0.5,5))
36 sum_stat_obs <- SumSt(USDEUR)
37 tick <- Sys.time()
38 ABC_AR <- ABC_rejection(model=RGKnewSum, prior=toy_prior,
39   summary_stat_target = sum_stat_obs, nb_simul=150000, tol
40   = 0.0067,
41   progress_bar = TRUE)
42 tock <- Sys.time()
43 tock - tick

```

### R code. Exchange rate log returns: Approximate Bayesian computation

```

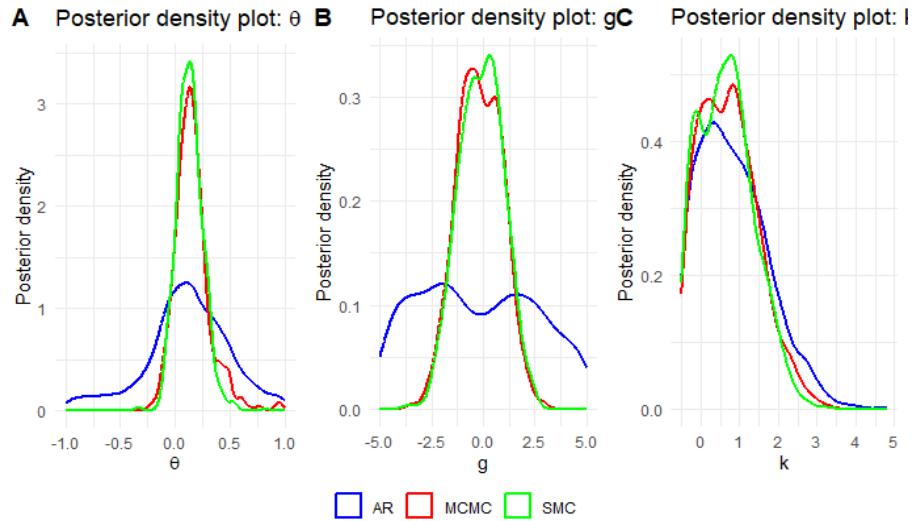
1 PostABCAR <- coda::mcmc(ABC_AR$param)
2 summary(PostABCAR)
3 tick <- Sys.time()
4 ABC_MCMC <- ABC_mcmc(method="Marjoram", model=RGKnewSum,
   prior=toy_prior, summary_stat_target=sum_stat_obs, n_rec
   = 100000, progress_bar = TRUE)
5 tock <- Sys.time()
6 tock - tick
7 PostABCMCMC <- coda::mcmc(ABC_MCMC[["param"]][order(ABC_MCMC
   [["dist"]])][1:1000],)
8 summary(PostABCMCMC)
9 tick <- Sys.time()
10 ABC_SMC <- ABC_sequential(method="Lenormand", model=
   RGKnewSum, prior=toy_prior, summary_stat_target=sum_stat
   _obs, nb_simul = 30000, alpha = 0.034, p_acc_min = 0.05,
   progress_bar = TRUE)
11 PostABCSMC <- coda::mcmc(ABC_SMC[["param"]])
12 summary(PostABCSMC)
13 tock <- Sys.time()
14 tock - tick
15 # Figures
16 library(ggplot2); library(latex2exp)
17 Sp <- 1000
18 df1 <- data.frame(Value = c(PostABCAR[1:Sp,1], PostABCMCMC
   [1:Sp,1], PostABCSMC[1:Sp,1]),
19 Distribution = factor(c(rep("AR", Sp), rep("MCMC", Sp), rep(
   "SMC", Sp))))
20 dentheta <- ggplot(df1, aes(x = Value, color = Distribution))
   + geom_density(linewidth = 1) + labs(title = TeX("$\theta$"),
   "Posterior density plot: $\theta$"), x = TeX("$\theta$"),
   y = "Posterior density") + scale_color_manual(values = c(
   "blue", "red", "green")) + theme_minimal() + theme(
   legend.title = element_blank())
21 df2 <- data.frame(Value = c(PostABCAR[1:Sp,4], PostABCMCMC
   [1:Sp,4], PostABCSMC[1:Sp,4]),
22 Distribution = factor(c(rep("AR", Sp), rep("MCMC", Sp), rep(
   "SMC", Sp))))
23 deng <- ggplot(df2, aes(x = Value, color = Distribution)) +
   geom_density(linewidth = 1) + labs(title = "Posterior
   density plot: g", x = "g", y = "Posterior density") +
24 scale_color_manual(values = c("blue", "red", "green")) +
   theme_minimal() + theme(legend.title = element_blank())
25 df3 <- data.frame(Value = c(PostABCAR[1:Sp,5], PostABCMCMC
   [1:Sp,5], PostABCSMC[1:Sp,5]),
26 Distribution = factor(c(rep("AR", Sp), rep("MCMC", Sp), rep(
   "SMC", Sp))))
```

**R code. Exchange rate log returns: Approximate Bayesian computation**

```

1 denk <- ggplot(df3, aes(x = Value, color = Distribution)) +
  geom_density(linewidth = 1) + labs(title = "Posterior
  density plot: k", x = "k", y = "Posterior density") +
2 scale_color_manual(values = c("blue", "red", "green")) +
  theme_minimal() + theme(legend.title = element_blank())
3 library(ggpubr)
4 ggarrange(dentheta, deng, denk, labels = c("A", "B", "C"),
  ncol = 3, nrow = 1,
5 legend = "bottom", common.legend = TRUE)

```



**FIGURE 14.1**

Posterior distributions using approximate Bayesian computation: Exchange rate log returns USD/EURO.

#### 14.1.2 Bayesian synthetic likelihood

Note that in ABC, in most cases, we target the posterior distribution  $\pi_\epsilon(\boldsymbol{\theta} | \boldsymbol{\eta}(\mathbf{y}))$  rather than  $\pi(\boldsymbol{\theta} | \mathbf{y})$ , as  $\boldsymbol{\eta}(\mathbf{y})$  is not a sufficient statistic,  $\mathbf{y} \in \mathbb{R}^N$ ,  $\boldsymbol{\theta} \in \mathbb{R}^K$ ,  $\boldsymbol{\eta}(\mathbf{y}) \in \mathbb{R}^L$ ,  $L \geq K$ . This can be beneficial since discarding information may improve the behavior of the likelihood or make the inference more robust

to model misspecification [341]. Given the intractability of  $p(\mathbf{y} \mid \boldsymbol{\theta})$ , it is highly likely that  $p(\boldsymbol{\eta}(\mathbf{y}) \mid \boldsymbol{\theta})$  is also intractable. [453] addressed this issue by introducing an auxiliary model for the summary statistics, assuming  $p_a(\boldsymbol{\eta}(\mathbf{y}) \mid \boldsymbol{\theta}) = N(\boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})$ .

Bayesian synthetic likelihood (BSL) arises when this auxiliary likelihood is combined with a prior distribution on the parameter [129, 341],

$$\pi_a(\boldsymbol{\theta} \mid \boldsymbol{\eta}(\mathbf{y})) \propto p_a(\boldsymbol{\eta}(\mathbf{y}) \mid \boldsymbol{\theta})\pi(\boldsymbol{\theta}),$$

where the subscript  $a$  indicates that this is an approximation due to the Gaussian assumption. This is referred to as the idealized BSL posterior. However, note that  $p_a(\boldsymbol{\eta}(\mathbf{y}) \mid \boldsymbol{\theta})$  is rarely available, as  $\boldsymbol{\mu}_{\boldsymbol{\theta}}$  and  $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}$  are generally unknown. Therefore, we estimate these quantities using simulations from the model given a realization of  $\boldsymbol{\theta}$ :<sup>2</sup>

$$\begin{aligned}\hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}} &= \frac{1}{M} \sum_{m=1}^M \boldsymbol{\eta}(\mathbf{z}^{(m)}), \\ \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}} &= \frac{1}{M-1} \sum_{m=1}^M (\boldsymbol{\eta}(\mathbf{z}^{(m)}) - \hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}})(\boldsymbol{\eta}(\mathbf{z}^{(m)}) - \hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}})^{\top}.\end{aligned}$$

Then, we have

$$\pi_{a,M}(\boldsymbol{\theta} \mid \boldsymbol{\eta}(\mathbf{y})) \propto p_{a,M}(\boldsymbol{\eta}(\mathbf{y}) \mid \boldsymbol{\theta})\pi(\boldsymbol{\theta}),$$

where  $p_{a,M}(\boldsymbol{\eta}(\mathbf{y}))$  uses the estimates, which in turn depends on the number of draws  $M$  to calculate  $\hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}}$  and  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}}$ . Note that although we can have unbiased estimators for the mean and covariance, in general,  $N(\boldsymbol{\eta}(\mathbf{y}) \mid \hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}}, \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}})$  is not an unbiased estimator of  $N(\boldsymbol{\eta}(\mathbf{y}) \mid \boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})$ . [10] show an unbiased estimator for BSL.

[141] show that  $\pi_{a,M}(\boldsymbol{\theta})$  converges asymptotically to a Gaussian distribution and that the  $100(1-\alpha)\%$  Bayesian credible regions using BSL have frequentist coverage of  $100(1-\alpha)\%$ . In addition, the posterior mean from BSL is also asymptotically Gaussian. These results require convenient estimation of the covariance matrix and that  $M \rightarrow \infty$  as  $N \rightarrow \infty$ ,  $M = C\lfloor N^{\gamma} \rfloor$ ,  $C > 0$ ,  $\gamma > 0$  and  $\lfloor x \rfloor$  denoting the integer floor of  $x$ . Thus, the choice of  $M$  does not drastically affect the desirable asymptotic properties of BSL [141]. In addition, [341] find in their examples that posterior inference depends only weakly on  $M$ . Thus,  $M$  can be chosen to balance computational efficiency such that the standard deviation of the log synthetic likelihood is between 1 and 3 [11].

A critical aspect of BSL is the estimation of the covariance matrix, which can be computationally demanding in high-dimensional settings. However, [141] propose an adjusted approach to BSL that allows the use of a simple, though potentially misspecified, estimator of the covariance matrix (see Equation 5 and the related discussion in their paper for details).

<sup>2</sup>There are better ways to calculate the covariance matrix, see [141].

If the normality assumption for summary statistics is too restrictive, [9] propose a robust BSL method based on a semi-parametric approach. Additionally, [142] show that when the model is misspecified, i.e., there is no compatibility between the assumed statistical model and the true data-generating process, BSL can lead to unreliable parameter inference. To address this issue, they propose a new BSL method that detects model misspecification and provides more reliable inference.

We can perform BSL using the Algorithm A39. We can use a random walk Metropolis-Hastings to set the proposal distribution. The covariance matrix of the random walk proposal can be tuned using an initial pilot run.

---

**Algorithm A39** Bayesian synthetic likelihood

---

```

1: for  $s = 1, \dots, S$  do
2:   Draw  $\boldsymbol{\theta}^c \sim q(\boldsymbol{\theta} | \boldsymbol{\theta}^{s-1})$ 
3:   for  $m = 1, \dots, M$  do
4:     Simulate  $\mathbf{z}^m = (z_1^m, z_2^m, \dots, z_N^m)^\top$  from the model,  $p(\cdot | \boldsymbol{\theta}^c)$ 
5:     Calculate  $\boldsymbol{\eta}(\mathbf{z}^m)$ 
6:   end for
7:   Calculate  $\hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}^c}$  and  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}^c}$ 
8:   Compute  $p_a^c(\boldsymbol{\eta}(\mathbf{y}) | \boldsymbol{\theta}^c) = N(\hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}^c}, \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}^c})$  and  $p_a^{s-1}(\boldsymbol{\eta}(\mathbf{y}) | \boldsymbol{\theta}^{s-1}) =$ 
     $N(\hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}^{s-1}}, \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}^{s-1}})$ 
9:   Compute the acceptance probability

$$\alpha(\boldsymbol{\theta}^{s-1}, \boldsymbol{\theta}^c) = \min \left\{ 1, \frac{p_a^c(\boldsymbol{\eta}(\mathbf{y}) | \boldsymbol{\theta}^c) \pi(\boldsymbol{\theta}^c) q(\boldsymbol{\theta}^{s-1} | \boldsymbol{\theta}^c)}{p_a^{s-1}(\boldsymbol{\eta}(\mathbf{y}) | \boldsymbol{\theta}^{s-1}) \pi(\boldsymbol{\theta}^{s-1}) q(\boldsymbol{\theta}^c | \boldsymbol{\theta}^{s-1})} \right\}$$

10:  Draw  $u$  from  $U(0, 1)$ 
11:  if  $u < \alpha$  then
12:    Set  $\boldsymbol{\theta}^s = \boldsymbol{\theta}^c$ ,  $\hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}^s} = \hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}^c}$  and  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}^s} = \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}^c}$ 
13:  else
14:    Set  $\boldsymbol{\theta}^s = \boldsymbol{\theta}^{s-1}$ ,  $\hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}^s} = \hat{\boldsymbol{\mu}}_{\boldsymbol{\theta}^{s-1}}$  and  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}^s} = \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\theta}^{s-1}}$ 
15:  end if
16: end for

```

---

An advantage of BSL over ABC is that it does not require selecting a tolerance parameter  $\epsilon$ . Furthermore, BSL is more computationally efficient than ABC when dealing with a high-dimensional vector of summary statistics as the acceptance rate of the former is asymptotically non-vanishing [141].

On the other hand, ABC is asymptotically more efficient than BSL, and it imposes very weak requirements on the choice of summary statistics, whereas BSL requires summary statistics that satisfy central limit theorems (CLTs), and consistent estimators of the covariance matrix. However, asymptotically, both approaches provide reliable inference, provided their respective requirements are met [296].

In practice, BSL may be more convenient than ABC when the summary

statistics are high-dimensional and satisfy CLT conditions. Otherwise, ABC may be the better alternative.

### Example: Simulation exercise

We simulate a dataset following the specification of the g-and-k distribution for the financial returns example, setting  $\theta_1 = 0.8$ ,  $a = 1$ ,  $b = 0.5$ ,  $g = -1$ , and  $k = 1$ , with a sample size of 500. We use the same priors and summary statistics as in that example.

We use the *BSL* package in **R** to perform posterior inference using BSL. Additionally, we implement our own BSL sampler to compare the results with the vanilla algorithm from the package. We run the BSL algorithms with  $M = 200$ ,  $S = 11,000$ , a burn-in of 1,000, and a thinning parameter of 10, using a random walk proposal distribution. This setting results in similar computational times between the two implementations (scratch and package) and is also comparable to the ABC implementation in Exercise 1. Take into account that BSL is a simulation-based method; therefore, the computational burden is mostly driven by the speed of simulating the process. The *BSL* package has parallel computing algorithms to speed up the processes, and users should refer to the *BSL* package for details.

The following code demonstrates this procedure. The summary statistics appear to be approximately normally distributed, as shown in Figure 14.2, where the red line represents the normal density and the black line represents the estimated density of the summary statistic. Figure 14.3 displays the posterior distributions from both algorithms for  $\theta_1$ ,  $g$ , and  $k$ . In general, the 95% credible intervals encompass the true parameter values. However, the posterior draws from the *BSL* package are more informative compared to our implementation of Algorithm A39. Additionally, BSL outperforms the results of ABC in Exercise 1, yielding more precise posterior distributions. This pattern has been observed in other settings [126, 296]. Both ABC and BSL produce posterior distributions centered at the true parameter values, although the posterior distribution of  $\theta_1$  deviates from this pattern.

#### *R* code. g-and-k simulation: Bayesian synthetic likelihood

```

1 ##### BSL: g-and-k simulation #####
2 rm(list = ls()); set.seed(010101); library(BSL)
3 # Simulate g-and-k data
4 RGKnew <- function(par) {
5   z <- NULL
6   theta <- par[1]; a <- par[2]; b <- par[3]; g <- par[4]

```

**R code. g-and-k simulation: Bayesian synthetic likelihood**

```

1   k <- par[5]; e <- rnorm(n + 1)
2   for(t in 2:(n + 1)){
3     zt <- e[t] + theta * e[t-1]
4     z <- c(z, zt)
5   }
6   zs <- z / (1 + theta^2)^0.5
7   x <- a + b * (1 + 0.8 * (1 - exp(-g * zs)) / (1 + exp(-g *
8     zs))) * (1 + zs^2)^k * zs
9   return(x)
10 }
11 # Summary statistics
12 SumSt <- function(y) {
13   Oct <- quantile(y, c(0.125, 0.25, 0.375, 0.5, 0.625, 0.75,
14     0.875))
15   eta1 <- Oct[6] - Oct[2]
16   eta2 <- (Oct[6] + Oct[2] - 2 * Oct[4]) / eta1
17   eta3 <- (Oct[7] - Oct[5] + Oct[3] - Oct[1]) / eta1
18   autocor <- acf(y, lag = 2, plot = FALSE)
19   autocor[["acf"]][2:3]
20   EtaY <- c(Oct, eta1, eta2, eta3, autocor[["acf"]][2:3])
21   return(EtaY)
22 }
23 # Prior function
24 LogPrior <- function(par){
25   LogPi <- log(par[1] > -1 & par[1] < 1 & par[2] > 0 & par
26     [2] < 5 & par[3] > 0 & par[3] < 5 & par[4] > -5 & par[4]
27     < 5 & par[5] > -0.5 & par[5] < 5)
28   return(LogPi)
29 }
30 # Population parameters
31 theta1 <- 0.8; a <- 1; b <- 0.5; g <- -1; k <- 1
32 parpop <- c(theta1, a, b, g, k)
33 K <- 5
34 n <- 500
35 y <- RGKnew(par = parpop)
36 # Algorithm parameters
37 M <- 200 # Number of iterations to calculate mu and sigma
38 S <- 11000 # Number of MCMC iterations
39 burnin <- 1000 # Burn in iterations
40 thin <- 10 # Thining parameter
41 keep <- seq(burnin + 1, S, thin)
42 par0 <- c(0.5, 2, 1, 0, 1)
43 Modelgk <- newModel(fnSim = RGKnew, fnSum = SumSt, theta0 =
44   par0, fnLogPrior = LogPrior, verbose = FALSE)
45 validObject(Modelgk)
46 # Check if the summary statistics are roughly normal
47 simgk <- simulation(Modelgk, n = M, theta = par0, seed = 10)
48 par(mfrow = c(4, 3))

```

### R code. g-and-k simulation: Bayesian synthetic likelihood

```

1  for (i in 1:12){
2    eval <- seq(min(simgk$ssx[, i]), max(simgk$ssx[, i]),
3                 0.001)
4    densnorm <- dnorm(eval, mean = mean(simgk$ssx[, i]), sd(
5      simgk$ssx[, i]))
6    plot(density(simgk$ssx[, i]), main = "", xlab = "")
7    lines(eval, densnorm, col = "red")
8  }
9 ###### Program from scratch ######
10 Lims <- matrix(c(-1, 0, 0, -5, -0.5, 1, rep(5, 4)), 5, 2)
11 parPost <- matrix(NA, S, K); parPost[1,] <- par0 ; EtaY <-
12   SumSt(y = y)
13 Zsl <- replicate(M, RGKnew(par = parPost[1,]))
14 EtasZsl <- t(apply(Zsl, 2, SumSt))
15 Usl <- colMeans(EtasZsl); SIGMASl <- var(EtasZsl)
16 dnormsl <- mvtnorm::dmvnorm(EtaY, Usl, SIGMASl, log = TRUE)
17 tune <- 0.005; CoVarRW <- tune*diag(K)
18 accept <- rep(0, S); tick1 <- Sys.time()
19 pb <- txtProgressBar(min = 0, max = S, style = 3)
20 for (s in 2:S){
21   parc <- MASS::mvrnorm(1, mu = parPost[s-1,], Sigma =
22     CoVarRW)
23   RestCheck <- NULL
24   for(j in 1:K){
25     if(parc[j] < Lims[j,1] | parc[j] > Lims[j,2]){
26       Rej <- 1
27     }else{Rej <- 0}
28     RestCheck <- c(RestCheck, Rej)
29   }
30   if(sum(RestCheck) != 0){
31     parPost[s,] <- parPost[s-1,]; accept[s] <- 0
32   }else{
33     Z <- replicate(M, RGKnew(par = parc))
34     EtasZ <- t(apply(Z, 2, SumSt))
35     Um <- colMeans(EtasZ); SIGMAM <- var(EtasZ)
36     dnormc <- mvtnorm::dmvnorm(EtaY, Um, SIGMAM, log = TRUE)
37     if(s > round(0.1*S, 0) & s < round(0.2*S, 0)){
38       CoVarRW <- var(parPost, na.rm = TRUE)
39     }
40     if(dnormc == -Inf){
41       parPost[s,] <- parPost[s-1,]; accept[s] <- 0
42     }
43   }
44 }
45 
```

### R code. g-and-k simulation: Bayesian synthetic likelihood

```

1 } else{
2   alpha <- min(1, exp(dnormc - dnormsl)); U <- runif(1)
3   if(U <= alpha){
4     parPost[s,] <- parc; Usl <- Um; SIGMAsl <- SIGMAM
5     dnormsl <- dnormc; accept[s] <- 1
6   } else{
7     parPost[s,] <- parPost[s-1,]; accept[s] <- 0
8   }
9 }
10 }
11 setTxtProgressBar(pb, s)
12 }
13 close(pb); tock1 <- Sys.time(); tock1 - tick1
14 mean(accept)
15 PostChainOwn <- coda::mcmc(parPost[keep,])
16 summary(PostChainOwn)
17 ##### BSL package
18 tick <- Sys.time()
19 Resultsgk <- bsl(y = y, n = M, M = S, model = Modelgk,
20 covRandWalk = CoVarRW,
21 method = "BSL", thetaNames = expression(theta, a, b, g, k),
22 plotOnTheFly = TRUE)
23 tock <- Sys.time(); tock - tick
24 PostChain <- coda::mcmc(Resultsgk@theta[keep,])
25 summary(PostChain)
26 Resultsgk@acceptanceRate
27 plot(Resultsgk@loglike[keep], type = "l")
28 sd(Resultsgk@loglike[keep])
29 # Figures
30 library(ggplot2); library(latex2exp)
31 Sp <- length(keep)
32 df1 <- data.frame(Value = c(PostChain[,1], PostChainOwn[,1]),
33 , Distribution = factor(c(rep("BSL", Sp), rep(
34 "BSLscratch", Sp))))
35 dentheta <- ggplot(df1, aes(x = Value, color = Distribution)
36 ) + geom_density(linewidth = 1) + labs(title = TeX(""
37 Posterior density plot: $theta$"), x = TeX("$theta$"), y
38 = "Posterior density") + geom_vline(xintercept = thetal1
39 , linetype = "dashed", color = "red", linewidth = 1) +
40 scale_color_manual(values = c("blue", "red")) + theme_
41 minimal() + theme(legend.title = element_blank())
42 df2 <- data.frame(Value = c(PostChain[,4], PostChainOwn[,4]),
43 , Distribution = factor(c(rep("BSL", Sp), rep(
44 "BSLscratch", Sp))))

```

### R code. g-and-k simulation: Bayesian synthetic likelihood

