Half Title

Introduction to Bayesian Inference: A GUIded tour using R

Title Page

Introduction to Bayesian Inference: A GUIded tour using R

by Andrés Ramírez-Hassan, PhD. Statistical Science.

LOC Page



Introduction



FIGURE 1
Supposedly portrait of Thomas Bayes.

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

Introduction to Bayesian inference: A GUIded tour using R mainly targets those who want to apply Bayesian regression analysis having a good conceptual and formal understanding, but not necessarily having time to develop programming skills. Thus, this book provides a graphical user interface (GUI) to carry out Bayesian regression in a very friendly environment. The book also provides the basic theory, and its code implementation using R software [78], some econometric/statistical applications to highlight the potential of Bayesian regression, and theory and computational exercises, for those who are interested in developing more complex models. In particular, this book contains the mathematical proofs step by step of the basic model, which are the base for obtaining the most relevant mathematical results of more complex models.

Our GUI is based on an interactive web application using shiny [18], and some packages in **R**. Users can estimate univariate, multivariate, time series, panel data (longitudinal) and Bayesian model average models using our GUI.

viii Introduction

In addition, it gives basic summaries and formal and graphical diagnostics of the posterior chains. Our GUI can be run in any operating system, and is freely available at https://github.com/besmarter/BSTApp.

Users can get simulated and real datasets in the folders $\mathbf{DataSim}$, and $\mathbf{DataApp}$, respectively. The former folder also includes the files that were used to simulate different processes, so, the population parameters are available, and as a consequence, these files can be used as a pedagogical tool to show some statistical properties. The latter folder contains the datasets used in our applications. Users should use these datasets as templates to structuring their own datasets. Simply type $\mathbf{shiny::runGitHub("besmarter/BSTApp"}$, $\mathbf{launch.browser=T)}$ in the \mathbf{R} package console or any \mathbf{R} code editor to run our $\mathbf{GUI.}^1$

This book has three parts. The first part covers theory (conceptual and mathematical), programming, and simulation foundations (chapters 1 to 5). The second part focuses on applications of regression analysis, with particular emphasis on the computational aspect of obtaining draws from the posterior distributions (chapters 6 to 11). The third part provides an introductory treatment of advanced methods in Bayesian inference (chapters 12 to 15). I show in some detail the mathematical deductions in the first part of the book, whereas I do not show any proof in the second and third parts. However, same mathematical steps in the first part can be used to find the results of parts two and three of the book. I also show three levels regarding computational implementation in the second part of the book: programming ourselves the algorithms, using Bayesian R packages, and using our GUI.

Chapter 1 begins with an introduction to formal concepts in Bayesian inference starting with the Bayes' rule, all its components with their formal definitions and basic examples. Then, it presents the basics of Bayesian inference based on decision theory under uncertainty. Chapter 2 presents conceptual differences between Bayesian and Frequentist statistical approaches, and a historical and philosophical perspective about Bayesian statistics and econometrics highlighting differences compared to the Frequentist approach. Chapter 3 presents the differences between the objective and subjective schools in Bayesian inference. Particular attention is put to elicitation techniques, that is, how to transform expert knowledge into prior probabilistic statements. In Chapter 4 I introduce conjugate families in basic statistical models, solving them analytically and computationally. Simulation based methods are shown in Chapter 5, these algorithms are very important in modern Bayesian inference as most realistic models do not have standard forms or analytical solutions. I present our graphical user interface in Chapter 6, and univariate and multivariate regression models are presented in chapters 7 and 8, respectively. Chapter 9 presents the state-space representation of time series models, and Chapter 10 presents Bayesian panel data (longitudinal) models. Chapter 11 introduces Bayesian model averaging. In the third part, there are

¹I strongly recommend to type the code line rather than copy and paste it.

Introduction ix

Chapter 12 introducing hierarchical models, Chapter 13 shows causal inference, Chapter 14 shows Bayesian methods in machine learning algorithms, and Chapter 15 describes some recent methodological developments such as approximate Bayesian computation (ABC), variational Bayes (VB), integrated nested Laplace approximations (INLA), and Bayesian exponential tilted empirical likelihood (BETEL).

About me

My name is Andrés Ramírez-Hassan, I am an applied and theory econometrician working as a Distinguished Professor in the School of Finance, Economics and Government at Universidad EAFIT (Medellín, Colombia). I got a PhD in Statistical Science, a masters degree in Finance, and another in Economics, and also 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. Having completed my PhD degree, much of my research has been in the area of Bayesian Econometrics with applications in crime, finance, health, sports and utilities. My work has been published (or is forthcoming) in 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 highly regarded international research outlets.