```

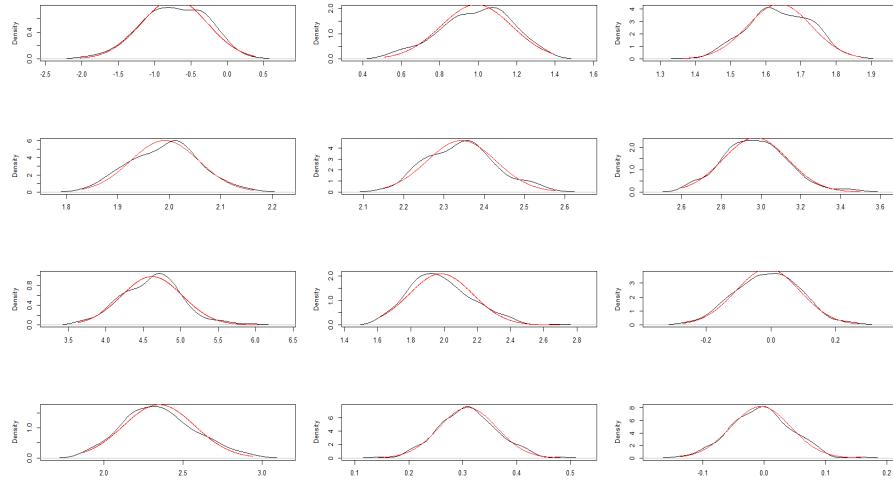
1 deng <- ggplot(df2, aes(x = Value, color = Distribution)) +
  geom_density(linewidth = 1) + labs(title = "Posterior
density plot: g", x = "g", y = "Posterior density") +
  geom_vline(xintercept = g, linetype = "dashed", color =
  "red", linewidth = 1) + scale_color_manual(values = c("blue",
  "red")) + theme_minimal() + theme(legend.title =
  element_blank())
2 df3 <- data.frame(Value = c(PostChain[,5], PostChainOwn[,5])
, Distribution = factor(c(rep("BSL", Sp), rep("BSLscratch", Sp))))
3 denk <- ggplot(df3, aes(x = Value, color = Distribution)) +
  geom_density(linewidth = 1) + labs(title = "Posterior
density plot: k", x = "k", y = "Posterior density") +
  geom_vline(xintercept = k, linetype = "dashed", color =
  "red", linewidth = 1) + scale_color_manual(values = c("blue",
  "red")) + theme_minimal() + theme(legend.title =
  element_blank())
4 library(ggpubr)
5 ggarrange(dentheta, deng, denk, labels = c("A", "B", "C"),
  ncol = 3, nrow = 1,
6 legend = "bottom", common.legend = TRUE)

```

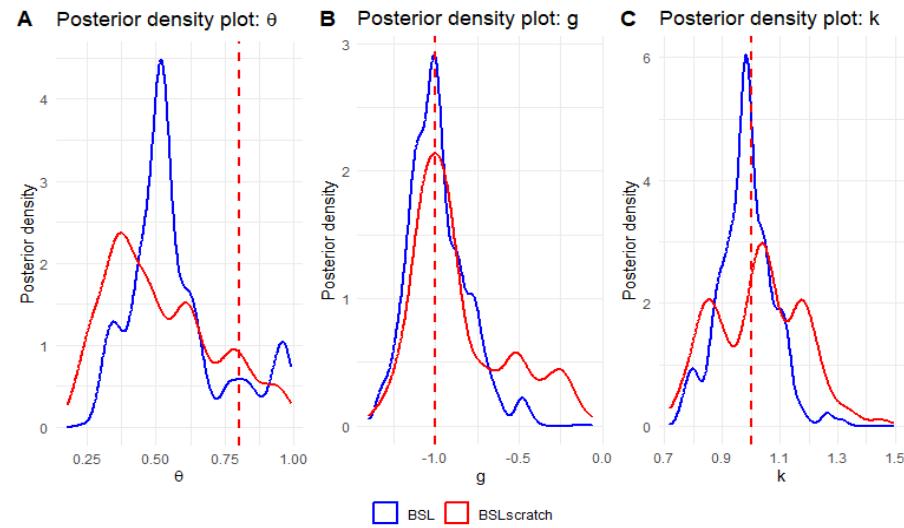
## 14.2 Optimization approaches

Traditional MCMC and importance sampling (IS) algorithms require pointwise evaluation of the likelihood function, which entails a massive number of operations when applied to very large datasets. Unfortunately, these algorithms are not designed to be *scalable*, at least in their standard form (see Chapter 12 for alternatives). Moreover, when the parameter space is large, they also lack *scalability* with respect to the number of parameters. Therefore, approximation methods should be considered even when the likelihood function has an analytical expression.

*Optimization approaches* are designed to scale efficiently with high-dimensional parameter spaces and large datasets. The key idea is to replace simulation with optimization. In this section, we introduce the most common optimization approaches within the Bayesian inferential framework.

**FIGURE 14.2**

Density estimates of summary statistics: g-and-k simulation.

**FIGURE 14.3**

Posterior distributions using Bayesian synthetic likelihood: g-and-k simulation.

#### 14.2.1 Integrated nested Laplace approximations

*Integrated Nested Laplace Approximations* (INLA) is a deterministic approach for Bayesian inference in latent Gaussian models (LGMs) [392]. In particular,

INLA approximates the marginal posterior distributions using a combination of Laplace approximations, low-dimensional deterministic integration and optimization steps in sparse covariance settings [393]. The advantages of INLA compared to MCMC are that it is fast and does not suffer from poor mixing.

The point of departure is a structured additive regression model, where the response variable  $y_i$  belongs to the exponential family such that the mean  $\mu_i$  is linked to a linear predictor  $\eta_i$  through the link function  $g(\cdot)$ , that is,  $\eta_i = g(\mu_i)$ , where

$$\mu_i = \alpha + \boldsymbol{\beta}^\top \mathbf{x}_i + \sum_{j=1}^J f^{(j)}(u_{ji}) + \epsilon_i,$$

where  $\alpha$  is the general intercept,  $f^{(j)}$  are unknown functions of the covariates  $\mathbf{u}_i$ ,  $\boldsymbol{\beta}$  is a  $K$ -dimensional vector of linear effects associated with regressors  $\mathbf{x}_i$ , and  $\epsilon_i$  is the unstructured error.

Note that latent Gaussian models encompass a wide range of relevant empirical models depending on the specific elements and structure involved in  $f^{(j)}$ , such as generalized linear models with unobserved heterogeneity, spatial and/or temporal dependence, and semi-parametric models.

Latent Gaussian models assign a Gaussian prior to  $\alpha$ ,  $\boldsymbol{\beta}$ ,  $f^{(j)}$ , and  $\epsilon_i$ . Let  $\mathbf{z}$  denote the vector of all latent Gaussian variables  $\{\alpha, \boldsymbol{\beta}, f^{(j)}, \eta_i\}$ , where the dimension is potentially  $P = 1 + K + J + N$  (although this is not always the case), and let  $\boldsymbol{\theta}$  be the  $m$ -dimensional vector of hyperparameters. Then, the density  $\pi(\mathbf{z} | \boldsymbol{\theta}_1)$  is Gaussian with mean zero and precision matrix (i.e., the inverse of the covariance matrix)  $\mathbf{Q}(\boldsymbol{\theta}_1)$ .

The distribution of  $\mathbf{y}$  is  $p(\mathbf{y} | \mathbf{z}, \boldsymbol{\theta}_2)$  such that  $y_i$  are conditional independent given  $z_i$  and  $\boldsymbol{\theta}_2$ ,  $\boldsymbol{\theta} = [\boldsymbol{\theta}_1^\top \boldsymbol{\theta}_2^\top]^\top$ ,  $i = 1, 2, \dots, N$ . Thus, the posterior distribution is

$$\begin{aligned} \pi(\mathbf{z}, \boldsymbol{\theta} | \mathbf{y}) &\propto \pi(\boldsymbol{\theta}) \times \pi(\mathbf{z} | \boldsymbol{\theta}_1) \times \prod_{i=1}^N p(y_i | \mathbf{z}, \boldsymbol{\theta}_2) \\ &\propto \pi(\boldsymbol{\theta}) |\mathbf{Q}(\boldsymbol{\theta}_1)|^{1/2} \exp \left\{ -\frac{1}{2} \mathbf{z}^\top \mathbf{Q}(\boldsymbol{\theta}_1) \mathbf{z} + \sum_{i=1}^N \log p(y_i | z_i, \boldsymbol{\theta}_2) \right\}. \end{aligned} \quad (14.1)$$

Most models in INLA assume a conditional independence structure within the high-dimensional latent Gaussian field; that is,  $\mathbf{z}$  is a Gaussian Markov random field (GMRF) with a sparse precision matrix  $\mathbf{Q}(\boldsymbol{\theta}_1)$ . A second key assumption is that the dimension of  $\boldsymbol{\theta}$  is small, for instance,  $m < 15$ . These assumptions are essential for enabling fast approximate inference.

The main aim in INLA is to approximate the marginal posterior distributions  $\pi(z_i | \mathbf{y})$  and  $\pi(\theta_l | \mathbf{y})$ ,  $l = 1, 2, \dots, m$ . The posterior distribution of  $\boldsymbol{\theta}$  is

$$\pi(\boldsymbol{\theta} | \mathbf{y}) = \frac{p(\mathbf{y}, \boldsymbol{\theta})}{p(\mathbf{y})}$$

$$\begin{aligned}
&= \frac{p(\mathbf{y}, \boldsymbol{\theta})}{p(\mathbf{y})} \times \frac{\pi(z | \boldsymbol{\theta}, \mathbf{y})}{\pi(z | \boldsymbol{\theta}, \mathbf{y})} \\
&= \frac{p(z, \boldsymbol{\theta}, \mathbf{y})}{p(\mathbf{y})\pi(z | \boldsymbol{\theta}, \mathbf{y})} \\
&\propto \frac{p(z, \boldsymbol{\theta}, \mathbf{y})}{\pi(z | \boldsymbol{\theta}, \mathbf{y})} \\
&\propto \frac{p(\mathbf{y} | z, \boldsymbol{\theta})\pi(z | \boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\pi(z | \boldsymbol{\theta}, \mathbf{y})}.
\end{aligned}$$

The numerator in the previous expression is easy to calculate (see Equation 14.1), but the denominator is generally not available in closed form and is difficult to compute. Thus, INLA approximates it at specific values  $\boldsymbol{\theta}_g$ ,

$$\pi_a(\boldsymbol{\theta}_g | \mathbf{y}) \propto \frac{p(\mathbf{y} | z, \boldsymbol{\theta}_g)\pi(z | \boldsymbol{\theta}_g)\pi(\boldsymbol{\theta}_g)}{\pi_{a,G}(z | \boldsymbol{\theta}_g, \mathbf{y})},$$

where  $\pi_{a,G}(z_i | \boldsymbol{\theta}_g, \mathbf{y})$  is a Gaussian approximation that matches the mode and covariance matrix of the full posterior  $\pi(z | \boldsymbol{\theta}, \mathbf{y})$ .

The approximation error of  $\pi_a(\boldsymbol{\theta} | \mathbf{y})$  is  $O(N^{-1})$  under standard conditions, meaning that the error, when multiplied by  $N$ , remains bounded [434, 392]. However, in many applications using INLA, the dimension of the latent Gaussian variables,  $P$ , increases with the sample size. In such cases, the error rate becomes  $O(P/N)$ . When  $P/N \rightarrow 1$ , which occurs in many models, the approximation error becomes  $O(1)$ : bounded, but potentially large.

Therefore, it is important to check the effective number of parameters, as the asymptotic error of  $\pi_a(\boldsymbol{\theta} | \mathbf{y})$  depends on the dimension of  $\mathbf{z}$ . According to [392], in most of their applications the effective number of parameters is small relative to the sample size in regions near the mode of  $\pi_a(\boldsymbol{\theta} | \mathbf{y})$ .

The marginal posterior  $\pi(z_i | \boldsymbol{\theta}, \mathbf{y})$  is more challenging to compute due to the potentially high dimension of  $\mathbf{z}$ . It may seem intuitive to use the Gaussian approximation  $\pi_{a,G}(z | \boldsymbol{\theta}, \mathbf{y})$ ; however, this is often not sufficiently accurate due to its lack of skewness. An alternative is to use the following expression

$$\begin{aligned}
\pi(z_i | \boldsymbol{\theta}, \mathbf{y}) &= \frac{p(z_i, \boldsymbol{\theta}, \mathbf{y})}{p(\boldsymbol{\theta}, \mathbf{y})} \\
&= \frac{p(z_i, \boldsymbol{\theta}, \mathbf{y})}{p(\boldsymbol{\theta}, \mathbf{y})} \times \frac{\pi(z | \boldsymbol{\theta}, \mathbf{y})}{\pi(z | \boldsymbol{\theta}, \mathbf{y})} \\
&= \frac{\pi(z | \boldsymbol{\theta}, \mathbf{y})}{p(z, \boldsymbol{\theta}, \mathbf{y})} \times p(z_i, \boldsymbol{\theta}, \mathbf{y}) \\
&= \frac{\pi(z | \boldsymbol{\theta}, \mathbf{y})}{\pi(z_{-i} | z_i, \boldsymbol{\theta}, \mathbf{y})} \\
&\propto \frac{p(\mathbf{y} | z, \boldsymbol{\theta})\pi(z | \boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\pi(z_{-i} | z_i, \boldsymbol{\theta}, \mathbf{y})},
\end{aligned}$$

where the second-to-last equality follows from the identity  $\pi(z_{-i} | z_i, \boldsymbol{\theta}, \mathbf{y}) = \frac{p(z, \boldsymbol{\theta}, \mathbf{y})}{p(z_i, \boldsymbol{\theta}, \mathbf{y})}$ , and  $-i$  denotes all elements of  $\mathbf{z}$  except the  $i$ -th.

We can approximate the denominator in the last expression using a simplified Laplace approximation, which corrects the Gaussian approximation for location and skewness via a Taylor series expansion about the mode of the Laplace approximation [392]. This follows the spirit of the approximations developed by [434, 435] for posterior moments and marginal densities, and is similar to those used in Chapter 10 when performing Bayesian model averaging (BMA) using the BIC information criterion. The discussion about the asymptotic error of the approximation  $\pi_a(z_i | \boldsymbol{\theta}, \mathbf{y})$  follows the same arguments as those for  $\pi_a(\boldsymbol{\theta} | \mathbf{y})$ .

We can obtain the marginal posterior distributions integrating with respect to  $\boldsymbol{\theta}$ ,

$$\begin{aligned}\pi(z_i | \mathbf{y}) &= \int_{\Theta} \pi(z_i | \boldsymbol{\theta}, \mathbf{y}) \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} \\ \pi(\theta_l | \mathbf{y}) &= \int_{\Theta} \pi(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}_{-l}.\end{aligned}$$

This integrals are solve numerically using a smart grid around the mode of  $\pi_a(\boldsymbol{\theta} | \mathbf{y})$ . In particular,

$$\pi_a(z_i | \mathbf{y}) = \sum_{g=1}^G \pi_a(z_i | \boldsymbol{\theta}_g, \mathbf{y}) \pi_a(\boldsymbol{\theta}_g | \mathbf{y}) \Delta_g,$$

where  $\pi_a(z_i | \boldsymbol{\theta}_g, \mathbf{y})$  is the approximation of  $\pi(z_i | \boldsymbol{\theta}, \mathbf{y})$  evaluated at  $\boldsymbol{\theta}_g$ . Then, we have the sum over the values of  $\boldsymbol{\theta}$  with area weights  $\Delta_g$ . If all support points are equidistant, then  $\Delta_g = 1$ .

Algorithm A40 presents the INLA algorithm. Note that stages 2 and 3 correspond to the nested Laplace approximations, whereas stage 4 involves the integration step.

---

**Algorithm A40** Integrated nested Laplace approximations (INLA)

---

- 1: Obtain the mode of  $\pi_a(\boldsymbol{\theta} | \mathbf{y})$ , that is  $\boldsymbol{\theta}^*$ , by maximizing  $\log \pi_a(\boldsymbol{\theta} | \mathbf{y})$
  - 2: Compute  $\pi_a(\boldsymbol{\theta}_g | \mathbf{y})$  for a set of high density points  $\boldsymbol{\theta}_g$ ,  $g = 1, 2, \dots, G$
  - 3: Compute the approximation  $\pi_a(z_i | \boldsymbol{\theta}_g, \mathbf{y})$ ,  $g = 1, 2, \dots, G$
  - 4: Compute  $\pi_a(z_i | \mathbf{y})$  and  $\pi_a(\theta_l | \mathbf{y})$  using numerical integration
- 

Thus, INLA can be used in LGMs that satisfy [393, 300]:

1. There is conditional independence, that is,  $y_i \perp y_s | \eta_i, \boldsymbol{\theta}$ , such that  $p(\mathbf{y} | \mathbf{z}, \boldsymbol{\theta}) = \prod_{i=1}^N p(y_i | \eta_i, \boldsymbol{\theta})$ .
2. The dimension of  $\boldsymbol{\theta}$  is small, less than 15.
3.  $\mathbf{z}$  is a GMRF with sparse precision matrix.
4. The linear predictor depends linearly on the smooth unknown functions of covariates.

5. The inference is on the marginal posterior distributions  $\pi(z_i | \mathbf{y})$  and  $\pi(\theta_l | \mathbf{y})$ .

[392] also discuss how to approximate the marginal likelihood and compute the deviance information criterion [420] for model selection, as well as how to perform predictive analysis.

The point of departure for INLA is the class of latent Gaussian models (LGMs); consequently, discrete latent classes are not supported. In addition, since INLA is based on approximations around the mode, it may struggle with models exhibiting multimodality, as there is no global exploration of the parameter space.

Implementing INLA from scratch is a complex task [300]; thus, applications are generally limited to the models implemented in the **INLA** package in **R**.<sup>3</sup> However, there are new packages that implement specialized models, as well as recent developments that combine INLA with MCMC algorithms (see Table 2 in [300] for details).

#### **Example: Poisson model with unobserved heterogeneity**

Let's simulate the model  $Y_i \sim \text{Poisson}(\lambda_i)$ , where  $\lambda_i = \exp\{1 + x_i + \epsilon_i\}$ , with  $\epsilon_i \sim \mathcal{N}(0, 0.5^2)$  and  $x_i \sim \mathcal{N}(0, 1^2)$  for  $i = 1, 2, \dots, 10,000$ . Note that  $\epsilon_i$  represents the unobserved heterogeneity.

The following code demonstrates how to perform inference on this model using the **INLA** package. Keep in mind that INLA specifies Gaussian priors in terms of precision, which is the inverse of the variance.

We present the results obtained using the three approximation strategies available in INLA: Simplified Laplace (default), Gaussian, and full Laplace. In this example, the results are practically identical across methods (see Figures 14.4 and 14.5), and all 95% credible intervals contain the true population parameters. Despite the large sample size, INLA performs inference in a matter of seconds.

#### ***R code. Integrated Nested Laplace Approximations: Poisson model with unobserved heterogeneity***

```

1 # Install Rtools according to your R version
2 # Check Rtools is properly installed
3 # Install INLA
4 # install.packages("INLA", repos=cgetOption("repos"), INLA =
  https://inla.r-inla-download.org/R/stable"), dep=TRUE)
5 rm(list = ls()); set.seed(010101); library(INLA)

```

<sup>3</sup>Visit <https://www.r-inla.org/> for documentation.

***R code. Integrated Nested Laplace Approximations: Poisson model with unobserved heterogeneity***

---

```

1 n <- 10000; x <- rnorm(n, sd = 1); u <- rnorm(n, sd = 0.5)
2 intercept <- 1; beta <- 1; id <- 1:n
3 y <- rpois(n, lambda = exp(intercept + beta * x + u))
4 my.data <- data.frame(y, x, id)
5 formula <- y ~ 1 + x + f(id, model="iid")
6 inla.sla <- inla(formula, data = my.data, family = "poisson"
                  , control.compute=list(return.marginals.predictor=TRUE))
7 inla.ga <- inla(formula, data = my.data, family = "poisson",
                  control.inla = list(strategy = "gaussian", int.strategy
                  = "eb"), control.compute=list(return.marginals.
                  predictor=TRUE))
8 inla.la <- inla(formula, data = my.data, family = "poisson",
                  control.inla = list(strategy = "laplace", int.strategy
                  = "grid", dz=0.1, diff.logdens=20), control.compute=list
                  (return.marginals.predictor=TRUE))
9 summary(inla.sla)
10 marg_sla <- inla.sla$ marginals.fixed$x
11 marg_ga <- inla.ga$ marginals.fixed$x
12 marg_la <- inla.la$ marginals.fixed$x
13 plot(marg_sla, type = "l", col = "blue", lwd = 2, xlab =
       expression(beta[x]), ylab = "Density", main = "Posterior
       of slope under different INLA strategies")
14 lines(marg_ga, col = "green", lwd = 2, lty = 2)
15 lines(marg_la, col = "red", lwd = 2, lty = 3)
16 abline(v = 1, col = "black", lty = 4, lwd = 2)
17 legend("topright", legend = c("Simplified Laplace", "
       Gaussian", "Full Laplace", "True = 1"), col = c("blue",
       "green", "red", "black"), lwd = 2, lty = c(1, 2, 3, 4),
       bty = "n")
18 # Summary beta.sla
19 inla.qmarginal(c(0.025, 0.5, 0.975), marg_sla) # 95%
       credible interval
20 # Variance
21 marg.prec.sla <- inla.sla$ marginals.hyperpar[["Precision for
      id"]]
22 marg.prec.ga <- inla.ga$ marginals.hyperpar[["Precision for
      id"]]
23 marg.prec.la <- inla.la$ marginals.hyperpar[["Precision for
      id"]]
24 marg.var.sla <- inla.tmarginal(function(x) 1/x, marg.prec.
      sla)
25 marg.var.ga <- inla.tmarginal(function(x) 1/x, marg.prec.ga
      )
26 marg.var.la <- inla.tmarginal(function(x) 1/x, marg.prec.la
      )

```

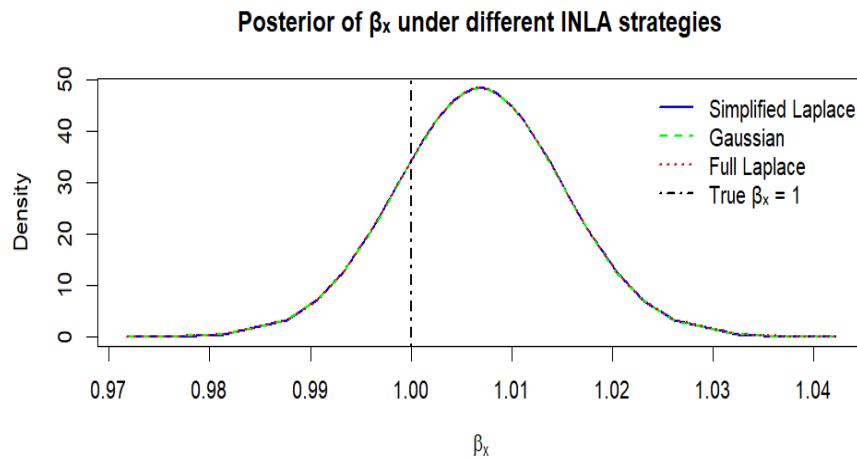
---

**R code. Integrated Nested Laplace Approximations: Poisson model with unobserved heterogeneity**

```

1 # Base plot
2 plot(marg.var.sla, type = "l", col = "blue", lwd = 2, xlab =
+   expression(sigma^2), ylab = "Density", main = "
+   Posterior of Random Effect Variance")
3 lines(marg.var.ga, col = "green", lwd = 2, lty = 2)
4 lines(marg.var.la, col = "red", lwd = 2, lty = 3)
5 abline(v = 0.25, col = "black", lty = 4, lwd = 2)
6 legend("topright", legend = c("Simplified Laplace", "
+   Gaussian", "Full Laplace", "True Variance (0.25)"), col
= c("blue", "green", "red", "black"), lwd = 2, lty = c
(1, 2, 3, 4), bty = "n")
7 # Summary variance sla
8 inla.qmarginal(c(0.025, 0.5, 0.975), marg.var.sla) # 95%
+   credible interval

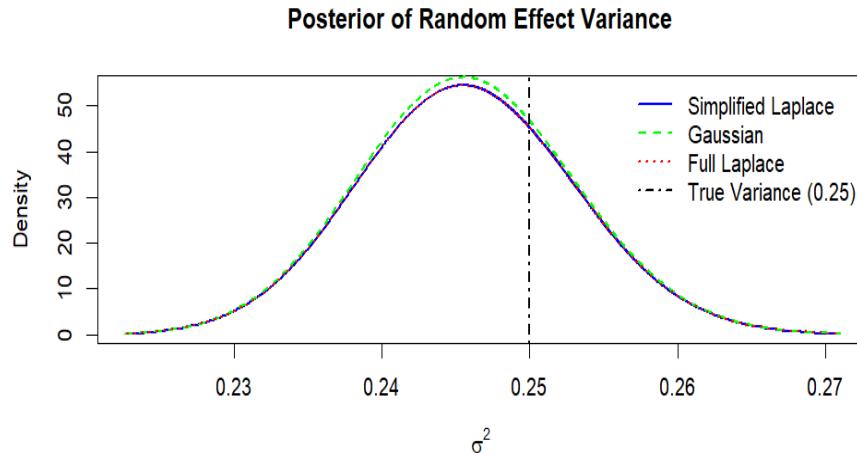
```

**FIGURE 14.4**

Density estimates of slope parameter: Three approximations in INLA.

**Example: Spatial econometrics model**

The starting point of spatial econometrics is the *contiguity* or *adjacency* matrix  $W$ , which defines which spatial polygons (regions) are considered

**FIGURE 14.5**

Density estimates of variance: Three approximations in INLA.

neighbors. This is an  $N \times N$  matrix, where each row and column corresponds to a spatial polygon, and a non-zero element indicates that two polygons are neighbors. By construction, the main diagonal is equal to zero, meaning that no polygon is a neighbor to itself. Given its structure, the contiguity matrix is sparse.

There are various ways to define contiguity between spatial units, with two of the most commonly used criteria being the *queen* and *rook* criteria, inspired by chess. In the queen criterion, two spatial units are considered neighbors if they share any part of a boundary or a point. In contrast, under the rook criterion, two units are considered neighbors only if they share a common boundary edge; they must touch along at least one side, and corner contact is not sufficient to establish a neighbor relationship.

However, users may apply any contiguity criteria they deem appropriate, depending on the context. For example, contiguity can be defined based on factors like the travel time by car between centroids or the proximity of main towns between municipalities. Typically, the contiguity matrix  $W$  is binary, with a value of 1 indicating that two regions are neighbors and 0 otherwise. Additionally, contiguity matrices are often row-standardized to facilitate the analysis of spatial stationarity.

[264, 47] present the most commonly used specifications in spatial econometrics, namely the spatial error model (SEM), the spatial autoregressive model (SAR), and the spatial Durbin model (SDM), which can, in turn, be combined to form the general nesting spatial (GNS) model [135, 353]. In par-

ticular, the SEM sets

$$\begin{aligned}\mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu} \\ \boldsymbol{\mu} &= \lambda\mathbf{W}\boldsymbol{\mu} + \boldsymbol{\epsilon},\end{aligned}$$

where  $\boldsymbol{\epsilon} \sim N(\mathbf{0}_N, \sigma^2 \mathbf{I}_N)$ . Then,  $\boldsymbol{\mu} = (\mathbf{I}_N - \rho\mathbf{W})^{-1}\boldsymbol{\epsilon}$ , which implies that

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

where  $\mathbf{e} \sim N(\mathbf{0}, \sigma^2(\mathbf{I}_N - \rho\mathbf{W})^{-1}(\mathbf{I}_N - \rho\mathbf{W}^\top)^{-1})$ . In this model  $\lambda$  controls the degree of spatial dependence.

The SAR sets

$$\begin{aligned}\mathbf{y} &= \rho\mathbf{W}\mathbf{y} + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu} \\ &= (\mathbf{I}_N - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},\end{aligned}$$

where  $\boldsymbol{\epsilon} = (\mathbf{I}_N - \rho\mathbf{W})^{-1}\boldsymbol{\mu}$ .

The SDM sets

$$\begin{aligned}\mathbf{y} &= \rho\mathbf{W}\mathbf{y} + \mathbf{X}\boldsymbol{\beta} + \mathbf{W}\mathbf{X}\boldsymbol{\delta} + \boldsymbol{\mu} \\ &= (\mathbf{I}_N - \rho\mathbf{W})^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{W}\mathbf{X}\boldsymbol{\delta}) + \boldsymbol{\epsilon},\end{aligned}$$

where  $\boldsymbol{\epsilon} = (\mathbf{I}_N - \rho\mathbf{W})^{-1}\boldsymbol{\mu}$ .

The degree of spatial dependence in the SAR and SDM models is determined by the parameter  $\rho$ . Note that the SDM model is similar to the SAR model, except that it also includes spatial lags of the regressors as additional covariates.

[328] and [20] show that a necessary condition for weak stationarity in spatial autoregressive processes with a row-standardized contiguity matrix is that the spatial autocorrelation coefficient must lie between  $1/\omega_{\min}$  and 1, where  $\omega_{\min}$  is the smallest (most negative) eigenvalue of the contiguity matrix.<sup>4</sup>

Note that these three models, conditional on the spatial parameter, are standard linear regressions, which can be easily estimated using INLA. [47] use the *INLABMA* package to estimate these models conditional on different values of the spatial parameters, and then perform Bayesian Model Averaging (BMA) using the resulting estimates. In particular, it is necessary to define a grid over the spatial parameters, perform Bayesian inference using INLA for each value in the grid, and then aggregate the posterior results using BMA.

We now perform Bayesian inference on a spatial econometric model using the *INLA* and *INLABMA* packages, based on the dataset provided by [361], who conducted a Bayesian analysis of electricity demand in the department of

---

<sup>4</sup>This is a necessary condition to ensure weak stationarity but is not sufficient due to edge and corner effects [190, ].

Antioquia (Colombia), accounting for spatial dependence between municipalities in the region using a Conditional Autoregressive (CAR) spatial model. In particular, we conduct inference using the following specification:

$$\begin{aligned}\log(\text{Electricity}_i) = & \beta_1 + \beta_2 \log(\text{Elect. price}_i) + \beta_3 \log(\text{Income}_i) \\ & + \beta_4 \log(\text{Subs. price}_i) + \mu_i,\end{aligned}$$

where  $\boldsymbol{\mu} = \lambda \mathbf{W} \boldsymbol{\mu} + \boldsymbol{\epsilon}$ ; *Electricity* is the average per capita annual consumption of electricity by individuals living in households of stratum one in each municipality of the department of Antioquia; *Elect. price* is the average price of electricity (per kWh); *Income* is the average per capita annual income; and *Subs. price* is the average price of an electricity substitute (see [361] for details).

The following code illustrates how to carry out a spatial econometric analysis. First, we download the files needed to plot the maps and construct the contiguity matrix. These files are available in the folder *DataApp/Antioquia* of our GitHub repository: <https://github.com/besmarter/BSTApp>. The initial part of the code demonstrates how to download the *GitHub* repository and extract the necessary files for mapping.

Figure 14.6 displays the average electricity consumption per municipality. We observe the presence of spatial clusters, particularly in the northwestern region, which corresponds to a low-altitude area along the Caribbean coast.

The second part of the code constructs the contiguity matrix using the queen criterion. Figure 14.7 illustrates the spatial links between municipalities. Following this, we estimate a standard Ordinary Least Squares (OLS) regression and conduct spatial autocorrelation tests on the residuals. The null hypothesis of both the global and local Moran's I tests is the absence of spatial autocorrelation in the OLS residuals. We reject the null hypothesis.

Finally, we perform Bayesian inference over a predefined grid of values for the spatial coefficient and apply Bayesian Model Averaging (BMA) using the *INLABMA* package. Figure 14.8 shows the BMA posterior distribution of the spatial coefficient in the SEM, which indicates the presence of spatial dependence. Figure 14.9 displays the posterior density of the own-price elasticity of electricity demand.

### *R code. Spatial econometric analysis: Spatial error model using INLA*

```
1 rm(list = ls()); set.seed(010101)
2 library(sf); library(tmap); library(classInt)
3 library(RColorBrewer); library(spdep); library(parallel)
4 library(INLABMA); library(sf); library(INLA)
```

**R code. Spatial econometric analysis: Spatial error model using INLA**

```

1 zip_url <- "https://github.com/BEsmarter-consultancy/BSTApp/
  archive/refs/heads/master.zip"
2 temp_zip <- tempfile(fileext = ".zip")
3 download.file(zip_url, temp_zip, mode = "wb")
4 temp_dir <- tempdir()
5 unzip(temp_zip, exdir = temp_dir)
6 antioquia_path <- file.path(temp_dir, "BSTApp-master", "
  DataApp", "Antioquia")
7 list.files(antioquia_path)
8 antioquia_path <- file.path(temp_dir, "BSTApp-master", "
  DataApp", "Antioquia")
9 shp_file <- file.path(antioquia_path, "Antioquia.shp")
10 antioquia <- st_read(shp_file)
11 print(antioquia)
12 plot(st_geometry(antioquia), main = "Map of Antioquia")
13 interval <- classIntervals(antioquia$CONS_OLD, 5, style = "
  quantile")
14 antioquia$cons_class <- cut(antioquia$CONS_OLD, breaks =
  interval$brks, include.lowest = TRUE)
15 plotcolors <- brewer.pal(5, "Reds")
16 tmap_mode("plot")
17 tm_shape(antioquia) +
18 tm_fill(fill = "cons_class", fill.scale = tm_scale(values =
  plotcolors), fill.legend = tm_legend(title = "
    Electricity consumption")) + tm_borders(col = "grey90")
  + tm_compass(type = "8star", position = c("right", "top"))
  ) + tm_layout(legend.outside = TRUE)
19 nb_object <- poly2nb(antioquia, queen = TRUE)
20 centroids <- st_centroid(st_geometry(antioquia))
21 coords <- st_coordinates(centroids)
22 plot(st_geometry(antioquia), border = "grey")
23 plot(nb_object, coords, add = TRUE, col = "red")
24 attach(antioquia)
25 fform <- L_CONS_OLD ~ L_P_OLD + L_ING_US + L_P_SUST
26 RegOLS <- lm(fform)
27 summary(RegOLS)
28 res <- RegOLS$residuals
29 NBList <- nb2listw(nb_object, style = "B")
30 moran_mc<- moran.mc(res, listw = NBList, 10000)
31 LM<-localmoran(as.vector(res), NBList)
32 sum(LM[,5]<0.05)
33 # Bayesian estimation
34 zero.variance <- list(prec = list(initial = 25, fixed = TRUE
  ))
35 ant.mat <- nb2mat(nb_object)
36 bmsp <- as(ant.mat, "CsparseMatrix")
37 antioquia$idx <- 1:nrow(antioquia)
38 rrho1 <- seq(0.5, 0.95, len = 10)

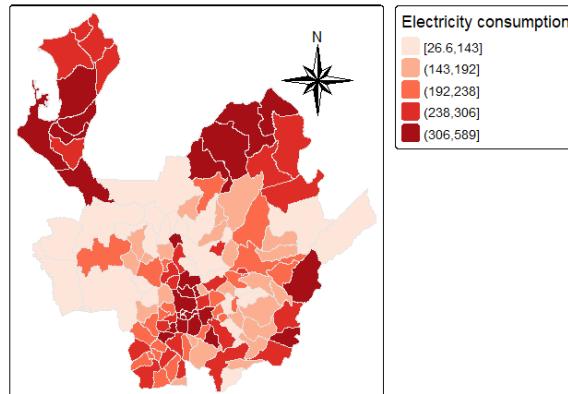
```

**R code. Spatial econometric analysis: Spatial error model using INLA**

```

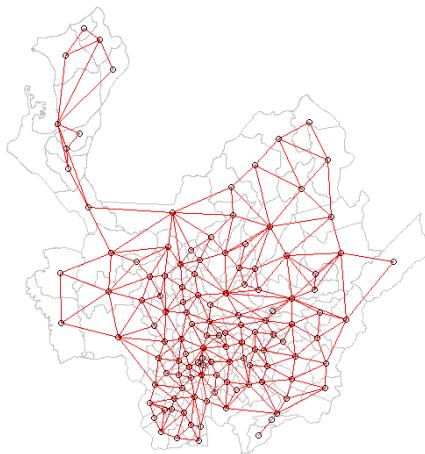
1 semmodels <- mclapply(rrho1, function(rho) {
2   sem.inla(fform, d = as.data.frame(antioquia), W = bmsp,
3   rho = rho, family = "gaussian", impacts = FALSE, control
4   .family = list(hyper = zero.variance), control.predictor
5   = list(compute = TRUE), control.compute = list(dic =
6   TRUE, cpo = TRUE), control.inla = list(print.joint.hyper
7   = TRUE))
8 })
9 bmasem <- INLABMA(semmodeles, rrho1, 0, impacts = FALSE)
10 #Display results
11 plot(bmasem$rho$marginal, type="l", col = "blue", lwd = 2,
12      xlab = expression(lambda), ylab = "Density", main =
13      "Spatial error model: Posterior spatial coefficient")
14 bmasem[["rho"]][["quantiles"]]
15 marg_sla <- bmasem[["marginals.fixed"]][["L_P_OLD"]]
16 plot(marg_sla, type = "l", col = "blue", lwd = 2, xlab =
17      expression(beta[x]), ylab = "Density", main = "Spatial
18      error model: Posterior price elasticity")
19 bmasem[["summary.fixed"]]

```

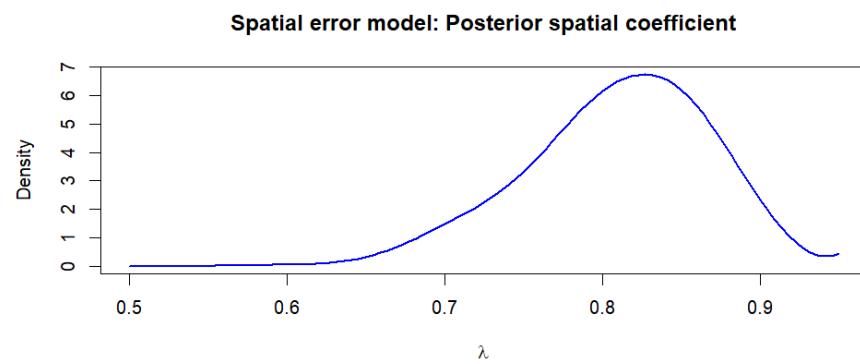


**FIGURE 14.6**

Per capita annual average electricity consumption: Stratum 1, Antioquia (Colombia).

**FIGURE 14.7**

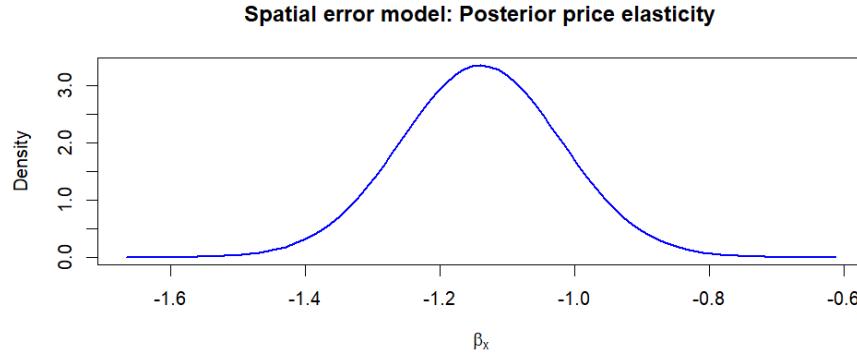
Contiguity map: Queen criterion, municipalities (Antioquia, Colombia).

**FIGURE 14.8**

BMA posterior density: Spatial error coefficient.

### 14.2.2 Variational Bayes

*Variational Bayes* (VB) is a method from machine learning [232, 443] that replaces  $\pi(\theta | \mathbf{y})$  with an approximation obtained through optimization using



**FIGURE 14.9**  
BMA posterior density: Own-price elasticity of electricity.

the calculus of variations, hence the name *variational* Bayes. This approach is useful when the posterior distribution is complex (e.g., multimodal) or when the parameter space is high-dimensional, making MCMC or IS algorithms computationally expensive.

The goal in VB is to approximate the posterior distribution using a distribution  $q(\boldsymbol{\theta})$  from a variational family  $\mathcal{Q}$  —a class of distributions that is computationally convenient yet flexible enough to closely approximate the true posterior [50]-. The distribution  $q$  is called the *variational approximation* to the posterior, and a particular  $q$  in  $\mathcal{Q}$  is defined by a specific set of *variational parameters*. Typically, this approximation is obtained by minimizing the Kullback–Leibler (KL) divergence between  $q(\boldsymbol{\theta})$  and  $\pi(\boldsymbol{\theta} | \mathbf{y})$ .

$$q^*(\boldsymbol{\theta}) := \arg \min_{q \in \mathcal{Q}} \text{KL}(q(\boldsymbol{\theta}) || \pi(\boldsymbol{\theta} | \mathbf{y})), \quad (14.2)$$

where  $\text{KL}(q(\boldsymbol{\theta}) || \pi(\boldsymbol{\theta} | \mathbf{y})) = \mathbb{E}_q \left[ \log \left( \frac{q(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta} | \mathbf{y})} \right) \right]$ .<sup>5</sup>

Note that in relatively complex models, the optimization in Equation 14.2 is not computable because it depends on the marginal likelihood  $p(\mathbf{y})$ , which is typically unknown due to the intractability of the integral involved. However, there is a solution to this problem. Let's see:

$$\begin{aligned} \log(p(\mathbf{y})) &= \log \left( \int_{\Theta} p(\mathbf{y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \right) \\ &= \log \left( \int_{\Theta} p(\mathbf{y}, \boldsymbol{\theta}) d\boldsymbol{\theta} \right) \end{aligned}$$

---

<sup>5</sup>The Kullback–Leibler (KL) divergence, also known as relative entropy, is non-negative: it equals 0 when  $q(\boldsymbol{\theta}) = \pi(\boldsymbol{\theta} | \mathbf{y})$ , and is positive otherwise. However, it does not satisfy the triangle inequality and is not symmetric, that is,  $\text{KL}(q(\boldsymbol{\theta}) || \pi(\boldsymbol{\theta} | \mathbf{y})) \neq \text{KL}(\pi(\boldsymbol{\theta} | \mathbf{y}) || q(\boldsymbol{\theta}))$ . Therefore, it is not a true distance (metric).

$$\begin{aligned}
&= \log \left( \int_{\Theta} \frac{p(\mathbf{y}, \boldsymbol{\theta})}{q(\boldsymbol{\theta})} q(\boldsymbol{\theta}) d\boldsymbol{\theta} \right) \\
&= \log \mathbb{E}_q \left( \frac{p(\mathbf{y}, \boldsymbol{\theta})}{q(\boldsymbol{\theta})} \right) \\
&\geq \mathbb{E}_q \log \left( \frac{p(\mathbf{y}, \boldsymbol{\theta})}{q(\boldsymbol{\theta})} \right) \\
&= \mathbb{E}_q \log(p(\mathbf{y}, \boldsymbol{\theta})) - \mathbb{E}_q \log(q(\boldsymbol{\theta})) \\
&= \text{ELBO}(q(\boldsymbol{\theta})),
\end{aligned}$$

where the inequality follows from Jensen's inequality, since  $\log(\cdot)$  is concave. The last term is the *evidence lower bound* (ELBO), which serves as a lower bound for the marginal likelihood. Note that the gap between the marginal likelihood and the ELBO is given by

$$\begin{aligned}
\log(p(\mathbf{y})) - \text{ELBO}(q(\boldsymbol{\theta})) &= \log(p(\mathbf{y})) - \mathbb{E}_q \log(p(\mathbf{y}, \boldsymbol{\theta})) + \mathbb{E}_q \log(q(\boldsymbol{\theta})) \\
&= \mathbb{E}_q \left( \log(q(\boldsymbol{\theta})) - \log \left( \frac{p(\mathbf{y}, \boldsymbol{\theta})}{p(\mathbf{y})} \right) \right) \\
&= \mathbb{E}_q \left( \log(q(\boldsymbol{\theta})) - \log \left( \frac{p(\mathbf{y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta})}{p(\mathbf{y})} \right) \right) \\
&= \mathbb{E}_q (\log(q(\boldsymbol{\theta})) - \log(\pi(\boldsymbol{\theta} | \mathbf{y}))) \\
&= \text{KL}(q(\boldsymbol{\theta}) || \pi(\boldsymbol{\theta} | \mathbf{y})).
\end{aligned}$$

Then,

$$\log(p(\mathbf{y})) = \text{KL}(q(\boldsymbol{\theta}) || \pi(\boldsymbol{\theta} | \mathbf{y})) + \text{ELBO}(q(\boldsymbol{\theta})),$$

which implies that maximizing the ELBO with respect to  $q(\boldsymbol{\theta})$  is equivalent to minimizing the KL divergence, since  $\log(p(\mathbf{y}))$  does not depend on  $q(\boldsymbol{\theta})$ . This avoids the need to compute the marginal likelihood and, consequently, makes the variational problem easier to solve. In addition, it provides a lower bound for the marginal likelihood, which can potentially be used for model selection. Thus, solving problem 14.2 is equivalent to solving

$$q^*(\boldsymbol{\theta}) := \arg \max_{q \in \mathcal{Q}} \text{ELBO}(q(\boldsymbol{\theta})). \quad (14.3)$$

Note that the ELBO can be also expressed as

$$\begin{aligned}
\text{ELBO}(q(\boldsymbol{\theta})) &= \mathbb{E}_q \log(p(\mathbf{y}, \boldsymbol{\theta})) - \mathbb{E}_q \log(q(\boldsymbol{\theta})) \\
&= \mathbb{E}_q \log(p(\mathbf{y} | \boldsymbol{\theta})) + \mathbb{E}_q \log(\pi(\boldsymbol{\theta})) - \mathbb{E}_q \log(q(\boldsymbol{\theta})) \\
&= \mathbb{E}_q \log(p(\mathbf{y} | \boldsymbol{\theta})) - \text{KL}(q(\boldsymbol{\theta}) || \pi(\boldsymbol{\theta})).
\end{aligned}$$

This means that when maximizing the ELBO, we seek the distribution  $q(\boldsymbol{\theta})$  that both maximizes the likelihood and minimizes the KL divergence between the variational distribution and the prior. In other words, we aim to strike

a balance between the prior and the likelihood, which aligns with the core principle of Bayesian inference deriving the posterior distribution.

The most common approach for specifying  $q(\boldsymbol{\theta})$  is to assume independence across blocks of  $\boldsymbol{\theta}$ , i.e.,  $q(\boldsymbol{\theta}) = \prod_{l=1}^K q_l(\boldsymbol{\theta}_l)$ . This is known as the *mean-field variational family*, a term that originates from statistical physics. Each  $q_l(\boldsymbol{\theta}_l)$  is parameterized by a set of *variational parameters*, and optimization is performed with respect to these parameters. Note that the mean-field approximation does not capture dependencies between parameters, although it can approximate the marginal distributions.

Let us now decompose the ELBO under the mean-field variational family [323],

$$\begin{aligned}
\text{ELBO}(q(\boldsymbol{\theta})) &= \mathbb{E}_q \log(p(\mathbf{y}, \boldsymbol{\theta})) - \mathbb{E}_q \log(q(\boldsymbol{\theta})) \\
&= \int_{\mathbb{R}^K} \log(p(\mathbf{y}, \boldsymbol{\theta})) q(\boldsymbol{\theta}) d\boldsymbol{\theta} - \int_{\mathbb{R}^K} \log(q(\boldsymbol{\theta})) q(\boldsymbol{\theta}) d\boldsymbol{\theta} \\
&= \int_{\mathbb{R}^K} \log(p(\mathbf{y}, \boldsymbol{\theta})) \prod_{l=1}^K q_l(\boldsymbol{\theta}_l) d\boldsymbol{\theta}_l - \int_{\mathbb{R}^K} \log\left(\prod_{l=1}^K q_l(\boldsymbol{\theta}_l)\right) \prod_{l=1}^K q_l(\boldsymbol{\theta}_l) d\boldsymbol{\theta}_l \\
&= \int_{\mathbb{R}} \underbrace{\left( \int_{\mathbb{R}^{K-1}} \log(p(\mathbf{y}, \boldsymbol{\theta})) \prod_{l \neq k} q_l(\boldsymbol{\theta}_l) d\boldsymbol{\theta}_l \right)}_{\mathbb{E}_{-k}(\log(p(\mathbf{y}, \boldsymbol{\theta})))} q_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k \\
&\quad - \int_{\mathbb{R}} \log(q_k(\boldsymbol{\theta}_k)) q_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k - \sum_{l \neq k} \int_{\mathbb{R}} \log(q_l(\boldsymbol{\theta}_l)) q_l(\boldsymbol{\theta}_l) d\boldsymbol{\theta}_l \\
&= \int_{\mathbb{R}} \log\{\exp(\mathbb{E}_{-k}(\log(p(\mathbf{y}, \boldsymbol{\theta}))))\} q_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k \\
&\quad - \int_{\mathbb{R}} \log(q_k(\boldsymbol{\theta}_k)) q_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k - \sum_{l \neq k} \int_{\mathbb{R}} \log(q_l(\boldsymbol{\theta}_l)) q_l(\boldsymbol{\theta}_l) d\boldsymbol{\theta}_l \\
&= \int_{\mathbb{R}} \log\left\{\frac{\exp(\mathbb{E}_{-k}(\log(p(\mathbf{y}, \boldsymbol{\theta}))))}{q_k(\boldsymbol{\theta}_k)}\right\} q_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k \\
&\quad - \sum_{l \neq k} \int_{\mathbb{R}} \log(q_l(\boldsymbol{\theta}_l)) q_l(\boldsymbol{\theta}_l) d\boldsymbol{\theta}_l \\
&= -\text{KL}(q_k(\boldsymbol{\theta}_k) || \exp(\mathbb{E}_{-k}(\log(p(\mathbf{y}, \boldsymbol{\theta})))) - \sum_{l \neq k} \int_{\mathbb{R}} \log(q_l(\boldsymbol{\theta}_l)) q_l(\boldsymbol{\theta}_l) d\boldsymbol{\theta}_l \\
&\leq -\sum_{l \neq k} \int_{\mathbb{R}} \log(q_l(\boldsymbol{\theta}_l)) q_l(\boldsymbol{\theta}_l) d\boldsymbol{\theta}_l,
\end{aligned}$$

where  $\mathbb{E}_{-k}$  denotes expectation with respect to the distribution  $\prod_{l \neq k} q_l(\boldsymbol{\theta}_l)$ .

Note that we maximize the ELBO when  $\text{KL}(q_k(\boldsymbol{\theta}_k) || \exp(\mathbb{E}_{-k}(\log(p(\mathbf{y}, \boldsymbol{\theta})))))$  equals 0, that is, when

$$q_k^*(\boldsymbol{\theta}_k) \propto \exp(\mathbb{E}_{-k}(\log(p(\mathbf{y}, \boldsymbol{\theta}))))$$

$$\begin{aligned}
&= \exp(\mathbb{E}_{-\bar{k}}(\log(p(\mathbf{y}, \boldsymbol{\theta}_{-\bar{k}}, \boldsymbol{\theta}_{\bar{k}})))) \\
&= \exp(\mathbb{E}_{-\bar{k}}(\log(p(\boldsymbol{\theta}_{\bar{k}} | \mathbf{y}, \boldsymbol{\theta}_{-\bar{k}})p(\mathbf{y}, \boldsymbol{\theta}_{-\bar{k}})))) \\
&\propto \exp(\mathbb{E}_{-\bar{k}}(\log(p(\boldsymbol{\theta}_{\bar{k}} | \mathbf{y}, \boldsymbol{\theta}_{-\bar{k}})))), \bar{k} = 1, 2, \dots, K.
\end{aligned} \tag{14.4}$$

Note the circular dependency inherent in  $q_k^*(\boldsymbol{\theta}_k)$ : it depends on  $q_{-\bar{k}}^*(\boldsymbol{\theta}_{-\bar{k}})$ . Therefore, this situation must be addressed algorithmically. One of the most common algorithms used to solve the optimization problem in 14.3 using 14.4 is the *coordinate ascent variational inference* (CAVI) algorithm [45]. The algorithm starts from an initial solution and iteratively cycles through each  $q_k^*(\boldsymbol{\theta}_k)$ , for  $k = 1, 2, \dots, K$ , updating each component in turn. Algorithm A41 outlines the basic structure of the CAVI algorithm.

---

**Algorithm A41** Variational Bayes: Coordinate ascent variational inference

---

```

1: Initialize the variational factors  $q_l^*(\boldsymbol{\theta}_l), l = 1, 2, \dots, K$ 
2: while ELBO >  $\epsilon$ ,  $\epsilon$  small do
3:   for  $l = 1, 2, \dots, L$  do
4:     Set  $q_l^*(\boldsymbol{\theta}_l) \propto \exp(\mathbb{E}_{-l}(\log(p(\boldsymbol{\theta}_l | \mathbf{y}, \boldsymbol{\theta}_{-l}))))$ 
5:   end for
6:   Compute ELBO( $q$ ) =  $\mathbb{E}_q \log(p(\mathbf{y}, \boldsymbol{\theta})) - \mathbb{E}_q \log(q(\boldsymbol{\theta}))$ 
7: end while
8: Return  $q(\boldsymbol{\theta})$ 

```

---

We should note that the CAVI algorithm is sensitive to the initial variational factors because it guarantees convergence to a local maximum, which may depend heavily on the initialization point. This issue can be mitigated by using multiple starting points and selecting the solution with the highest ELBO [49]. In addition, it is well known that VB tends to overestimate the precision of the posterior distribution, although it does not necessarily suffer in terms of accuracy [50].

An important feature of the method is that it is deterministic rather than stochastic; that is, it does not require approximating the posterior distribution via sampling from simpler distributions. This often makes VB faster than MCMC methods, which are inherently stochastic. Instead, VB solves a deterministic optimization problem to find the *best variational distribution* within a chosen family, typically from the exponential family, by optimizing the variational parameters to minimize the KL divergence from the posterior distribution. Once the *variational parameters* are obtained, we can sample from the *variational distribution* to conduct estimation, hypothesis testing, prediction, and other tasks.

A limitation of the CAVI algorithm is that it requires evaluation at each data point, making it non-scalable in large data settings. In such situations, we can use *stochastic variational inference* (SVI) [209], an algorithm that optimizes the ELBO using natural gradients combined with stochastic optimization. SVI is particularly effective when each complete conditional belongs to the exponential family (see Section 3.1), which includes most of the models

discussed in this book. It is especially useful for conditionally conjugate models that include *local* latent variables ( $\mathbf{z}_i$ ) associated with specific data points, and *global* parameters ( $\phi$ ) shared across the entire dataset. That is, we define  $\boldsymbol{\theta} = [\mathbf{z}^\top \phi^\top]^\top$ . An example is the probit model using data augmentation [428].

Given the global-local exchangeable structure, the joint distribution can be expressed as  $p(\phi, \mathbf{z}, \mathbf{y}) = \pi(\phi | \boldsymbol{\alpha}) \prod_{i=1}^N p(\mathbf{z}_i, \mathbf{y}_i | \phi)$ , where  $\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1^\top \ \boldsymbol{\alpha}_2]^\top$  represents the hyperparameters of the prior distribution of the global parameters  $\phi$ .

Assuming that the joint distribution  $p(\mathbf{z}_i, \mathbf{y}_i | \phi)$  is in the canonical form of the exponential family, and that each prior distribution is conjugate, the complete conditional distribution of the global parameters is also in the exponential family. The posterior parameters are given by  $\boldsymbol{\alpha}_n = [\boldsymbol{\alpha}_1^\top + \sum_{i=1}^N t(\mathbf{z}_i, \mathbf{y}_i)^\top, \ \boldsymbol{\alpha}_2 + N]^\top$ , where  $t(\mathbf{z}_i, \mathbf{y}_i)$  is a sufficient statistic of  $\mathbf{z}_i$  and  $\mathbf{y}_i$  (see Section 3.2).

Additionally, the structure of the model implies that  $\mathbf{z}_i$  is independent of  $\mathbf{z}_{-i}$  and  $\mathbf{y}_{-i}$  given  $\mathbf{y}_i$  and  $\phi$ . Thus, if  $p(\mathbf{z}_i | \mathbf{y}_i, \phi)$  is in the canonical form of the exponential family, then the local variational update is given by  $\boldsymbol{\psi}_i = \mathbb{E}_{\xi}[\boldsymbol{\eta}(\phi, \mathbf{y}_i)]$ , where the parameter set  $\boldsymbol{\eta}(\phi, \mathbf{y}_i)$  is a function of the conditional set,  $\xi$ , and  $\boldsymbol{\psi}_i$  are the variational parameters of the variational approximations  $q(\phi | \xi)$  and  $q(\mathbf{z}_i | \boldsymbol{\psi}_i)$  for the posterior distributions of  $\phi$  and  $\mathbf{z}_i$ , respectively. The global variational update is given by  $\xi = [\boldsymbol{\alpha}_1^\top + \sum_{i=1}^N \mathbb{E}_{\boldsymbol{\psi}_i} t(\mathbf{z}_i, \mathbf{y}_i)^\top, \ \boldsymbol{\alpha}_2 + N]^\top$ .

SVI focuses on the global parameters, updating them using the ELBO natural gradients with respect to  $\xi$ , where natural gradients are the usual gradients premultiplied by the inverse covariance matrix of the sufficient statistic. These gradients are easily calculated in exponential families. In particular, the updates are given by

$$\xi_t = \xi_{t-1} + \epsilon_t g(\xi_{t-1}),$$

where  $\epsilon_t$  is the step size, and  $g(\xi_t) = \mathbb{E}_{\boldsymbol{\psi}_i}[\boldsymbol{\alpha}_n - \xi_t]$  is the natural gradient of the global variational parameters [50]. Consequently,

$$\xi_t = (1 - \epsilon_t)\xi_{t-1} + \epsilon_t \mathbb{E}_{\boldsymbol{\psi}}[\boldsymbol{\alpha}_n].$$

However, calculating this update requires using the entire dataset, which is computationally burdensome. Therefore, we should use stochastic optimization, which follows noisy but cheap-to-compute unbiased gradients to optimize the function. The key idea is to construct the natural gradient using only one random draw from the dataset, and then scale it (by multiplying it with the sample size) to approximate the sample information:

$$\begin{aligned}\hat{\xi} &= \boldsymbol{\alpha} + N(\mathbb{E}_{\boldsymbol{\psi}_i^*}[t(\mathbf{z}_i, \mathbf{y}_i)]^\top, 1)^\top \\ \xi_t &= (1 - \epsilon_t)\xi_{t-1} + \epsilon_t \hat{\xi},\end{aligned}$$

where  $(\mathbf{z}_i, \mathbf{y}_i)$  is a random draw from the sample and its corresponding latent

variable. Finally, the step size schedule to update the global parameters are given by

$$\epsilon_t = t^{-\kappa}, \quad 0.5 < \kappa \leq 1.$$

This step size schedule satisfies the Robbins and Monro conditions necessary for stochastic optimization [370].

Algorithm A42 shows the stochastic variational inference algorithm.

---

**Algorithm A42** Variational Bayes: Stochastic variational inference

---

- 1: Initialize the variational global parameter  $\phi_0$
  - 2: Set the step size schedule  $\epsilon_t = t^{-\kappa}, \quad 0.5 < \kappa \leq 1$
  - 3: **while** TRUE **do**
  - 4:     Choose randomly a data point  $\mathbf{y}_i$  using a discrete uniform distribution,  $i \sim U(1, 2, \dots, N)$
  - 5:     Optimize its local variational parameters  $\psi_i^* = \mathbb{E}_{\xi_{t-1}}[\eta(\phi, \mathbf{y}_i)]$
  - 6:     Compute the coordinates updates as though  $\mathbf{y}_i$  were repeated  $N$  times,  

$$\hat{\xi} = \alpha + N(\mathbb{E}_{\psi_i^*}[t(\mathbf{z}_i, \mathbf{y}_i)]^\top, 1)^\top$$
  - 7:     Update the global variational parameters,  $\xi_t = (1 - \epsilon_t)\xi_{t-1} + \epsilon_t \hat{\xi}$
  - 8: **end while**
  - 9: Return  $q_\xi(\phi)$
- 

[468] analyzes the asymptotic properties of the VB posterior by decomposing the convergence rates of VB into the convergence rate of the true posterior and the approximation error induced by the variational family. These authors show that the VB posterior concentrates entirely in a neighborhood of the true posterior distribution. In addition, if the loss function is convex, there exists a point estimator that converges at the same rate. [468] also shows that, for specific cases such as sparse linear models, the concentration rate of the VB posterior is faster than the concentration rate of the exact posterior.

Variational Bayes (VB) shares some similarities with Gibbs sampling. In Gibbs sampling, we iteratively sample from the conditional posterior distributions, whereas in VB we iteratively update the variational parameters of the variational family. The former is stochastic, while the latter is deterministic, and consequently faster in complex models or large datasets. VB also bears resemblance to Expectation Propagation (EP): both are deterministic algorithms that approximate the posterior distribution by minimizing the Kullback-Leibler (KL) divergence. However, VB minimizes  $\text{KL}(q(\theta) \parallel \pi(\theta \mid \mathbf{y}))$ , whereas EP minimizes  $\text{KL}(\pi(\theta \mid \mathbf{y}) \parallel q(\theta))$ . As a result, VB tends to approximate the mode of the posterior by maximizing the ELBO, while EP focuses on matching the mean and variance through moment matching [45, 160]. Although EP often provides better uncertainty quantification, it tends to be less stable and does not scale well to large datasets.

VB also shares some conceptual features with the Expectation-Maximization (EM) algorithm. In both methods, there is an initial step involving the computation of expectations —used in VB to derive the variational

distributions—and an iterative step that maximizes a target function (the ELBO in VB, the expected complete-data log-likelihood in EM). However, EM yields point estimates, whereas VB yields full posterior approximations.

### Example: Linear regression

Let's perform variational Bayes inference in the linear regression model with conjugate family. In particular,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu},$$

where  $\boldsymbol{\mu} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ . This implies that:

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}).$$

The conjugate priors for the parameters are

$$\begin{aligned} \boldsymbol{\beta} \mid \sigma^2 &\sim N(\boldsymbol{\beta}_0, \sigma^2 \mathbf{B}_0), \\ \sigma^2 &\sim IG(\alpha_0/2, \delta_0/2). \end{aligned}$$

Then, the posterior distributions are  $\boldsymbol{\beta} \mid \sigma^2, \mathbf{y}, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \sigma^2 \mathbf{B}_n)$ , and  $\sigma^2 \mid \mathbf{y}, \mathbf{X} \sim IG(\alpha_n/2, \delta_n/2)$ , where  $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}$ ,  $\boldsymbol{\beta}_n = \mathbf{B}_n(\mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}})$ ,  $\hat{\boldsymbol{\beta}}$  is the MLE,  $\alpha_n = \alpha_0 + N$  and  $\delta_n = \delta_0 + \mathbf{y}^\top \mathbf{y} + \boldsymbol{\beta}_0^\top \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 - \boldsymbol{\beta}_n^\top \mathbf{B}_n^{-1} \boldsymbol{\beta}_n$  (see Section 3.3).

Let's use the mean-field variational family  $q(\boldsymbol{\beta}, \sigma^2) = q(\boldsymbol{\beta})q(\sigma^2) \approx \pi(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}, \mathbf{X})$ .<sup>6</sup> Then,

$$\begin{aligned} \log q^*(\boldsymbol{\beta}) &\propto \mathbb{E}_{\sigma^2} [\log p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2 \mid \mathbf{X})] \\ &= \mathbb{E}_{\sigma^2} [\log p(\mathbf{y} \mid \boldsymbol{\beta}, \sigma^2, \mathbf{X}) + \log \pi(\boldsymbol{\beta} \mid \sigma^2) + \log \pi(\sigma^2)] \\ &= \mathbb{E}_{\sigma^2} \left( -\frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \frac{K}{2} \log \sigma^2 \right. \\ &\quad \left. - \frac{1}{2\sigma^2} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^\top \mathbf{B}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \right) + c_1 \\ &= -0.5 \mathbb{E}_{\sigma^2} \left( \frac{1}{\sigma^2} \right) [(\boldsymbol{\beta} - \boldsymbol{\beta}_n)^\top \mathbf{B}_n^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_n)] + c_2. \end{aligned}$$

The last equality follows the same steps to obtain the posterior distribution of  $\boldsymbol{\beta}$ ,  $c_1$  and  $c_2$  are arbitrary constants. This expression implies that  $q^*(\boldsymbol{\beta})$  is  $N(\boldsymbol{\beta}_n, (\mathbb{E}_{\sigma^2} (\frac{1}{\sigma^2}))^{-1} \mathbf{B}_n)$ .

In addition,

$$\log q^*(\sigma^2) \propto \mathbb{E}_{\boldsymbol{\beta}} [\log p(\mathbf{y}, \boldsymbol{\beta}, \sigma^2 \mid \mathbf{X})]$$

---

<sup>6</sup>We can also use the variational family  $q(\boldsymbol{\beta}, \sigma^2) = q(\boldsymbol{\beta} \mid \sigma^2)q(\sigma^2)$ , which may be more straightforward to manipulate.

$$\begin{aligned}
&= \mathbb{E}_{\beta}[\log p(\mathbf{y} | \beta, \sigma^2, \mathbf{X}) + \log \pi(\beta | \sigma^2) + \log \pi(\sigma^2)] \\
&= \mathbb{E}_{\beta} \left( -\frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^{\top} (\mathbf{y} - \mathbf{X}\beta) - \frac{K}{2} \log \sigma^2 \right. \\
&\quad \left. - \frac{1}{2\sigma^2} (\beta - \beta_0)^{\top} \mathbf{B}_0^{-1} (\beta - \beta_0) - (\alpha_0/2 + 1) \log \sigma^2 - \frac{\delta_0}{2\sigma^2} \right) + c_1 \\
&= -\mathbb{E}_{\beta} \left[ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^{\top} (\mathbf{y} - \mathbf{X}\beta) + (\beta - \beta_0)^{\top} \mathbf{B}_0^{-1} (\beta - \beta_0) + \delta_0 \right] \\
&\quad - \left( \frac{N + K + \alpha_0}{2} + 1 \right) \log \sigma^2 + c_1.
\end{aligned}$$

This means that  $q^*(\sigma^2)$  is  $IG(\alpha_n/2, \delta_n/2)$ , where  $\alpha_n = N + K + \alpha_0$ , and

$$\begin{aligned}
\delta_n &= \mathbb{E}_{\beta} [(\mathbf{y} - \mathbf{X}\beta)^{\top} (\mathbf{y} - \mathbf{X}\beta) + (\beta - \beta_0)^{\top} \mathbf{B}_0^{-1} (\beta - \beta_0) + \delta_0] \\
&= \mathbb{E}_{\beta} [(\beta - \beta_n)^{\top} \mathbf{B}_n^{-1} (\beta - \beta_n)] - \beta_n^{\top} \mathbf{B}_n^{-1} \beta_n + \mathbf{y}^{\top} \mathbf{y} + \beta_0^{\top} \mathbf{B}_0^{-1} \beta_0 + \delta_0.
\end{aligned}$$

This result implies that  $\mathbb{E}_{\sigma^2} \left( \frac{1}{\sigma^2} \right) = \alpha_n/\delta_n$ , since the inverse of a gamma-distributed random variable (in the rate parametrization) follows an inverse-gamma distribution. Therefore,  $\text{Var}(\beta) = \frac{\delta_n}{\alpha_n} \mathbf{B}_n$ . Note that

$$\begin{aligned}
\mathbb{E}_{\beta} [(\beta - \beta_n)^{\top} \mathbf{B}_n^{-1} (\beta - \beta_n)] &= \text{tr} \{ \mathbb{E}_{\beta} [(\beta - \beta_n)^{\top} \mathbf{B}_n^{-1} (\beta - \beta_n)] \} \\
&= \mathbb{E}_{\beta} \{ \text{tr} [(\beta - \beta_n)^{\top} \mathbf{B}_n^{-1} (\beta - \beta_n)] \} \\
&= \mathbb{E}_{\beta} \{ \text{tr} [(\beta - \beta_n)(\beta - \beta_n)^{\top} \mathbf{B}_n^{-1}] \} \\
&= \text{tr} \{ \mathbb{E}_{\beta} [(\beta - \beta_n)(\beta - \beta_n)^{\top}] \mathbf{B}_n^{-1} \} \\
&= \text{tr} \{ \text{Var}(\beta) \mathbf{B}_n^{-1} \},
\end{aligned}$$

where we use that the trace of a scalar is the scalar itself, that expectation and trace can be interchanged since both are linear operators, and that the trace operator is invariant under cyclic permutations. Then,

$$\begin{aligned}
\delta_n &= \text{tr} \{ \text{Var}(\beta) \mathbf{B}_n^{-1} \} - \beta_n^{\top} \mathbf{B}_n^{-1} \beta_n + \mathbf{y}^{\top} \mathbf{y} + \beta_0^{\top} \mathbf{B}_0^{-1} \beta_0 + \delta_0 \\
&= \left( \frac{\alpha_0 + N + K}{\alpha_0 + K} \right) \left( -\beta_n^{\top} \mathbf{B}_n^{-1} \beta_n + \mathbf{y}^{\top} \mathbf{y} + \beta_0^{\top} \mathbf{B}_0^{-1} \beta_0 + \delta_0 \right),
\end{aligned}$$

where the last equality follows from  $\text{tr} \{ \text{Var}(\beta) \mathbf{B}_n^{-1} \} = \delta_n/\alpha_n \text{tr} \{ \mathbf{I}_K \}$ .

In addition (Exercise 4),

$$\begin{aligned}
\text{ELBO}(q(\beta, \sigma^2)) &= \mathbb{E}_{\beta, \sigma^2} [\log p(\mathbf{y}, \beta, \sigma^2 | \mathbf{X})] - \mathbb{E}_{\beta, \sigma^2} [\log q(\beta, \sigma^2)] \\
&= -\frac{N}{2} \log(2\pi) + \frac{\alpha_0}{2} \log(\delta_0/2) - \frac{\alpha_n}{2} \log(\delta_n/2) - 0.5 \log |\mathbf{B}_0| \\
&\quad + 0.5 \log |\mathbf{B}_n| - \log \Gamma(\alpha_0/2) + \log \Gamma(\alpha_n/2) \\
&\quad - K/2 \log(\alpha_n/\delta_n) + 0.5 \text{tr} \{ \text{Var}(\beta) \mathbf{B}_n^{-1} \}.
\end{aligned}$$

Note that the first two lines of the ELBO share the same structure as the log marginal likelihood  $p(\mathbf{y})$  in Section 3.3.

The following code presents a simulation setting with a sample size of 500 and two regressors drawn from standard normal distributions. The population parameters are set to 1, and we use non-informative priors with 10,000 posterior draws.

First, we perform VB inference using the *Laplace'sDemon* package, and then implement it from scratch. Although we have analytical solutions in this setting, we apply our own CAVI algorithm to assess its performance. We also compare the results with those from the Gibbs sampler, the marginal likelihood, and the ELBO.

In general, all calculations seem to perform well: the 95% credible intervals contain the population parameters, and the posterior means are very close to them.

### R code. Variational Bayes: Linear regression model

```

1 set.seed(010101)
2 library(Laplace'sDemon)
3 N <- 500; K <- 2
4 sig2 <- 1; B <- rep(1, K + 1)
5 X <- cbind(1, matrix(rnorm(N*K), N, K))
6 e <- rnorm(N, 0, sig2^0.5)
7 y <- X%*%B + e
8 ##### Data List Preparation #####
9 mon.names <- "mu[1]"
10 parm.names <- as.parm.names(list(beta=rep(0,K + 1), sigma2
11 =0))
12 pos.beta <- grep("beta", parm.names)
13 pos.sig2 <- grep("sigma2", parm.names)
14 PGF <- function(Data) {
15   beta <- rnorm(Data$K)
16   sigma2 <- runif(1)
17   return(c(beta, sigma2))
18 }
19 MyData <- list(K=K, PGF=PGF, X=X, mon.names=mon.names,
20 parm.names=parm.names, pos.beta=pos.beta, pos.sig2=pos.
21 sigma2, y=y)
22 ##### Model Specification #####
23 b0 <- 0; B0 <- 1000; a0 <- 0.01; d0 <- 0.01

```

***R code. Variational Bayes: Linear regression model***

---

```

1 Model <- function(parm, Data){
2   ### Parameters
3   beta <- parm[Data$pos.beta]
4   sigma2 <- interval(parm[Data$pos.sigma2], 1e-100, Inf)
5   parm[Data$pos.sigma2] <- sigma2
6   ### Log-Prior
7   beta.prior <- sum(dnormv(beta, b0, B0, log=TRUE))
8   sigma2.prior <- dinvgamma(sigma2, a0/2, d0/2, log=TRUE)
9   ### Log-Likelihood
10  mu <- tcrossprod(Data$X, t(beta))
11  LL <- sum(dnorm(Data$y, mu, sigma2^0.5, log=TRUE))
12  ### Log-Posterior
13  LP <- LL + beta.prior + sigma2.prior
14  Modelout <- list(LP=LP, Dev=-2*LL, Monitor=mu[1],
15    yhat=rnorm(length(mu), mu, sigma2^0.5), parm=parm)
16  return(Modelout)
17 }
18 Initial.Values <- rep(0,K+2); S <- 10000
19 Fit <- VariationalBayes(Model, Initial.Values, Data=MyData,
20   Covar=NULL,
21   Iterations=S, Method="Salimans2", Stop.Tolerance=1e-2, CPUs
22   =1)
23 print(Fit)
24 PosteriorChecks(Fit)
25 caterpillar.plot(Fit, Params="beta")
26 plot(Fit, MyData, PDF=FALSE)
27 Pred <- predict(Fit, Model, MyData, CPUs=1)
28 summary(Pred, Discrep="Chi-Square")
29 ##### MCMC #####
30 # Posterior parameters
31 b0 <- rep(b0, K+1); B0 <- B0*diag(K+1)
32 bhat <- solve(t(X) %*% X) %*% t(X) %*% y
33 Bn <- as.matrix(Matrix::forceSymmetric(solve(solve(B0) + t(X)
34   ) %*% X)))
35 bn <- Bn %*% (solve(B0) %*% b0 + t(X) %*% X %*% bhat)
36 dn <- as.numeric(d0 + t(y) %*% y + t(b0) %*% solve(B0) %*% b0 - t(bn) %
37   *%solve(Bn) %*% bn)
38 an <- a0 + N
39 # Posterior draws
40 sig2 <- MCMCpack::rinvgamma(S, an/2, dn/2)
41 summary(coda::mcmc(sig2))
42 Betas <- t(sapply(1:S, function(s){MASS::mvrnorm(1, bn, sig2
43   [s]*Bn)}))
44 summary(coda::mcmc(Betas))

```

---

### *R code. Variational Bayes: Linear regression model*

---

```

1 ##### VB from scratch #####
2 dnVB <- ((a0+N+K)/(a0+N))*dn; anVB <- a0 + N + K
3 BnVB <- Bn; bnVB <- bn
4 sig2VB <- MCMCpack::rinvgamma(S,anVB/2,dnVB/2)
5 BetasVB <- MASS::mvrnorm(S, mu = bnVB, Sigma = (dn/an)*BnVB)
6 summary(coda::mcmc(sig2VB)); summary(coda::mcmc(BetasVB))
7 ELBO <- -N/2*log(2*pi) + a0/2*log(d0/2) - anVB/2*log(dnVB/2)
   - 0.5*log(det(B0)) + 0.5*log(det(BnVB)) - lgamma(a0/2)
   + lgamma(anVB/2) - K/2*log(anVB/dnVB) + K/2
8 LogMarLik <- -N/2*log(2*pi) + a0/2*log(d0/2) - an/2*log(dn/
   2) - 0.5*log(det(B0)) + 0.5*log(det(Bn)) - lgamma(a0/2)
   + lgamma(an/2)
9 ELBO; LogMarLik; ELBO <- LogMarLik
10 # CAVI
11 ELBOfunc <- function(d,B){
12   ELBOi <- -N/2*log(2*pi) + a0/2*log(d0/2) - anVB/2*log(d/2)
   - 0.5*log(det(B0)) + 0.5*log(det(B)) - lgamma(a0/2) +
   lgamma(anVB/2) - K/2*log(anVB/d) + 0.5*(anVB/d)*sum(diag
   (B%*%solve(Bn)))
13   return(ELBOi)
14 }
15 d <- 100; B <- diag(K)
16 Esig2inv <- anVB/d; ELBOs <- rep(-Inf, S); epsilon <- 1e-5
17 for(s in 2:S){
18   B <- Esig2inv^(-1)*Bn
19   EbQb <- sum(diag(B%*%solve(Bn))) - t(bn)%*%solve(Bn)%*%bn
   + t(y)%*%y + t(b0)%*%solve(B0)%*%b0
20   d <- EbQb + d0
21   Esig2inv <- as.numeric(anVB/d)
22   ELBOs[s] <- ELBOfunc(d = d, B = B)
23   if (ELBOs[s] < ELBOs[s - 1]) { message("Lower bound
     decreases!\n")}
24   # Check for convergence
25   if (ELBOs[s] - ELBOs[s - 1] < epsilon) { break }
26   # Check if VB converged in the given maximum iterations
27   if (s == S) {warning("VB did not converge!\n")}
28 }
29 sig2VBscratch <- MCMCpack::rinvgamma(S,anVB/2,d/2)
30 BetasVBscratch <- MASS::mvrnorm(S, mu = bnVB, Sigma = B)
31 summary(coda::mcmc(sig2VBscratch)); summary(coda::mcmc(
  BetasVBscratch))

```

---

---

### 14.3 Summary

In this chapter, we present approximate methods designed for situations where the likelihood function does not have a closed-form expression (e.g., ABC and BSL), or where the sample size and parameter space are large (e.g., INLA and VI), making traditional MCMC and importance sampling (IS) methods ineffective. We provide the theoretical foundations and include applications to illustrate the potential of these methods.

*Simulation-based algorithms* suffer from the *curse of dimensionality* in the parameter space, while *optimization-based approaches* typically require explicit evaluation of the likelihood function. To address situations where both challenges arise simultaneously, recent developments, referred to as *hybrid methods*, combine these strategies with the *pseudo-marginal* MCMC approach [15, 296]. For instance, when pseudo-marginal MCMC is combined with variational Bayes, it becomes possible to perform inference in high-dimensional spaces even when the likelihood function is intractable [436].

---

### 14.4 Exercises

#### 1. g-and-k distribution for financial returns continues

Simulate a dataset following the specification given in the g-and-k distribution for financial returns example. Set  $\theta_1 = 0.8$ ,  $a = 1$ ,  $b = 0.5$ ,  $g = -1$ , and  $k = 1$ , with a sample size of 500, and use the same priors as in that example. Implement the ABC accept/reject algorithm from scratch using one million prior draws, selecting the 1,000 draws with the smallest distance.<sup>7</sup>

- Perform a linear regression adjustment using the posterior draws of our ABC-AR algorithm (ABC-AR-Adj).
- Compare the results with those obtained using the ABC-AR implementation in the *EasyABC* package, ensuring that the computational time is relatively similar between both implementations.
- Compare the posterior results of ABC-AR, ABC-AR-Adj and EasyABC with the population values.

#### 2. Simulation: g-and-k distribution continues

---

<sup>7</sup>Note that this setting does not satisfy the asymptotic requirements for Bayesian consistency. However, it serves as a pedagogical exercise.

Perform the simulation example of the Bayesian synthetic likelihood presented in the book, using the same population parameters and setting  $M = 500$ ,  $S = 6,000$ , with burn-in and thinning parameters set to 1,000 and 5, respectively. Use the *BSL* package in **R** to perform inference using the vanilla, unbiased, semi-parametric, and misspecified (mean and variance) versions of BSL. Compare the posterior distributions of the methods with the true population parameters.

3. Simulate a multinomial logit model (see Section 6.5 in the book) with 3 alternatives, 2 alternative-specific regressors, and 1 individual-specific regressor. The population parameters for the alternative-specific regressors are  $-0.3$  and  $1.2$ , while the population values for the individual-specific regressor are  $0.3$ ,  $0$ , and  $0.5$ . All regressors are assumed to follow a standard normal distribution, and the sample size is 1,000.

Perform inference using the *INLA* package, and note that the Poisson trick should be used for multinomial models in this exercise (see [405] for details).

4. Get the expression for the ELBO in the linear regression model with conjugate family.

#### 5. Linear regression continues

Perform inference in the linear regression example using stochastic variational inference via the automatic differentiation variational inference (ADVI) approach [253], implemented in *rstan* package. In particular, consider the model  $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \mu_i$ , assuming non-informative independent priors:  $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$  and  $\sigma^2 \sim \text{IG}(\alpha_0/2, \delta_0/2)$ , where  $\boldsymbol{\beta}_0 = \mathbf{0}_3$ ,  $\mathbf{B}_0 = 1000\mathbf{I}_3$ , and  $\alpha_0 = \delta_0 = 0.01$ . The sample size is one million.

6. Let's retake the mixture regression model of Chapter 11, that is, the simple regression mixture with two components such that  $z_{i1} \sim \text{Ber}(0.5)$ , consequently,  $z_{i2} = 1 - z_{i1}$ , and assume one regressor,  $x_i \sim N(0, 1)$ ,  $i = 1, 2, \dots, 1,000$ . Then,

$$p(y_i | \mathbf{x}_i) = 0.5\phi(y_i | 2 + 1.5x_i, 1^2) + 0.5\phi(y_i | -1 + 0.5x_i, 0.8^2).$$

Let's set  $\alpha_{h0} = \delta_{h0} = 0.01$ ,  $\boldsymbol{\beta}_{h0} = \mathbf{0}_2$ ,  $\mathbf{B}_{h0} = \mathbf{I}_2$ , and  $\mathbf{a}_0 = [1/2 \ 1/2]^\top$ .

Perform VB inference in this model using the CAVI algorithm.



---

## Bibliography

---

- [1] Arild Aakvik, James J. Heckman, and Edward J. Vytlacil. Estimating treatment effects for discrete outcomes when responses to treatment vary: An application to norwegian vocational rehabilitation programs. *Journal of Econometrics*, 125(1-2):15–51, 2005.
- [2] Alberto Abadie. Semiparametric difference-in-differences estimators. *Review of Economic Studies*, 72(1):1–19, 2005.
- [3] Alberto Abadie, Joshua D. Angrist, and Guido W. Imbens. Instrumental variables estimates of the effect of subsidized training on the quantiles of trainee earnings. *Econometrica*, 70(1):91–117, 2002.
- [4] D. Acemoglu, S. Johnson, and J. Robinson. The colonial origins of comparative development: An empirical investigation. *The American Economic Review*, 91(5):1369–1401, 2001.
- [5] J. H. Albert and S. Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669–679, 1993.
- [6] Jim Albert. *Bayesian Computation with R*. Use R! Springer, New York, NY, 2nd edition, 2009.
- [7] Izhar Asael Alonso and Cristian Cruz. varstan: An R package for Bayesian time series models with Stan. *ARXIV Preprint*, 2020.
- [8] Sungbae An and Frank Schorfheide. Bayesian analysis of dsge models. *Econometric reviews*, 26(2-4):113–172, 2007.
- [9] Ziwen An, David J Nott, and Christopher Drovandi. Robust bayesian synthetic likelihood via a semi-parametric approach. *Statistics and Computing*, 30(3):543–557, 2020.
- [10] Ziwen An, Leah F South, and Christopher Drovandi. Bsl: An r package for efficient parameter estimation for simulation-based models via bayesian synthetic likelihood. *Journal of Statistical Software*, 101:1–33, 2022.
- [11] Ziwen An, Leah F South, David J Nott, and Christopher C Drovandi. Accelerating bayesian synthetic likelihood with the graphical lasso. *Journal of Computational and Graphical Statistics*, 28(2):471–475, 2019.

- [12] Isaiah Andrews and Anna Mikusheva. Optimal decision rules for weak GMM. *Econometrica*, 90(2):715–748, March 2022.
- [13] C. Andrieu, A. Doucet, and R. Holenstein. Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3):269–342, 2010.
- [14] Christophe Andrieu, Arnaud Doucet, and Roman Holenstein. Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3):269–342, 2010.
- [15] Christophe Andrieu and Gareth O. Roberts. The pseudo-marginal approach for efficient monte carlo computations. *Annals of Statistics*, 37(2):697–725, 2009.
- [16] Joshua D. Angrist, Guido W. Imbens, and Donald B. Rubin. Identification of causal effects using instrumental variables. *Journal of the American Statistical Association*, 91(434):444–455, 1996.
- [17] Joshua D. Angrist and Jörn-Steffen Pischke. Better late than nothing: Some comments on deaton (2009) and heckman and urzua (2009). *Journal of Economic Literature*, 48(2):399–423, 2010.
- [18] Joshua D. Angrist and Jörn-Steffen Pischke. *Mastering 'Metrics: The Path from Cause to Effect*. Princeton University Press, Princeton, NJ, 2014.
- [19] Joshua D. Angrist and Jörn-Steffen Pischke. *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton University Press, Princeton, NJ, 2009.
- [20] Luc Anselin. A note of small sample properties of estimators in a first order spatial autoregressive model. *Environment and Planning*, 14:1023–1030, 1982.
- [21] C. Antoniak. Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *The Annals of Statistics*, 2(6):1152–1174, 1974.
- [22] R. Baath. *Package bayesboot*, 2018.
- [23] Andrew Baker, Brantly Callaway, Scott Cunningham, Andrew Goodman-Bacon, and Pedro H. C. Sant'Anna. Difference-in-differences designs: A practitioner's guide. *arXiv preprint arXiv:2503.13323*, June 2025. Revised June 18, 2025.
- [24] Jack Baker, Paul Fearnhead, Emily B Fox, and Christopher Nemeth. sgmc: An r package for stochastic gradient markov chain monte carlo. *Journal of Statistical Software*, 91:1–27, 2019.

- [25] M. Barbieri and J. Berger. Optimal predictive model selection. *The Annals of Statistics*, 32(3):870–897, 2004.
- [26] Rémi Bardenet, Arnaud Doucet, and Chris Holmes. Towards scaling up markov chain monte carlo: an adaptive subsampling approach. In *International conference on machine learning*, pages 405–413. PMLR, 2014.
- [27] Rémi Bardenet, Arnaud Doucet, and Chris Holmes. On markov chain monte carlo methods for tall data. *Journal of Machine Learning Research*, 18(47):1–43, 2017.
- [28] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- [29] Sanjib Basu and Siddhartha Chib. Marginal likelihood and bayes factors for dirichlet process mixture models. *Journal of the American Statistical Association*, 98(461):224–235, 2003.
- [30] M. Bayarri and J. Berger. P-values for composite null models. *Journal of American Statistical Association*, 95:1127–1142, 2000.
- [31] M. J. Bayarri and J. Berger. The interplay of bayesian and frequentist analysis. *Statistical Science*, 19(1):58–80, 2004.
- [32] T. Bayes. An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, 53:370–416, 1763.
- [33] Thomas Bayes. LII. an essay towards solving a problem in the doctrine of chances. by the late rev. mr. bayes, FRS communicated by Mr. Price, in a letter to John Canton, AMFR S. *Philosophical transactions of the Royal Society of London*, 53:370–418, 1763.
- [34] Mark A Beaumont, Wenyang Zhang, and David J Balding. Approximate Bayesian computation in population genetics. *Genetics*, 162(4):2025–2035, 2002.
- [35] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [36] Daniel J Benjamin, James O Berger, Magnus Johannesson, Brian A Nosek, E-J Wagenmakers, Richard Berk, Kenneth A Bollen, Björn Brembs, Lawrence Brown, Colin Camerer, et al. Redefine statistical significance. *Nature human behaviour*, 2(1):6–10, 2018.

- [37] J. Berger. *Statistical Decision Theory and Bayesian Analysis*. Springer, third edition edition, 1993.
- [38] J. Berger. The case for objective bayesian analysis. *Bayesian Analysis*, 1(3):385–402, 2006.
- [39] James O Berger. *Statistical decision theory and Bayesian analysis*. Springer Science & Business Media, 2013.
- [40] James O Berger and Luis R Pericchi. The intrinsic bayes factor for model selection and prediction. *Journal of the American Statistical Association*, 91(433):109–122, 1996.
- [41] J. Bernardo and A. Smith. *Bayesian Theory*. Wiley, Chichester, 1994.
- [42] José M Bernardo and Adrian FM Smith. *Bayesian theory*, volume 405. John Wiley & Sons, 2009.
- [43] Giorgio Bertorelle, Andrea Benazzo, and S Mona. ABC as a flexible framework to estimate demography over space and time: Some cons, many pros. *Molecular Ecology*, 19(13):2609–2625, 2010.
- [44] Peter J Bickel and Joseph A Yahav. Some contributions to the asymptotic theory of bayes solutions. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 11(4):257–276, 1969.
- [45] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Springer, 2006.
- [46] Pier Giovanni Bissiri, Chris C Holmes, and Stephen G Walker. A general framework for updating belief distributions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(5):1103–1130, 2016.
- [47] Roger Bivand, Virgilio Gómez-Rubio, and Håvard Rue. Spatial data analysis with r-inla with some extensions. *Journal of statistical software*, 63:1–31, 2015.
- [48] D. Blackwell and J. MacQueen. Ferguson distributions via Pólya urn schemes. *The Annals of Statistics*, 1:353–355, 1973.
- [49] David M Blei and Michael I Jordan. Variational inference for dirichlet process mixtures. *Bayesian Analysis*, 1(1):121–143, 2006.
- [50] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- [51] Michael GB Blum. Approximate Bayesian computation: a nonparametric perspective. *Journal of the American Statistical Association*, 105(491):1178–1187, 2010.

- [52] Tim Bollerslev. Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 31(3):307–327, 1986.
- [53] Howard D. Bondell and Brian J. Reich. Simultaneous regression shrinkage, variable selection, and supervised clustering of predictors with oscar. *Biometrics*, 64(1):115–123, 2008.
- [54] G. E. P. Box. Science and statistics. *Journal of the American Statistical Association*, 71:791–799, 1976.
- [55] George E. P. Box. Science and statistics. *Journal of the American Statistical Association*, 71(356):791–799, 1976.
- [56] George E. P. Box and Gwilym M. Jenkins. *Time Series Analysis: Forecasting and Control*. Holden-Day, San Francisco, 1st edition, 1976.
- [57] George EP Box. Robustness in the strategy of scientific model building. In *Robustness in statistics*, pages 201–236. Elsevier, 1979.
- [58] Zach Branson, Maxime Rischard, Luke Bornn, and Luke W. Miratrix. A nonparametric Bayesian methodology for regression discontinuity designs. *Journal of Statistical Planning and Inference*, 202:14–30, 2019. Preprint: arXiv:1704.04858.
- [59] Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
- [60] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [61] Leo Breiman, Jerome H Friedman, Richard A Olshen, and Charles J Stone. *Classification and Regression Trees*. Wadsworth International Group, Belmont, CA, 1984.
- [62] Christoph Breunig, Ruixuan Liu, and Zhengfei Yu. Double robust bayesian inference on average treatment effects. *Econometrica*, 93(2):539–568, 2025.
- [63] Eric Brochu, Vlad M. Cora, and Nando de Freitas. A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *arXiv preprint arXiv:1012.2599*, 2010.
- [64] Kenneth P Burnham and David R Anderson. Multimodel inference: understanding aic and bic in model selection. *Sociological methods & research*, 33(2):261–304, 2004.
- [65] Colin Cameron and Pravin Trivedi. *Microeometrics: Methods and Applications*. Cambridge, 2005.

- [66] Olivier Cappé, Simon J Godsill, and Eric Moulines. An overview of existing methods and recent advances in sequential monte carlo. *Proceedings of the IEEE*, 95(5):899–924, 2007.
- [67] O. Cappé, S. J. Godsill, and E. Moulines. An overview of existing methods and recent advances in sequential monte carlo. *Proceedings of the IEEE*, 95(5):899–924, 2007.
- [68] Bradley P Carlin, Alan E Gelfand, and Adrian FM Smith. Hierarchical bayesian analysis of changepoint problems. *Journal of the royal statistical society: series C (applied statistics)*, 41(2):389–405, 1992.
- [69] Bob Carpenter, Andrew Gelman, Matthew D Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan: A probabilistic programming language. *Journal of statistical software*, 76:1–32, 2017.
- [70] J. Carpenter, P. Clifford, and P. Fearnhead. Improved particle filter for nonlinear problems. *IEE Proceedings-Radar, Sonar and Navigation*, 146(1):2–7, 1999.
- [71] Chris K Carter and Robert Kohn. On gibbs sampling for state space models. *Biometrika*, 81(3):541–553, 1994.
- [72] G. Casella and E. Moreno. Objective Bayesian variable selection. *Journal of the American Statistical Association*, 101(473):157–167, 2006.
- [73] George Casella and Roger Berger. *Statistical inference*. CRC Press, 2024.
- [74] Matias D. Cattaneo, Nicolas Idrobo, and Rocío Titiunik. A practical introduction to regression discontinuity designs: Foundations. *Cambridge Elements: Quantitative and Computational Methods for Social Science*, 2019.
- [75] DM Ceperley and Mark Dewing. The penalty method for random walks with uncertain energies. *The Journal of chemical physics*, 110(20):9812–9820, 1999.
- [76] Joshua Chan, Gary Koop, Dale J Poirier, and Justin L Tobias. *Bayesian econometric methods*, volume 7. Cambridge University Press, 2019.
- [77] W. Chang. *Web Application Framework for R: Package shiny*. R Studio, 2018.
- [78] Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient hamiltonian monte carlo. In *International conference on machine learning*, pages 1683–1691. PMLR, 2014.

- [79] Xiaohong Chen, Timothy M. Christensen, and Elie Tamer. Monte carlo confidence sets for identified sets. *Econometrica*, 86(6):1965–2018, 2018.
- [80] V. Chernozhukov and H. Hong. An MCMC approach to classical estimation. *Journal of Econometrics*, 115:293–346, 2003.
- [81] Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21:C1–C68, 2018.
- [82] Victor Chernozhukov and Christian Hansen. The effects of 401(k) participation on the wealth distribution: An instrumental quantile regression analysis. *The Review of Economics and Statistics*, 86(3):735–751, 2004.
- [83] Victor Chernozhukov and Christian Hansen. An iv model of quantile treatment effects. *Econometrica*, 73(1):245–261, 2005.
- [84] Victor Chernozhukov, Christian Hansen, Nathan Kallus, Martin Spindler, and Vasilis Syrgkanis. Applied causal inference powered by ml and ai. arXiv preprint 2403.02467, arXiv, 2024.
- [85] Victor Chernozhukov and Han Hong. An mcmc approach to classical estimation. *Journal of econometrics*, 115(2):293–346, 2003.
- [86] S. Chib. Bayes inference in the Tobit censored regression model. *Journal of Econometrics*, 51:79–99, 1992.
- [87] S. Chib and B. Carlin. On MCMC sampling in hierarchical longitudinal models. *Statistics and Computing*, 9:17–26, 1999.
- [88] Siddhartha Chib. Bayes regression with autoregressive errors: A gibbs sampling approach. *Journal of econometrics*, 58(3):275–294, 1993.
- [89] Siddhartha Chib. Marginal likelihood from the gibbs output. *Journal of the american statistical association*, 90(432):1313–1321, 1995.
- [90] Siddhartha Chib. Analysis of treatment response data without the joint distribution of potential outcomes. *Journal of Econometrics*, 140(2):401–412, 2007.
- [91] Siddhartha Chib and Edward Greenberg. Bayes inference in regression models with arma (p, q) errors. *Journal of Econometrics*, 64(1-2):183–206, 1994.
- [92] Siddhartha Chib and Edward Greenberg. Understanding the metropolis-hastings algorithm. *The american statistician*, 49(4):327–335, 1995.

- [93] Siddhartha Chib, Edward Greenberg, and Ivan Jeliazkov. Estimation of semiparametric models in the presence of endogeneity and sample selection. *Journal of Computational and Graphical Statistics*, 18(2):321–348, 2009.
- [94] Siddhartha Chib, Edward Greenberg, and Anna Simoni. Nonparametric bayes analysis of the sharp and fuzzy regression discontinuity designs. *Econometric Theory*, 39(3):481–533, 2023.
- [95] Siddhartha Chib and Liana Jacobi. Bayesian fuzzy regression discontinuity analysis and returns to compulsory schooling. *Journal of Applied Econometrics*, 31(6):1026–1047, 2016.
- [96] Siddhartha Chib and Ivan Jeliazkov. Marginal likelihood from the metropolis–hastings output. *Journal of the American Statistical Association*, 96(453):270–281, 2001.
- [97] Siddhartha Chib and Todd A Kuffner. Bayes factor consistency. *arXiv preprint arXiv:1607.00292*, 2016.
- [98] Siddhartha Chib, Minchul Shin, and Anna Simoni. Bayesian estimation and comparison of moment condition models. *Journal of the American Statistical Association*, 113(524):1656–1668, 2018.
- [99] Hugh A Chipman, Edward I George, and Robert E McCulloch. Bayesian cart model search. *Journal of the American Statistical Association*, 93(443):935–948, 1998.
- [100] Hugh A Chipman, Edward I George, and Robert E McCulloch. Bart: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4(1):266–298, 2010.
- [101] Taeryon Choi and Mark J. Schervish. On posterior consistency in nonparametric regression problems. *Journal of Multivariate Analysis*, 98(10):1969–1987, 2007.
- [102] Gerda Claeskens and Nils Lid Hjort. Model selection and model averaging. *Cambridge books*, 2008.
- [103] M. Clyde and E. George. Model uncertainty. *Statistical Science*, 19(1):81–94, 2004.
- [104] T. Conley, C. Hansen, and P. Rossi. Plausibly exogenous. *The Review of Economics and Statistics*, 94(1):260–272, 2012.
- [105] Timothy G Conley, Christian B Hansen, Robert E McCulloch, and Peter E Rossi. A semi-parametric bayesian approach to the instrumental variable problem. *Journal of Econometrics*, 144(1):276–305, 2008.

- [106] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, Hoboken, NJ, 2nd edition, 2006.
- [107] John G. Cragg and Stephen G. Donald. Testing identifiability and specification in instrumental variable models. *Econometric Theory*, 9(2):222–240, 1993.
- [108] Richard K. Crump, V. Joseph Hotz, Guido W. Imbens, and Oscar A. Mitnik. Dealing with limited overlap in estimation of Average Treatment Effects. *Biometrika*, 96(1):187–199, March 2009.
- [109] Scott Cunningham. *Causal Inference: The Mixtape*. Yale University Press, New Haven, CT, 2021.
- [110] Johan Dahlin and Thomas B Schön. Getting started with particle metropolis-hastings for inference in nonlinear dynamical models. *Journal of Statistical Software*, 88:1–41, 2019.
- [111] A. P. Dawid, M. Musio, and S. E. Fienberg. From statistical evidence to evidence of causality. *Bayesian Analysis*, 11(3):725–752, 2016.
- [112] Bruno De Finetti. Foresight: its logical laws, its subjective sources. In H. E. Kyburg and H. E. Smokler, editors, *Studies in Subjective Probability*. Krieger, New York, 1937. p.55–118.
- [113] Piet De Jong and Neil Shephard. The simulation smoother for time series models. *Biometrika*, 82(2):339–350, 1995.
- [114] Angus Deaton. Instruments, randomization, and learning about development. *Journal of Economic Literature*, 48(2):424–455, 2010.
- [115] M. H. DeGroot. *Probability and statistics*. Addison-Wesley Publishing Co., London, 1975.
- [116] Pierre Del Moral, Arnaud Doucet, and Ajay Jasra. An adaptive sequential Monte Carlo method for approximate Bayesian computation. *Statistics and Computing*, 22:1009–1020, 2012.
- [117] Marco Del Negro and F. Schorfheide. Forecasting with bayesian var models. In John Geweke, Gary Koop, and Herman van Dijk, editors, *The Oxford Handbook of Bayesian Econometrics*, pages 224–254. Oxford University Press, 2011.
- [118] Persi Diaconis and Donald Ylvisaker. Conjugate priors for exponential families. *The Annals of statistics*, 7(2):269–281, 1979.
- [119] J. M. Dickey and E. Gunel. Bayes factors from mixed probabilities. *Journal of the Royal Statistical Society: Series B (Methodology)*, 40:43–46, 1978.

- [120] James M Dickey. The weighted likelihood ratio, linear hypotheses on normal location parameters. *The Annals of Mathematical Statistics*, pages 204–223, 1971.
- [121] Peter. Diggle, P. Heagerty, Liang K-Y., and S. Zeger. *Analysis of longitudinal data*. Oxford University Press, 2002.
- [122] Peng Ding. Bayesian robust inference of sample selection using selection-t models. *Journal of Multivariate Analysis*, 124:451–464, 2014.
- [123] Thomas Doan, Robert Litterman, and Christopher Sims. Forecasting and conditional projection using realistic prior distributions. *Econometric reviews*, 3(1):1–100, 1984.
- [124] Arnaud Doucet, Nando De Freitas, and Neil Gordon. An introduction to sequential monte carlo methods. *Sequential Monte Carlo methods in practice*, pages 3–14, 2001.
- [125] Arnaud Doucet, Adam M Johansen, et al. A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of nonlinear filtering*, 12(656-704):3, 2009.
- [126] Christopher Drovandi and David T Frazier. A comparison of likelihood-free methods with and without summary statistics. *Statistics and Computing*, 32(3):42, 2022.
- [127] Christopher C Drovandi and Anthony N Pettitt. Estimation of parameters for macroparasite population evolution using approximate Bayesian computation. *Biometrics*, 67(1):225–233, 2011.
- [128] Christopher C Drovandi and Anthony N Pettitt. Likelihood-free bayesian estimation of multivariate quantile distributions. *Computational Statistics & Data Analysis*, 55(9):2541–2556, 2011.
- [129] Christopher C Drovandi, Anthony N Pettitt, and Anthony Lee. Bayesian indirect inference using a parametric auxiliary model. *Statistical Science*, 30(1):72—95, 2015.
- [130] Simon Duane, Anthony D Kennedy, Brian J Pendleton, and Duncan Roweth. Hybrid monte carlo. *Physics letters B*, 195(2):216–222, 1987.
- [131] Y. D. Edwards and G. M. Allenby. Multivariate analysis of multiple response data. *Journal of Marketing Research*, 40:321–334, 2003.
- [132] B. Efron. Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, 7:1–26, 1979.
- [133] Bradley Efron and Trevor Hastie. *Computer age statistical inference*, volume 5. Cambridge University Press, 2016.

- [134] Bradley Efron and Trevor Hastie. *Computer age statistical inference, student edition: algorithms, evidence, and data science*, volume 6. Cambridge University Press, 2021.
- [135] J Paul Elhorst et al. *Spatial econometrics: from cross-sectional data to spatial panels*, volume 479. Springer, 2014.
- [136] Walter Enders. *Applied Econometric Time Series*. Wiley, Hoboken, NJ, 4th edition, 2014.
- [137] M. Escobar and M. West. Bayesian density estimation and inference using mixtures. *Journal of the American Statistical Association*, 90(430):577–588, 1995.
- [138] T. Ferguson. A Bayesian analysis of some nonparametric problems. *The Annals of Statistics*, 1(2):209–230, 1973.
- [139] Carmen Fernandez, Eduardo Ley, and Mark FJ Steel. Benchmark priors for bayesian model averaging. *Journal of Econometrics*, 100(2):381–427, 2001.
- [140] R. Fisher. *Statistical Methods for Research Workers*. Hafner, New York, 13th edition, 1958.
- [141] David Frazier, David J Nott, Christopher Drovandi, and Robert Kohn. Bayesian inference using synthetic likelihood: Asymptotics and adjustments. *Journal of the American Statistical Association*, 118(544):2821–2832, 2023.
- [142] David T Frazier and Christopher Drovandi. Robust approximate bayesian inference with synthetic likelihood. *Journal of Computational and Graphical Statistics*, 30(4):958–976, 2021.
- [143] David T Frazier, Worapree Maneesoonthorn, Gael M Martin, and Brendan PM McCabe. Approximate Bayesian forecasting. *International Journal of Forecasting*, 35(2):521–539, 2019.
- [144] David T Frazier, Gael M Martin, Christian P Robert, and Judith Rousseau. Asymptotic properties of approximate Bayesian computation. *Biometrika*, 105(3):593–607, 2018.
- [145] David T Frazier, Christian P Robert, and Judith Rousseau. Model misspecification in approximate Bayesian computation: consequences and diagnostics. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(2):421–444, 2020.
- [146] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.

- [147] Jerome H. Friedman. Multivariate adaptive regression splines. *Annals of Statistics*, 19(1):1–67, 1991.
- [148] Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5):1189–1232, 2001.
- [149] Sylvia Frühwirth-Schnatter. Data augmentation and dynamic linear models. *Journal of time series analysis*, 15(2):183–202, 1994.
- [150] Sylvia Frühwirth-Schnatter. *Finite mixture and Markov switching models*. Springer, 2006.
- [151] Wenjiang J. Fu. Penalized regression: The bridge versus the lasso. *Journal of Computational and Graphical Statistics*, 7(3):397–416, 1998.
- [152] George M. Furnival and R. W. Wilson. Regressions by leaps and bounds. *Technometrics*, 16(4):499–511, 1974.
- [153] George M. Furnival and Robert W. Wilson. Regressions by leaps and bounds. *Technometrics*, 16(4):499–511, 1974.
- [154] Jacob R. Gardner, Geoff Pleiss, David Bindel Wu, Kilian Q. Weinberger, and Andrew Gordon Wilson. Product kernel interpolation for scalable gaussian processes. In *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 84, pages 1407–1416. PMLR, 2018.
- [155] P. Garthwaite, J. Kadane, and A. O'Hagan. Statistical methods for eliciting probability distributions. *Journal of American Statistical Association*, 100(470):680–701, 2005.
- [156] A. E. Gelfand and A. F. M. Smith. Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association*, 85:398–409, 1990.
- [157] Alan E Gelfand and Dipak K Dey. Bayesian model choice: asymptotics and exact calculations. *Journal of the Royal Statistical Society: Series B (Methodological)*, 56(3):501–514, 1994.
- [158] A. Gelman and X. Meng. Model checking and model improvement. In Gilks, Richardson, and Speigelhalter, editors, *In Markov chain Monte Carlo in practice*. Springer US, 1996. Chapter 6, pp. 157–196.
- [159] A. Gelman, X. Meng, and H. Stern. Posterior predictive assessment of model fitness via realized discrepancies. *Statistica sinica*, pages 733–760, 1996.
- [160] Andrew Gelman, John B Carlin, Hal S Stern, David Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. Chapman and Hall/CRC, 2021.

- [161] Andrew Gelman et al. Prior distributions for variance parameters in hierarchical models (comment on article by browne and draper). *Bayesian analysis*, 1(3):515–534, 2006.
- [162] Andrew Gelman, Walter R Gilks, and Gareth O Roberts. Weak convergence and optimal scaling of random walk metropolis algorithms. *The annals of applied probability*, 7(1):110–120, 1997.
- [163] Andrew Gelman and Guido Imbens. Why ask why? forward causal inference and reverse causal questions. Technical report, National Bureau of Economic Research, 2013.
- [164] Andrew Gelman and Donald B. Rubin. Inference from iterative simulation using multiple sequences. *Statistical Science*, 7(4):457–472, 1992.
- [165] S Geman and D. Geman. Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741, 1984.
- [166] E. George and R. McCulloch. Variable selection via Gibbs sampling. *Journal of the American Statistical Association*, 88(423):881–889, 1993.
- [167] E. George and R. McCulloch. Approaches for Bayesian variable selection. *Statistica Sinica*, 7:339–373, 1997.
- [168] Edward I George and Robert E McCulloch. Variable selection via gibbs sampling. *Journal of the American Statistical Association*, 88(423):881–889, 1993.
- [169] Edward I George and Robert E McCulloch. Approaches for Bayesian variable selection. *Statistica Sinica*, 7:339–373, 1997.
- [170] J. Geweke. *Bayesian Statistics*, chapter Evaluating the accuracy of sampling-based approaches to calculating posterior moments. Clarendon Press, Oxford, UK., 1992.
- [171] John Geweke. Using simulation methods for bayesian econometric models: inference, development, and communication. *Econometric reviews*, 18(1):1–73, 1999.
- [172] John Geweke. Getting it right: Joint distribution tests of posterior simulators. *Journal of the American Statistical Association*, 99(467):799–804, 2004.
- [173] John Geweke. *Contemporary Bayesian econometrics and statistics*, volume 537. John Wiley & Sons, 2005.
- [174] John Geweke, Gary Koop, and Herman K van Dijk. *The Oxford handbook of Bayesian econometrics*. Oxford University Press, USA, 2011.

- [175] Charles J Geyer. Practical markov chain monte carlo. *Statistical science*, pages 473–483, 1992.
- [176] Raffaella Giacomini and Tetsuya Kitagawa. Robust bayesian inference for set-identified models. *Econometrica*, 89(4):1519–1556, 2021.
- [177] Richard D Gill and James M Robins. Causal inference for complex longitudinal data: the continuous case. *Annals of Statistics*, pages 1785–1811, 2001.
- [178] Ryan Giordano, Runjing Liu, Michael I Jordan, and Tamara Broderick. Evaluating sensitivity to the stick-breaking prior in bayesian nonparametrics. *Bayesian Analysis*, 1(1):1–34, 2022.
- [179] I. J. Good. The bayes/non bayes compromise: A brief review. *Journal of the American Statistical Association*, 87(419):597–606, September 1992.
- [180] S. N. Goodman. Toward evidence-based medical statistics. 1: The P value fallacy. *Annals of internal medicine*, 130(12):995–1004, 1999.
- [181] N. J. Gordon, D. J. Salmond, and A. F. Smith. Novel approach to nonlinear/non-gaussian bayesian state estimation. *IEE Proceedings F (Radar and Signal Processing)*, 140(2):107–113, 1993.
- [182] P. J. Green. Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika*, 82:711–732, 1995.
- [183] David S Greenberg, Marcel Nonnenmacher, and Jakob H Macke. Automatic posterior transformation for likelihood-free inference. In *International Conference on Machine Learning*, pages 2404–2414, 2019.
- [184] Edward Greenberg. *Introduction to Bayesian econometrics*. Cambridge University Press, 2012.
- [185] Damodar N. Gujarati and Dawn C. Porter. *Basic Econometrics*. McGraw-Hill Education, New York, NY, 5th edition, 2009.
- [186] Paul Gustafson. Local robustness in bayesian analysis. In *Robust Bayesian Analysis*, pages 71–88. Springer, 2000.
- [187] Trygve Haavelmo. The statistical implications of a system of simultaneous equations. *Econometrica*, 11(1):1–12, 1943.
- [188] Jinyong Hahn, Petra Todd, and Wilbert van der Klaauw. Identification and estimation of treatment effects with a regression-discontinuity design. *Econometrica*, 69(1):201–209, jan 2001.
- [189] P Richard Hahn, Jared S Murray, and Carlos M Carvalho. Bayesian regression tree models for causal inference: Regularization, confounding, and heterogeneous effects (with discussion). *Bayesian Analysis*, 15(3):965–1056, 2020.

- [190] Robert Haining. *Spatial data analysis in the social and environmental sciences*. Cambridge University Press, United of Kingdom, first edition, 1990.
- [191] Johannes Edmund Handschin and David Q Mayne. Monte carlo techniques to estimate the conditional expectation in multi-stage non-linear filtering. *International journal of control*, 9(5):547–559, 1969.
- [192] Lars Peter Hansen. Large sample properties of generalized method of moments estimators. *Econometrica*, 50(4):1029–1054, 1982.
- [193] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *Annals of statistics*, 50(2):949, 2022.
- [194] Trevor Hastie and Robert Tibshirani. Bayesian backfitting (with discussion). *Journal of the American Statistical Association*, 95(452):1228–1240, 2000.
- [195] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, New York, 2 edition, 2009.
- [196] W. Hastings. Monte Carlo sampling methods using Markov chains and their application. *Biometrika*, 57:97–109, 1970.
- [197] Fumio Hayashi. *Econometrics*. Princeton University Press, 2000.
- [198] James J. Heckman. Sample selection bias as a specification error. *Econometrica*, 47(1):153–161, 1979.
- [199] James J. Heckman. Varieties of selection bias. *American Economic Review*, 80(2):313–318, 1990.
- [200] James J. Heckman, Hedibert F. Lopes, and Rémi Piatek. Treatment effects: A bayesian perspective. *Econometric Reviews*, 33(1-4):36–67, 2014.
- [201] James J. Heckman and Edward Vytlacil. Structural equations, treatment effects, and econometric policy evaluation. *Econometrica*, 73(3):669–738, 2005.
- [202] P. Heidelberger and P. D. Welch. Simulation run length control in the presence of an initial transient. *Operations Research*, 31(6):1109–1144, 1983.
- [203] Lütkepohl Helmut. *New introduction to multiple time series analysis*. Springer, 2005.
- [204] Miguel A. Hernán and James M. Robins. *Causal Inference: What If*. Chapman & Hall/CRC, Boca Raton, FL, 2020.

- [205] Jennifer L Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.
- [206] Keisuke Hirano, Guido W. Imbens, Donald B. Rubin, and Xiaohong Zhou. Assessing the effect of an influenza vaccine in an encouragement design. *Biostatistics*, 1(1):69–88, 2000.
- [207] Arthur E. Hoerl and Robert W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- [208] J. Hoeting, D. Madigan, A. Raftery, and C. Volinsky. Bayesian model averaging: A tutorial. *Statistical Science*, 14(4):382–417, 1999.
- [209] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *the Journal of machine Learning research*, 14(1):1303–1347, 2013.
- [210] Paul W. Holland. Statistics and causal inference. *Journal of the American Statistical Association*, 81(396):945–960, 1986.
- [211] Daniel Hosszejni and Gregor Kastner. Modeling univariate and multivariate stochastic volatility in r with `stochvol` and `factorstochvol`. *Journal of Statistical Software*, 100(12):1–34, 2021.
- [212] Zhen Huang and Andrew Gelman. Sampling for bayesian computation with large datasets. Technical report, Department of Statistics, Columbia University, 2005. Technical Report, Department of Statistics, Columbia University.
- [213] Joseph G. Ibrahim and Purushottam W. Laud. On bayesian analysis of generalized linear models using jeffreys’s prior. *Journal of the American Statistical Association*, 86(416):981–986, 1991.
- [214] Guido W. Imbens. Instrumental variables: An econometrician’s perspective. *Statistical Science*, 29(3):323–358, 2014. Also IZA Discussion Paper No. 8048 and NBER Working Paper No. 19983.
- [215] Guido W. Imbens and Joshua D. Angrist. Identification and estimation of local average treatment effects. *Econometrica*, 62(2):467–475, 1994.
- [216] Guido W. Imbens and Thomas Lemieux. Regression discontinuity designs: A guide to practice. *Journal of Econometrics*, 142(2):615–635, 2008.
- [217] Guido W Imbens and Donald B Rubin. Bayesian inference for causal effects in randomized experiments with noncompliance. *The annals of statistics*, pages 305–327, 1997.

- [218] Guido W. Imbens and Donald B. Rubin. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press, New York, NY, 2015.
- [219] H. Ishwaran and J. S. Rao. Spike and slab variable selection: Frequentist and Bayesian strategies. *The Annals of Statistics*, 33(2):730–773, 2005.
- [220] Liana Jacobi, Chun Fung Kwok, Andrés Ramírez-Hassan, and Nhung Nghiem. Posterior manifolds over prior parameter regions: Beyond pointwise sensitivity assessments for posterior statistics from mcmc inference. *Studies in Nonlinear Dynamics & Econometrics*, 28(2):403–434, 2024.
- [221] Liana Jacobi, Dan Zhu, and Mark Joshi. Estimating posterior sensitivities with application to structural analysis of bayesian vector autoregressions, 2022.
- [222] H. Jeffreys. Some test of significance, treated by the theory of probability. *Proceedings of the Cambridge philosophy society*, 31:203–222, 1935.
- [223] H. Jeffreys. *Theory of Probability*. Oxford University Press, London, 1961.
- [224] Harold Jeffreys. An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 186(1007):453–461, 1946.
- [225] M. Jetter and A. Ramírez Hassan. Want export diversification? Educate the kids first. *Economic Inquiry*, 53(4):1765–1782, 2015.
- [226] Michael Jetter, Rafat Mahmood, Christopher F. Parmeter, and Andrés Ramírez-Hassan. Post-cold war civil conflict and the role of history and religion: A stochastic search variable selection approach. *Economic Modelling*, 114:105907, 2022.
- [227] V. E. Jhonson and D. Rossell. Bayesian model selection in high-dimensional settings. *Journal of the American Statistical Association*, 107(498):649–660, 2012.
- [228] Bai Jiang. Approximate bayesian computation with kullback-leibler divergence as data discrepancy. In *International conference on artificial intelligence and statistics*, pages 1711–1721. PMLR, 2018.
- [229] Wenhua Jiang and Martin A. Tanner. Gibbs posterior for variable selection in high-dimensional classification and data mining. *Annals of Statistics*, 36(5):2207–2231, 2008.

- [230] Valen E Johnson and David Rossell. Bayesian model selection in high-dimensional settings. *Journal of the American Statistical Association*, 107(498):649–660, 2012.
- [231] M. I. Jordan, Z. Ghahramani, T. Jaakkola, and L. Saul. Introduction to variational methods for graphical models. *Machine Learning*, 37:183–233, 1999.
- [232] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37:183–233, 1999.
- [233] J. B. Kadane. Predictive and structural methods for eliciting prior distributions. In A Zellner, editor, *Bayesian Analysis in Econometrics and Statistics: Essays in honor of Harold Jeffreys*, pages 89–93. North-Holland Publishing Company,, Amsterdam, 1980.
- [234] Joseph Kadane and Lara Wolfson. Experiences in elicitation. *The Statistician*, 47(1):3–19, 1998.
- [235] Daniel Kahneman. *Thinking, fast and slow*. Macmillan, 2011.
- [236] N. Kantas, A. Doucet, S. S. Singh, and J. M. Maciejowski. An overview of sequential monte carlo methods for parameter estimation in general state-space models. *IFAC Proceedings Volumes*, 42(10):774–785, 2009.
- [237] Nikolas Kantas, Arnaud Doucet, Sumeetpal S Singh, Jan Maciejowski, Nicolas Chopin, et al. On particle methods for parameter estimation in state-space models. *Statistical science*, 30(3):328–351, 2015.
- [238] G. Karabatsos. A menu-driven software package of Bayesian nonparametric (and parametric) mixed models for regression analysis and density estimation. *Behavior Research Methods*, 49:335–362, 2016.
- [239] A. Karl and A. Lenkoski. Instrumental variable Bayesian model averaging via conditional Bayes factor. Technical report, Heidelberg University, 2012.
- [240] R. Kass. Statistical inference: the big picture. *Statistical science*, 26(1):1–9, 2011.
- [241] Robert E. Kass and Adrian E. Raftery. Bayes factorss. *Journal of American Statistical Association*, 90(430):773–795, 1995.
- [242] Gregor Kastner and Sylvia Frühwirth-Schnatter. Ancillarity-sufficiency interweaving strategy (asis) for boosting mcmc estimation of stochastic volatility models. *Computational Statistics & Data Analysis*, 76:408–423, 2014.

- [243] Chang-Jin Kim and Charles R Nelson. Has the us economy become more stable? a bayesian approach based on a markov-switching model of the business cycle. *Review of Economics and Statistics*, 81(4):608–616, 1999.
- [244] Frank Kleibergen and Richard Paap. Generalized reduced rank tests using the singular value decomposition. *Journal of Econometrics*, 133(1):97–126, 2006.
- [245] Bas J. K. Kleijn and Aad W. van der Vaart. Misspecification in infinite-dimensional bayesian statistics. *Annals of Statistics*, 34(2):837–877, 2006.
- [246] Augustine Kong, Jun S Liu, and Wing Hung Wong. Sequential imputations and bayesian missing data problems. *Journal of the American statistical association*, 89(425):278–288, 1994.
- [247] G Koop, R León-Gonzalez, and R Strachan. Bayesian model averaging in the instrumental variable regression model. *Journal of Econometrics*, 171:237–250, 2012.
- [248] Gary Koop, Dimitris Korobilis, et al. Bayesian multivariate time series methods for empirical macroeconomics. *Foundations and Trends® in Econometrics*, 3(4):267–358, 2010.
- [249] Gary Koop and Dale J Poirier. Learning about the across-regime correlation in switching regression models. *Journal of Econometrics*, 78(2):217–227, 1997.
- [250] Gary M Koop. *Bayesian econometrics*. John Wiley & Sons Inc., 2003.
- [251] Hideo Kozumi and Genya Kobayashi. Gibbs sampling methods for Bayesian quantile regression. *Journal of Statistical Computation and Simulation*, 81(11):1565—1578, 2011.
- [252] Fabian Krueger. *bvars: Bayesian Analysis of a Vector Autoregressive Model with Stochastic Volatility and Time-Varying Parameters*, 2022. R package version 1.1.
- [253] Alp Kucukelbir, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M Blei. Automatic differentiation variational inference. *Journal of machine learning research*, 18(14):1–45, 2017.
- [254] Minjung Kyung, Jeff Gill, Malay Ghosh, and George Casella. Penalized regression, standard errors, and bayesian lassos. *Bayesian Analysis*, 5(2):369–411, 2010.
- [255] Tony Lancaster. *An introduction to modern Bayesian econometrics*. Blackwell Oxford, 2004.
- [256] P. Laplace. *Théorie Analytique des Probabilités*. Courcier, 1812.

- [257] Pierre Simon Laplace. Mémoire sur la probabilité de causes par les événements. *Mémoire de l'académie royale des sciences*, 1774.
- [258] Nicole A. Lazar. Bayesian empirical likelihood. *Biometrika*, 90(2):319–326, 2003.
- [259] Nicole A Lazar. A review of empirical likelihood. *Annual Review of Statistics and its Application*, 8(1):329–344, 2021.
- [260] E.L. Lehmann and George Casella. *Theory of Point Estimation*. Springer, second edition edition, 2003.
- [261] Alex Lenkoski, Theo S. Eicher, and Adrian Raftery. Two-stage Bayesian model averaging in endogeneous variable models. *Econometric Reviews*, 33, 2014.
- [262] Alex Lenkoski, Anna Karl, and Andreas Neudecker. *Package ivbma*, 2013.
- [263] Maxime Lenormand, Franck Jabot, and Guillaume Deffuant. Adaptive approximate Bayesian computation for complex models. *Computational Statistics*, 28(6):2777–2796, 2013.
- [264] James LeSage and Robert Kelley Pace. *Introduction to spatial econometrics*. Chapman and Hall/CRC, 2009.
- [265] Christoph Leuenberger and Daniel Wegmann. Bayesian computation and model selection without likelihoods. *Genetics*, 184(1):243–252, 2010.
- [266] Arthur Lewbel and Krishna Pendakur. Tricks with hicks: The easi demand system. *American Economic Review*, 99(3):827–863, 2009.
- [267] Eduardo Ley and Mark FJ Steel. On the effect of prior assumptions in Bayesian model averaging with applications to growth regression. *Journal of Applied Econometrics*, 24(4):651–674, 2009.
- [268] Wei Li, Sangwoo Ahn, and Max Welling. Scalable mcmc for mixed membership stochastic blockmodels. In JMLR Workshop and Conference Proceedings, editors, *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 51, pages 723–731, 2016.
- [269] D. V. Lindley. The philosophy of statistics. *The Statistician*, 49(3):293–337, 2000.
- [270] D. V. Lindley and L. D. Phillips. Inference for a Bernoulli process (a Bayesian view). *American Statistician*, 30:112–119, 1976.
- [271] Dennis V Lindley. A statistical paradox. *Biometrika*, 44(1/2):187–192, 1957.

- [272] Antonio R Linero. Bayesian regression trees for high-dimensional prediction and variable selection. *Journal of the American Statistical Association*, 113(522):626–636, 2018.
- [273] Jarno Lintusaari, Michael U Gutmann, Ritabrata Dutta, Samuel Kaski, and Jukka Corander. Fundamentals and recent developments in approximate bayesian computation. *Systematic biology*, 66(1):e66–e82, 2017.
- [274] Robert B Litterman. Forecasting with bayesian vector autoregressions—five years of experience. *Journal of Business & Economic Statistics*, 4(1):25–38, 1986.
- [275] Roderick J. A. Little and Donald B. Rubin. *Statistical Analysis with Missing Data*. Wiley Series in Probability and Statistics. John Wiley & Sons, 3rd edition, 2019.
- [276] J. S. Liu and R. Chen. Blind deconvolution via sequential imputations. *Journal of the American Statistical Association*, 90(430):567–576, 1995.
- [277] Pangpang Liu and Yichuan Zhao. A review of recent advances in empirical likelihood. *Wiley Interdisciplinary Reviews: Computational Statistics*, 15(3):e1599, 2023.
- [278] Ruixuan Liu and Zhengfei Yu. Quasi-bayesian estimation and inference with control functions. *arXiv preprint arXiv:2402.17374*, 2024. Revised May 2025.
- [279] Albert Y. Lo. A large sample study of the Bayesian bootstrap. *The Annals of Statistics*, 15(1):360–375, 1987.
- [280] Jan-Matthis Lueckmann, Pedro J Goncalves, Gabriele Bassetto, Kaan Öcal, Marcel Nonnenmacher, and Jakob H Macke. Flexible statistical inference for mechanistic models of neural dynamics. *Advances in Neural Information Processing Systems*, 30, 2017.
- [281] Gruber Luis and Kastner Gregor. *bayesianVARs: MCMC Estimation of Bayesian Vectorautoregressions*, 2024. R package version 1.0.5.
- [282] Yu Luo, Daniel J. Graham, and Emma J. McCoy. Semiparametric bayesian doubly robust causal estimation. *Journal of Statistical Planning and Inference*, 225:171–187, 2023.
- [283] Yi-An Ma, Tianqi Chen, and Emily Fox. A complete recipe for stochastic gradient mcmc. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2917–2925, 2015.
- [284] Dougal Maclaurin and Ryan P. Adams. Firefly monte carlo: Exact mcmc with subsets of data. In *The Twenty-Fourth International Joint Conference on Artificial Intelligence*, pages 4279–4295. IJCAI, 2015.

- [285] D. Madigan and A. E. Raftery. Model selection and accounting for model uncertainty in graphical models using Occam’s window. *Journal of the American Statistical Association*, 89(428):1535–1546, 1994.
- [286] D. Madigan, J. C. York, and D. Allard. Bayesian graphical models for discrete data. *International Statistical Review*, 63(2):215–232, 1995.
- [287] David Madigan, Jeremy York, and Denis Allard. Bayesian graphical models for discrete data. *International Statistical Review*, 63:215–232, 1995.
- [288] Charles F. Manski. Anatomy of the selection problem. *The Journal of Human Resources*, 24(3):343–360, 1989.
- [289] Charles F. Manski. Nonparametric bounds on treatment effects. *American Economic Review*, 80(2):319–323, 1990.
- [290] Charles F. Manski. *Identification Problems in the Social Sciences*. Harvard University Press, Cambridge, MA, 1995.
- [291] Charles F. Manski. *Partial Identification of Probability Distributions*. Springer, New York, 2003.
- [292] Paul Marjoram, John Molitor, Vincent Plagnol, and Simon Tavaré. Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences of the USA*, 100(26):15324–15328, 2003.
- [293] Jacob Marschak and William H. Andrews. Random simultaneous equations and the theory of production. *Econometrica*, 12(3-4):143–205, 1944.
- [294] Andrew D. Martin, Kevin M. Quinn, and Jong Hee Park. MCMCpack: Markov chain Monte Carlo in R. *Journal of Statistical Software*, 42(9):1–21, 2011.
- [295] Andrew D. Martin, Kevin M. Quinn, and Jong Hee Park. *Package MCMCpack*, 2018.
- [296] Gael M Martin, David T Frazier, and Christian P Robert. Approximating bayes in the 21st century. *Statistical Science*, 39(1):20–45, 2024.
- [297] Gael M Martin, David T Frazier, and Christian P Robert. Computing bayes: From then ‘til now. *Statistical Science*, 39(1):3–19, 2024.
- [298] Gael M Martin, Rubén Loaiza-Mayo, Worapree Maneesoonthorn, David T Frazier, and Andrés Ramírez-Hassan. Optimal probabilistic forecasts: When do they work? *International Journal of Forecasting*, 38(1):384–406, 2022.

- [299] Osvaldo A. Martin, Ravin Kumar, and Junpeng Lao. *Bayesian Modeling and Computation in Python*. publisher=Chapman and Hall/CRC, Boca Raton, December 2021.
- [300] Sara Martino and Andrea Riebler. Integrated nested laplace approximations (inla). *arXiv preprint arXiv:1907.01248*, 2019.
- [301] Tyler H. McCormick, Adrian Raftery, David Madigan, Sevvandi Kan-danaarachchi, and Hana Sevcikova. *dma: Dynamic Model Averaging*, 2018. R package version 1.4-0.
- [302] Tyler H McCormick, Adrian E Raftery, David Madigan, and Randall S Burd. Dynamic logistic regression and dynamic model averaging for binary classification. *Biometrics*, 68(1):23–30, 2012.
- [303] P. McCullagh and J. A. Nelder. *Generalized Linear Models*. London: Chapman and Hall, 1989.
- [304] R. McCulloch and P. Rossi. An exact likelihood analysis of the multinomial probit model. *Journal of Econometrics*, 64:207–240, 1994.
- [305] Robert E McCulloch, Nicholas G Polson, and Peter E Rossi. A bayesian analysis of the multinomial probit model with fully identified parameters. *Journal of econometrics*, 99(1):173–193, 2000.
- [306] Sharon Bertsch McGrayne. *The Theory That Would Not Die: How Bayes’ Rule Cracked the Enigma Code, Hunted Down Russian Submarines, & Emerged Triumphant from Two Centuries of C.* Yale University Press, 2011.
- [307] Karel Mertens and Morten O Ravn. A reconciliation of svar and narrative estimates of tax multipliers. *Journal of Monetary Economics*, 68:S1–S19, 2014.
- [308] N. Metropolis, A. Rosenbluth, M. Rosenbluth, A. Teller, and E. Teller. Equations of state calculations by fast computing machines. *J. Chem. Phys*, 21:1087–1092, 1953.
- [309] J. Miller and M. Harrison. Inconsistency of Pitman-Yor process mixtures for the number of components. *Journal of Machine Learning Research*, 15:3333–3370, 2014.
- [310] Stanislav Minsker. Scalable and robust bayesian inference via the median posterior. *Proceedings of the 31st International Conference on Machine Learning (ICML)*, 37:1656–1664, 2015.
- [311] Stanislav Minsker, Subhabrata Srivastava, Lin Lin, and David B. Dun-son. Robust and scalable bayes via a median of subset posterior mea-sures. *Journal of Machine Learning Research*, 18(1):4488–4527, 2017.

- [312] Franz X. Mohr. *bvartools: Bayesian Inference of Vector Autoregressive and Error Correction Models*, 2024. R package version 0.2.4.
- [313] Hyungsik Roger Moon and Frank Schorfheide. Bayesian and frequentist inference in partially identified models. *Econometrica*, 80(2):755–782, 2012.
- [314] P. Müller, F. Quintana, A. Jara, and T. Hanson. *Bayesian Nonparametric Data Analysis*. Springer, New York, 2015.
- [315] Ulrich K. Müller. Risk of Bayesian inference in misspecified models, and the sandwich covariance matrix. *Econometrica*, 81(5):1805–1849, September 2013.
- [316] Radford M. Neal. *Bayesian Learning for Neural Networks*, volume 118 of *Lecture Notes in Statistics*. Springer, 1996.
- [317] Radford M. Neal. Mcmc using hamiltonian dynamics. In Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng, editors, *Handbook of Markov Chain Monte Carlo*, chapter 5, pages 113–162. Chapman and Hall/CRC, 2011.
- [318] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [319] Willie Neiswanger, Chong Wang, and Eric P Xing. Asymptotically exact, embarrassingly parallel mcmc. In *Proceedings of the Thirtieth International Conference on Machine Learning*, pages 623–632. PMLR, 2013.
- [320] J.A. Nelder and R.W.M. Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society: Series A (General)*, 135(3):370–384, 1972.
- [321] Christopher Nemeth and Paul Fearnhead. Stochastic gradient markov chain monte carlo. *Journal of the American Statistical Association*, 116(533):433–450, 2021.
- [322] J. Neyman and E. Pearson. On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society, Series A*, 231:289–337, 1933.
- [323] Duy Nguyen. An in depth introduction to variational bayes note. Available at SSRN 4541076, 2023.
- [324] Kuschnig Nikolas, Vashold Lukas, Tomass Nirai, McCracken Michael, and Ng Serena. *BVAR: Hierarchical Bayesian Vector Autoregression*, 2022. R package version 1.0.5.

- [325] Agostino Nobile. Comment: Bayesian multinomial probit models with a normalization constraint. *Journal of Econometrics*, 99(2):335–345, 2000.
- [326] James P Normington, Eric F Lock, Thomas A Murray, and Caroline S Carlin. Bayesian variable selection in hierarchical difference-in-differences models. *Statistical methods in medical research*, 31(1):169–183, 2022.
- [327] Robert B. O’Hara and Mikko J. Sillanpää. A review of bayesian variable selection methods: What, how and which. *Bayesian Analysis*, 4(1):85–118, 2009.
- [328] K. Ord. Estimation methods for models of spatial interaction. *Journal of the American Statistical Association*, 70(349):120–126, 1975.
- [329] George Papamakarios and Iain Murray. Fast  $\epsilon$ -free inference of simulation models with bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- [330] T. Park and G. Casella. The Bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.
- [331] G. Parmigiani and L. Inoue. *Decision theory principles and approaches*. John Wiley & Sons, 2008.
- [332] Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–710, 1995.
- [333] Judea Pearl and Dana Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic Books, New York, 2018.
- [334] Luis Pericchi and Carlos Pereira. Adaptative significance levels using optimal decision rules: Balancing by weighting the error probabilities. *Brazilian Journal of Probability and Statistics*, 2015.
- [335] Giovanni Petris. An R package for dynamic linear models. *Journal of Statistical Software*, 36(12):1–16, 2010.
- [336] Giovanni Petris, Sonia Petrone, and Patrizia Campagnoli. Dynamic linear models. In *Dynamic Linear Models with R*, pages 31–84. Springer, 2009.
- [337] Michael K Pitt, Ralph dos Santos Silva, Paolo Giordani, and Robert Kohn. On some properties of markov chain monte carlo simulation methods based on the particle filter. *Journal of Econometrics*, 171(2):134–151, 2012.
- [338] Geoff Pleiss, Jacob R. Gardner, Kilian Q. Weinberger, and Andrew Gordon Wilson. Constant-time predictive distributions for gaussian processes. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, volume 80, pages 4111–4120. PMLR, 2018.

- [339] Martyn Plummer, Nicky Best, Kate Cowles, Karen Vines, Deepayan Sarkar, Douglas Bates, Russell Almond, and Arni Magnusson. *Output Analysis and Diagnostics for MCMC*, 2016.
- [340] Andy Pole, Mike West, and Jeff Harrison. *Applied Bayesian forecasting and time series analysis*. Chapman and Hall/CRC, 2018.
- [341] Leah F Price, Christopher C Drovandi, Anthony Lee, and David J Nott. Bayesian synthetic likelihood. *Journal of Computational and Graphical Statistics*, 27(1):1–11, 2018.
- [342] Jonathan K Pritchard, Mark T Seielstad, Anna Perez-Lezaun, and Marcus W Feldman. Population growth of human Y chromosomes: a study of Y chromosome microsatellites. *Molecular Biology and Evolution*, 16(12):1791–1798, 1999.
- [343] Matias Quiroz, Robert Kohn, Mattias Villani, and Minh-Ngoc Tran. Speeding up mcmc by efficient data subsampling. *Journal of the American Statistical Association*, 2019.
- [344] Matias Quiroz, Mattias Villani, Robert Kohn, Minh-Ngoc Tran, and Khue-Dung Dang. Subsampling mcmc—an introduction for the survey statistician. *Sankhya: The Indian Journal of Statistics, Series A*, 80:33–69, 2018.
- [345] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2021.
- [346] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2023.
- [347] A. Raftery. Bayesian model selection in social research. *Sociological Methodology*, 25:111–163, 1995.
- [348] Adrian Raftery, Jennifer Hoeting, Chris Volinsky, Ian Painter, and Ka Yee Yeung. *Package BMA*, 2012.
- [349] Adrian E Raftery. Bayesian model selection in structural equation models. *Sage Focus Editions*, 154:163–163, 1993.
- [350] Adrian E Raftery, Miroslav Kárný, and Pavel Ettler. Online prediction under model uncertainty via dynamic model averaging: Application to a cold rolling mill. *Technometrics*, 52(1):52–66, 2010.
- [351] Adrian E. Raftery, David Madigan, and Jennifer A. Hoeting. Bayesian model averaging for linear regression models. *Journal of the American Statistical Association*, 92(437):179–191, 1997.
- [352] A.E. Raftery and S.M. Lewis. One long run with diagnostics: Implementation strategies for Markov chain Monte Carlo. *Statistical Science*, 7:493–497, 1992.

- [353] A. Ramírez Hassan. The interplay between the Bayesian and frequentist approaches: A general nesting spatial panel data model. *Spatial Economic Analysis*, 12(1):92–112, 2017.
- [354] A. Ramírez Hassan, J. Cardona Jiménez, and R. Cadavid Montoya. The impact of subsidized health insurance on the poor in Colombia: Evaluating the case of Medellín. *Economia Aplicada*, 17(4):543–556, 2013.
- [355] A. Ramírez-Hassan and R. Guerra-Urzola. Bayesian treatment effects due to a subsidized health program: The case of preventive health care utilization in medellín (Colombia). *Empirical Economics*, 60:1477–1506, 2021.
- [356] Andrés Ramírez-Hassan. Dynamic variable selection in dynamic logistic regression: an application to internet subscription. *Empirical Economics*, 59(2):909–932, 2020.
- [357] Andrés Ramírez-Hassan and Daniela A. Carvajal-Rendón. Specification uncertainty in modeling internet adoption: A developing city case analysis. *Utilities Policy*, 70:101218, 2021.
- [358] Andrés Ramírez-Hassan and David T. Frazier. Testing model specification in approximate bayesian computation using asymptotic properties. *Journal of Computational and Graphical Statistics*, 33(3):1–14, 2024.
- [359] Andrés Ramírez-Hassan, Gustavo A. García, Estefanía Saravia, Juan Fernando Duque, and Daniel Londoño. What kind of schools parents choose when they have more options? effects of school transport subsidies. *Socio-Economic Planning Sciences*, 87:101509, 2023.
- [360] Andrés Ramírez-Hassan and Alejandro López-Vera. Welfare implications of a tax on electricity: A semi-parametric specification of the incomplete easi demand system. *Energy Economics*, 131:1–13, 2024.
- [361] Andrés Ramírez Hassan and Santiago Montoya Blandón. Welfare gains of the poor: An endogenous bayesian approach with spatial random effects. *Econometric Reviews*, 38(3):301–318, 2019.
- [362] F. Ramsey. Truth and probability. In Routledge and Kegan Paul, editors, *The Foundations of Mathematics and other Logical Essays*. New York: Harcourt, Brace and Company, London, 1926. Ch. VII, p.156–198.
- [363] A. Ramírez-Hassan and M. Graciano-Londoño. A GUIded tour of Bayesian regression. *The R Journal*, 13(2):135–152, 2020.
- [364] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, 2006.

- [365] Glen D Rayner and Helen L MacGillivray. Numerical maximum likelihood estimation for the g-and-k and generalized g-and-h distributions. *Statistics and Computing*, 12(1):57–75, 2002.
- [366] Lewis J Rendell, Adam M Johansen, Anthony Lee, and Nick Whiteley. Global consensus monte carlo. *Journal of Computational and Graphical Statistics*, 30(2):249–259, 2020.
- [367] Silvio R Rendon. Fixed and random effects in classical and bayesian regression. *Oxford Bulletin of Economics and Statistics*, 75(3):460–476, 2013.
- [368] Sylvia Richardson and Peter J Green. On bayesian analysis of mixtures with an unknown number of components (with discussion). *Journal of the Royal Statistical Society: series B (statistical methodology)*, 59(4):731–792, 1997.
- [369] Maxime Rischard, Zach Branson, Luke Bornn, and Luke W. Miratrix. Do school districts affect nyc house prices? identifying border differences using a bayesian nonparametric approach to geographic regression discontinuity designs. *Journal of the American Statistical Association*, 115(531):1149–1165, 2020.
- [370] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [371] Christian P Robert and George Casella. *Introducing monte carlo methods with R*, volume 18. Springer, 2010.
- [372] Christian P. Robert and George Casella. *Monte Carlo Statistical Methods*. Springer, New York, 2nd edition, 2011.
- [373] G. O. Roberts, A. Gelman, and W. R. Gilks. Weak convergence and optimal scaling of random walk Metropolis algorithms. *The Annals of Applied Probability*, 7(1):110–120, 1997.
- [374] Gareth O. Roberts and Richard L. Tweedie. Exponential convergence of langevin distributions and their discrete approximations. *Bernoulli*, 2(4):341–363, 1996.
- [375] Veronika Ročková. On semi-parametric inference for BART. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 8137–8146. PMLR, 2020.
- [376] Paul R. Rosenbaum and Donald B. Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- [377] P. Rossi. *Package bayesm*, 2017.

- [378] Peter E. Rossi. *Bayesian Non- and Semi-Parametric Methods and Applications*. Princeton University Press, Princeton, NJ, 2014.
- [379] Peter E Rossi, Greg M Allenby, and Rob McCulloch. *Bayesian statistics and marketing*. John Wiley & Sons, 2012.
- [380] Jonathan Roth, Pedro H. C. Sant’Anna, Alyssa Bilinski, and John Poe. What’s trending in difference-in-differences? a synthesis of the recent econometrics literature. *Journal of Econometrics*, 235(2):2218–2244, 2023.
- [381] Veronika Ročková and Edward I. George. The spike-and-slab lasso. *Journal of the American Statistical Association*, 113(521):431–444, 2018.
- [382] Veronika Ročková and Enakshi Saha. On theory for bart. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 2839–2848. PMLR, 2019.
- [383] Veronika Ročková and Stephanie van der Pas. Posterior concentration for bayesian regression trees and forests. *The Annals of Statistics*, 48(4):2103–2130, 2020.
- [384] Andrew Donald Roy. Some thoughts on the distribution of earnings. *Oxford economic papers*, 3(2):135–146, 1951.
- [385] Donald B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688–701, 1974.
- [386] Donald B. Rubin. Bayesian inference for causal effects: The role of randomization. *The Annals of Statistics*, 6(1):34–58, 1978.
- [387] Donald B Rubin. Bayesianly justifiable and relevant frequency calculations for the applies statistician. *The Annals of Statistics*, 12(4):1151–1172, 1984.
- [388] Donald B. Rubin. Using the sir algorithm to simulate posterior distributions. In J. M. Bernardo, M. H. DeGroot, D. V. Lindley, and A. F. M. Smith, editors, *Bayesian Statistics 3*, pages 395–402. Oxford University Press, 1988.
- [389] Donald B. Rubin. Neyman (1923) and causal inference in experiments and observational studies. *Statistical Science*, 5(4):472–480, 1990.
- [390] Donald B. Rubin. Teaching statistical inference for causal effects in experiments and observational studies. *Journal of Educational and Behavioral Statistics*, 29(3):343–367, 2004.

- [391] Donnald B. Rubin. The Bayesian bootstrap. *The Annals of Statistics*, 9(1):130–134, 1981.
- [392] Håvard Rue, Sara Martino, and Nicolas Chopin. Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 71(2):319–392, 2009.
- [393] Håvard Rue, Andrea Riebler, Sigrunn H Sørbye, Janine B Illian, Daniel P Simpson, and Finn K Lindgren. Bayesian computing with inla: a review. *Annual Review of Statistics and Its Application*, 4(1):395–421, 2017.
- [394] Olli Saarela, Léo R. Belzile, and David A. Stephens. A bayesian view of doubly robust causal inference. *Biometrika*, 103(3):667–681, 2016.
- [395] John Denis Sargan. The estimation of economic relationships using instrumental variables. *Econometrica*, 26(3):393–415, 1958.
- [396] L. J. Savage. *The foundations of statistics*. John Wiley & Sons, Inc., New York, 1954.
- [397] Susanne M. Schennach. Bayesian exponentially tilted empirical likelihood. *Biometrika*, 92(1):31–46, 2005.
- [398] Susanne M. Schennach. Point estimation with exponentially tilted empirical likelihood. *The Annals of Statistics*, 35(2):634–672, 2007.
- [399] Robert Schlaifer and Howard Raiffa. *Applied statistical decision theory*. Wiley New York, 1961.
- [400] Gideon Schwarz. Estimating the dimension of a model. *The annals of statistics*, pages 461–464, 1978.
- [401] Steven L. Scott, Alexander W. Blocker, Fernando V. Bonassi, Hugh A. Chipman, Edward I. George, and Robert E. McCulloch. Bayes and big data: The consensus monte carlo algorithm. *International Journal of Management Science and Engineering Management*, 11(2):78–88, 2016.
- [402] Steven L Scott, Alexander W Blocker, Fernando V Bonassi, Hugh A Chipman, Edward I George, and Robert E McCulloch. Bayes and big data: The consensus monte carlo algorithm. In *Big Data and Information Theory*, pages 8–18. Routledge, 2022.
- [403] Thomas Sellke, MJ Bayarri, and James O Berger. Calibration of p values for testing precise null hypotheses. *The American Statistician*, 55(1):62–71, 2001.
- [404] Steve Selvin. A problem in probability (letter to the editor). *The American Statistician*, 29(1):67–71, 1975.

- [405] Francesco Serafini. Multinomial logit models with inla. *R-INLA Tutorial*. <https://inla.r-inla-download.org/r-inla.org/doc/vignettes/multinomial.pdf>, 2019.
- [406] M. Serna Rodríguez, A. Ramírez Hassan, and A. Coad. Uncovering value drivers of high performance soccer players. *Journal of Sport Economics*, 20(6):819–849, 2019.
- [407] J. Sethuraman. A constructive definition of Dirichlet priors. *Statistica Sinica*, 4:639–650, 1994.
- [408] Burr Settles. *Active Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2012.
- [409] Neil Shephard. Partial non-gaussian state space. *Biometrika*, 81(1):115–131, 1994.
- [410] Robert H. Shumway and David S. Stoffer. *Time Series Analysis and Its Applications: With R Examples*. Springer, Cham, Switzerland, 4th edition, 2017.
- [411] Susan J Simmons, Fang Fang, Qijun Fang, and Karl Ricanek. Markov chain Monte Carlo model composition search strategy for quantitative trait loci in a Bayesian hierarchical model. *World Academy of Science, Engineering and Technology*, 63:58–61, 2010.
- [412] Christopher A Sims. Macroeconomics and reality. *Econometrica: journal of the Econometric Society*, pages 1–48, 1980.
- [413] Sisson, Scott A, Fan, Yanan, Beaumont, and Mark. *Handbook of Approximate Bayesian Computation*. CRC Press, 2018.
- [414] Scott A Sisson, Yanan Fan, and Mark M Tanaka. Sequential Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences of the USA*, 104(6):1760–1765, 2007.
- [415] A. F. M. Smith. A General Bayesian Linear Model. *Journal of the Royal Statistical Society. Series B (Methodological)*, 35(1):67–75, 1973.
- [416] Adrian FM Smith and Alan E Gelfand. Bayesian statistics without tears: a sampling–resampling perspective. *The American Statistician*, 46(2):84–88, 1992.
- [417] Ercan Solak, Roderick Murray-Smith, W. E. Leithead, D. J. Leith, and C. E. Rasmussen. Derivative observations in gaussian process models of dynamic systems. In *Advances in Neural Information Processing Systems*, pages 1033–1040. MIT Press, 2003.
- [418] Qifan Song, Yan Sun, Mao Ye, and Faming Liang. Extended stochastic gradient markov chain monte carlo for large-scale bayesian variable selection. *Biometrika*, 107(4):997–1004, 2020.

- [419] Rodney A. Sparapani, Charles Spanbauer, and Robert McCulloch. Non-parametric machine learning and efficient computation with bayesian additive regression trees: the bart r package. *Journal of Statistical Software*, 97(1):1–66, 2021.
- [420] David J Spiegelhalter, Nicola G Best, Bradley P Carlin, and Angelika Van Der Linde. Bayesian measures of model complexity and fit. *Journal of the royal statistical society: Series b (statistical methodology)*, 64(4):583–639, 2002.
- [421] Douglas Staiger and James H. Stock. Instrumental variables regression with weak instruments. *Econometrica*, 65(3):557–586, 1997.
- [422] Stan Development Team. shinystan: Interactive visual and numerical diagnostics and posterior analysis for Bayesian models., 2017. R package version 2.3.0.
- [423] Stan Development Team. Stan modeling language users guide and reference manual, 2024., 2024.
- [424] Stephen Stigler. Richard price, the first bayesian. *Statistical Science*, 33(1):117–125, 2018.
- [425] James H. Stock and Motohiro Yogo. Asymptotic distributions of instrumental variables statistics with many instruments. In Donald W. K. Andrews and James H. Stock, editors, *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*, pages 109–120. Cambridge University Press, 2005.
- [426] Andrew M. Stuart and Aretha L. Teckentrup. Posterior consistency for gaussian process approximations of bayesian posterior distributions. *Mathematics of Computation*, 87(310):721–753, 2018.
- [427] Nicholas Syring and Ryan Martin. Calibrating general posterior credible regions. *Biometrika*, 106(2):479–486, 2019.
- [428] M. A. Tanner and W. H. Wong. The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association*, 82(398):528–540, 1987.
- [429] Simon Tavaré, David J Balding, Robert C Griffiths, and Peter Donnelly. Inferring coalescence times from DNA sequence data. *Genetics*, 145(2):505–518, 1997.
- [430] Yee Whye Teh, Alexandre H. Thiery, and Sebastian J. Vollmer. Consistency and fluctuations for stochastic gradient langevin dynamics. *Journal of Machine Learning Research*, 17(1):193–225, 2016.
- [431] Samuel Thomas and Wanzhu Tu. Learning hamiltonian monte carlo in r. *The American Statistician*, 75(4):403–413, 2021.

- [432] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [433] Luke Tierney. Markov chains for exploring posterior distributions. *The Annals of Statistics*, pages 1701–1728, 1994.
- [434] Luke Tierney and Joseph B Kadane. Accurate approximations for posterior moments and marginal densities. *Journal of the American Statistical Association*, 81(393):82–86, 1986.
- [435] Luke Tierney, Robert E Kass, and Joseph B Kadane. Fully exponential laplace approximations to expectations and variances of nonpositive functions. *Journal of the american statistical association*, 84(407):710–716, 1989.
- [436] Minh-Ngoc Tran, David J. Nott, and Robert Kohn. Variational bayes with intractable likelihood. *Journal of Computational and Graphical Statistics*, 26(4):873–882, 2017.
- [437] Regina Tüchler. Bayesian variable selection for logistic models using auxiliary mixture sampling. *Journal of Computational and Graphical Statistics*, 17(1):76–94, 2008.
- [438] A. Tversky and D. Kahneman. Judgement under uncertainty: heuristics and biases. *Science*, 185:1124–1131, 1974.
- [439] Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- [440] Martijn Van Hasselt. Bayesian inference in a sample selection model. *Journal of Econometrics*, 165(2):221–232, 2011.
- [441] Isabella Verdinelli and Larry Wasserman. Computing bayes factors using a generalization of the savage-dickey density ratio. *Journal of the American Statistical Association*, 90(430):614–618, 1995.
- [442] Sebastian J. Vollmer, Konstantinos C. Zygalakis, and Yee Whye Teh. Exploration of the (non-)asymptotic bias and variance of stochastic gradient langevin dynamics. *Journal of Machine Learning Research*, 17(159):1–48, 2016.
- [443] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- [444] Stephen G Walker. Modern bayesian asymptotics. *Statistical Science*, pages 111–117, 2004.

- [445] Stephen G. Walker. New approaches to bayesian consistency. *Annals of Statistics*, 32(5):2028–2043, 2004.
- [446] Xiangyu Wang and David B Dunson. Parallelizing mcmc via weierstrass sampler. *arXiv preprint arXiv:1312.4605*, 2013.
- [447] Ronald L. Wasserstein and Nicole A. Lazar. The ASA’s statement on p-values: context, process and purpose. *The American Statistician*, 2016.
- [448] Daniel Wegmann, Christoph Leuenberger, and Laurent Excoffier. Efficient approximate Bayesian computation coupled with Markov chain Monte Carlo without likelihood. *Genetics*, 182(4):1207–1218, 2009.
- [449] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688. Citeseer, 2011.
- [450] Mike West and Jeff Harrison. *Bayesian forecasting and dynamic models*. Springer Science & Business Media, 2006.
- [451] Andrew G. Wilson and Hannes Nickisch. Kernel interpolation for scalable structured gaussian processes (kiss-gp). In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, volume 37, pages 1775–1784. PMLR, 2015.
- [452] R. Winkelmann. Health care reform and the number of doctor visits - An econometric analysis. *Journal of Applied Econometrics*, 19(4):455–472, 2004.
- [453] Simon N Wood. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310):1102–1104, 2010.
- [454] P. Woodward. BugsXLA: Bayes for the common man. *Journal of Statistical Software*, 14(5):1–18, 2005.
- [455] Jeffrey M Wooldridge. *Econometric analysis of cross section and panel data*. MIT press, 2010.
- [456] Jeffrey M. Wooldridge. *Introductory Econometrics: A Modern Approach*. Cengage Learning, Mason, Ohio: South-Western, fifth edition, 2012.
- [457] Jeffrey M. Wooldridge. *Introductory Econometrics: A Modern Approach*. Cengage Learning, Boston, MA, 6th edition, 2016.
- [458] Tomasz Woźniak. Bayesian vector autoregressions. *Australian Economic Review*, 49(3):365–380, 2016.
- [459] Tomasz Woźniak. *bsvars: Bayesian Estimation of Structural Vector Autoregressive Models*, 2024. R package version 3.1.

- [460] Changye Wu and Christian P Robert. Average of recentered parallel mcmc for big data. *arXiv preprint arXiv:1706.04780*, 2017.
- [461] Wang Xiaolei and Tomasz Woźniak. *bsvarSIGNS: Bayesian SVARs with Sign, Zero, and Narrative Restrictions*, 2024. R package version 1.0.1.
- [462] Yuling Yao. Bayesian aggregation. In *Wiley StatsRef: Statistics Reference Online*, pages 1–13. John Wiley & Sons, Ltd, 2021.
- [463] Yuling Yao, Aki Vehtari, Daniel Simpson, and Andrew Gelman. Using stacking to average bayesian predictive distributions. *Bayesian Analysis*, 13(3):917–1003, 2018.
- [464] Andrew Yiu, Robert J. B. Goudie, and Brian D. M. Tom. Inference under unequal probability sampling with the bayesian exponentially tilted empirical likelihood. *Biometrika*, 107(4):857–873, 2020.
- [465] A. Zellner. *Introduction to Bayesian inference in econometrics*. John Wiley & Sons Inc., 1996.
- [466] Arnold Zellner. On assessing prior distributions and bayesian regression analysis with g-prior distributions. *Bayesian inference and decision techniques*, 1986.
- [467] Arnold Zellner. Bayesian method of moments (bmom) analysis of mean and regression models. In Jack C. Lee, Wesley O. Johnson, and Arnold Zellner, editors, *Modelling and Prediction Honoring Seymour Geisser*, pages 61–72. Springer, New York, NY, 1996.
- [468] Fengshuo Zhang and Chao Gao. Convergence rates of variational posterior distributions. *The Annals of Statistics*, 48(4):2180–2207, 2020.
- [469] Tong Zhang. From  $varepsilon$ -Entropy to KL-Entropy: Analysis of Minimum Information Complexity Density Estimation. *Annals of Statistics*, 34(5):2180–2210, 2006.
- [470] S. Ziliak. Guinnessometrics; the Economic Foundation of student’s t. *Journal of Economic Perspectives*, 22(4):199–216, 2008.
- [471] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.



# A

---

---

## *Appendix*

---

**TABLE A.1**

Libraries and commands in BEsmarter GUI.

Univariate models			
Model	Library	Command	Reference
Normal	MCMCpack	MCMCregress	[295]
Logit	MCMCpack	MCMClogit	[295]
Probit	bayesm	rbprobitGibbs	[377]
Multinomial(Mixed) Probit	bayesm	rmnpGibbs	[377]
Multinomial(Mixed) Logit	bayesm	rmnlIndepMetrop	[377]
Ordered Probit	bayesm	rordprobitGibbs	[377]
Negative Binomial(Poisson)	bayesm	rnegbinRw	[377]
Tobit	MCMCpack	MCMCtobit	[295]
Quantile	MCMCpack	MCMCquantreg	[295]
Bayesian bootstrap	bayesboot	bayesboot	[22]
Multivariate models			
Model	Library	Command	Reference
Multivariate	bayesm	rmultireg	[377]
Seemingly Unrelated Regression	bayesm	rsurGibbs	[377]
Instrumental Variable	bayesm	rivGibbs	[377]
Bivariate Probit	bayesm	rmvpGibbs	[377]
Time series models			
Model	Library	Command	Reference
Normal dynamic linear model	dlm	dlmGibbsDIGs	[335]
ARMA	bayesforecast	stan_sarima	[7]
Stochastic volatility models	stochvol	svsample	[211]
VAR	bvarools	draw_posterior	[312]
Hierarchical longitudinal models			
Model	Library	Command	Reference
Normal	MCMCpack	MCMChregress	[295]
Logit	MCMCpack	MCMChlogit	[295]
Poisson	MCMCpack	MCMChpoisson	[295]
Bayesian model averaging			
Model	Library	Command	Reference
Normal (BIC)	BMA	bicreg	[348]
Normal (MC <sup>3</sup> )	BMA	MC3.REG	[348]
Normal (instrumental variables)	ivbma	ivbma	[262]
Normal (Dynamic BMA)	dma	dma	[301]
Logit (BIC)	BMA	bic.glm	[348]
Gamma (BIC)	BMA	bic.glm	[348]
Poisson (BIC)	BMA	bic.glm	[348]
Diagnostics			
Diagnostic	Library	Command	Reference
Trace plot	coda	traceplot	[339]
Autocorrelation plot	coda	autocorr.plot	[339]
Geweke test	coda	geweke.diag	[339]
Raftery & Lewis test	coda	raftery.diag	[339]
Heidelberger & Welch test	coda	heidel.diag	[339]

**TABLE A.2**Datasets templates in folder *DataSim*.

<b>Univariate models</b>		
<b>Model</b>	<b>Data set file</b>	<b>Data set simulation</b>
Normal	11SimNormalmodel.csv	11SimNormal.R
Logit	12SimLogitmodel.csv	12SimLogit
Probit	13SimProbitmodel.csv	13SimProbit.R
Multinomial(Mixed) Probit	14SimMultProbm model.csv	14SimMultinomialProbit.R
Multinomial(Mixed) Logit	15SimMultLogitmodel.csv	15SimMultinomialLogit.R
Ordered Probit	16SimOrderedProbitmodel.csv	16SimOrderedProbit.R
Negative Binomial(Poisson)	17SimNegBinmodel.csv	17SimNegBin.R
Tobit	18SimTobitmodel.csv	18SimTobit.R
Quantile	19SimQuantilemodel.csv	19SimQuantile.R
Bayesian bootstrap	110SimBootstrapmodel.csv	110SimBootstrapmodel.R
<b>Multivariate models</b>		
<b>Model</b>	<b>Data set file</b>	<b>Data set simulation</b>
Multivariate	21SimMultivariate.csv	21SimMultReg.R
Seemingly Unrelated Regression	22SimSUR.csv	22SimSUR.R
Instrumental Variable	23SimIV.csv	23SimIV.R
Bivariate Probit	24SimMultProbit.csv	24SimMultProbit.R
<b>Time series models</b>		
<b>Model</b>	<b>Data set file</b>	<b>Data set simulation</b>
Dynamic Linear Models	31SimDynamicLinearModel.csv	31SimDynamicLinearModel.R
ARMA	32SimARMAmodels.csv	32SimARMAmodels.R
Stochastic volatility Models	33SimStochasticVolatility.csv	33SimStochasticVolatility.R
VAR	34SimVARmodels.csv	34SimVARmodels.R
<b>Hierarchical longitudinal models</b>		
<b>Model</b>	<b>Data set file</b>	<b>Data set simulation</b>
Normal	41SimLogitudinalNormal.csv	41SimLogitudinalNormal.R
Logit	42SimLogitudinalLogit.csv	42SimLogitudinalLogit.R
Poisson	43SimLogitudinalPoisson.csv	43SimLogitudinalPoisson.R
<b>Bayesian model averaging</b>		
<b>Model</b>	<b>Data set file</b>	<b>Data set simulation</b>
Normal (BIC)	511SimNormalBMA.csv	511SimNormalBMA.R
Normal (MC <sup>3</sup> )	512SimNormalBMA.csv	512SimNormalBMA.R
Normal (instrumental variables)	513SimNormalBMAivYXW.csv 513SimNormalBMAivZ.csv	513SimNormalBMAiv.R
Normal (Dynamic BMA)	514SimDynamicBMA.csv 514SimModels.csv	514SimDynamicBMA.R
Logit (BIC)	52SimLogitBMA.csv	52SimLogitBMA.R
Gamma (BIC)	53SimGammaBMA.csv	53SimGammaBMA.R
Poisson (BIC)	53SimPoissonBMA.csv	53SimPoissonBMA.R

**TABLE A.3**  
Real datasets in folder *DataApp*.

Univariate models		
Model	Data set file	Dependent variable
Normal	1ValueFootballPlayers.csv	log(Value)
Logit	2HealthMed.csv	Hosp
Probit	2HealthMed.csv	Hosp
Multinomial(Mixed) Probit	Fishing.csv	mode
Multinomial(Mixed) Logit	Fishing.csv	mode
Ordered Probit	2HealthMed.csv	MedVisPrevOr
Negative Binomial(Poisson)	2HealthMed.csv	MedVisPrev
Tobit	1ValueFootballPlayers.csv	log(ValueCens)
Quantile	1ValueFootballPlayers.csv	log(Value)
Bayesian bootstrap	1ValueFootballPlayers.csv	log(Value)
Multivariate models		
Model	Data set file	Dependent variable
Multivariate	4Institutions.csv	logpcGDP95 and PAER
Seemingly Unrelated Regression	5Institutions.csv	logpcGDP95 and PAER
Instrumental Variable	6Institutions.csv	logpcGDP95 and PAER
Bivariate Probit	7HealthMed.csv	$y = [\text{Hosp } \text{SHI}]'$
Multivariate models		
Model	Data set file	Dependent variable
Multivariate	4Institutions.csv	logpcGDP95 and PAER
Seemingly Unrelated Regression	5Institutions.csv	logpcGDP95 and PAER
Instrumental Variable	6Institutions.csv	logpcGDP95 and PAER
Bivariate Probit	7HealthMed.csv	$y = [\text{Hosp } \text{SHI}]^\top$
Time series models		
Model	Data set file	Dependent variable
Dynamic Linear Models	16INTDEF.csv	$\Delta(i_3)$
ARMA	16INTDEF.csv	$\Delta(i_3)$
Stochastic volatility models	17ExcRate.csv	USDEUR
VAR	18USAfiscal	$y_t = [\Delta(ttr_t) \Delta(gs_t) \Delta(gdp_t)]^\top$
Hierarchical longitudinal models		
Model	Data set file	Dependent variable
Normal	8PublicCap.csv	log(gsp)
Logit	9VisitDoc.csv	DocVis
Poisson	9VisitDoc.csv	DocNum
Bayesian model averaging		
Model	Data set file	Dependent variable
Normal (BIC)	10ExportDiversificationHHI.csv	avghhi
Normal ( $MC^3$ )	10ExportDiversificationHHI.csv	avghhi
Normal (instrumental variables)	11ExportDiversificationHHI.csv 12ExportDiversificationHHIInstr.csv	avghhi and avgldpcap
Logit (BIC)	13InternetMed.csv	internet
Gamma (BIC)	14ValueFootballPlayers.csv	log market value
Poisson (BIC)	15Fertile2.csv	ceb