I founded **BEsmarter** –**B**ayesian **E**conometrics: simulations, **m**odels and applications to **research**, **teaching** and **e**ncoding 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
- And where one of the main concerns of science is to solve social problems

mail: aramir21@gmail.com / aramir21@eafit.edu.co website: http://www.besmarter-team.org / https://sites.google.com/view/arh-bayesian

 ${\bf x}$ Introduction



FIGURE 2

This book is licensed under the $Creative\ Commons\ Attribution-NonCommercial-Share Alike\ 4.0$ International License.

Contents

In	ntroduction	vii
Fc	preword	xi
Pı	reface	xiii
$\mathbf{S}_{\mathbf{J}}$	ymbols	$\mathbf{x}\mathbf{v}$
Ι	Foundations: Theory, simulation methods and programming	1
1	Basic formal concepts 1.1 The Bayes' rule 1.2 Bayesian framework: A brief summary of theory 1.2.1 Example: Health insurance 1.3 Bayesian reports: Decision theory under uncertainty 1.3.1 Example: Health insurance continues 1.4 Summary 1.5 Exercises	3 8 14 27 30 32 32
2	Conceptual differences of the Bayesian and Frequentist approaches 2.1 The concept of probability 2.2 Subjectivity is not the key 2.3 Estimation, hypothesis testing and prediction 2.4 The likelihood principle 2.5 Why is not the Bayesian approach that popular? 2.6 A simple working example 2.6.1 Example: Math test 2.7 Summary 2.8 Exercises	35 36 37 41 42 44 46 47 48
3	Objective and subjective Bayesian approaches	51
4	Cornerstone models: Conjugate families 4.1 Motivation of conjugate families	53 53 54

viii	Contents

	4.2	Conjugate prior to exponential family	58
	4.3	4.2.1 Examples: Theorem 4.2.1	59
		model	75
	4.4	Multivariate linear regression: The conjugate normal-normal/inve	
		Wishart model	87
	4.5	<u>Summary</u>	91
	4.6	Exercises	91
5	Sim	ulation methods	95
	5.1	The inverse transform method	95
	5.2	Method of composition	95
	5.3	Accept and reject algorithm	95
	5.4	Importance sampling	95
	5.5	Markov chain Monte Carlo methods	95
		5.5.1 Some theory	95
		5.5.2 Gibbs sampler	95
		5.5.3 Metropolis-Hastings	95
	5.6	Sequential Monte Carlo	95
	5.7	Hamiltonian Monte Carlo	95
	5.8	Convergence diagnostics	95
Η	\mathbf{R}	egression models: A GUIded tour	97
6	Gra	phical user interface	99
	6.1	Introduction	99
	6.2	Univariate models	100
	6.3	Multivariate models	104
	6.4	Time series model	106
	6.5	Longitudinal (panel) models	106
	6.6	Bayesian model average	106
	6.7	Warning	109
7	Uni	variate models	111
	7.1	The Gaussian linear model	111
	7.2	The logit model	116
	7.3	The probit model	121
	7.4	The multinomial probit model	124
	7.5	The multinomial logit model	128
	7.6	Ordered probit model	132
	7.7	Negative binomial model	136
	7.8	Tobit model	141
	7.9	Quantile regression	145
	7.10	Bayesian bootstrap regression	148
		Summary	150
		Exercises	150

Contents	ix
8 Multivariate models 8.1 Multivariate regression 8.2 Seemingly unrelated regression 8.3 Instrumental variable 8.4 Multivariate probit model 8.5 Summary 8.6 Exercises	153 153 159 164 168 174 174
9 Time series models	177
10 Panel data models	179
11 Bayesian model average 11.1 Calculating the marginal likelihood	181 181 181 181 181
gramming	183
gramming 12 Hierarchical models 12.1 Finite mixtures	183185185
gramming 12 Hierarchical models 12.1 Finite mixtures	183 185 185 185
gramming 12 Hierarchical models 12.1 Finite mixtures	183 185 185 187 189 189 189

Foreword

Preface

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

To instructors and students

This book is divided in three parts, foundations (chapters 1 to 5), regression analysis (chapters 6 to 11), and Advanced methods (chapters 12 to 15). Our graphical user interface (GUI) targets the second part. This can be download at https://github.com/besmarter/BSTApp. Instructors and students can have all codes, simulated and real data sets are also there. To install our GUI just type shiny::runGitHub("besmarter/BSTApp", launch.browser=T) in the R package console or any R code editor, and execute it.

Students should have some basic knowledge in probability theory and statistics, particularly, regression analysis. It is strongly recommended to have some familiarity with standard univariate and multivariate probability distributions.

I included some formal and computational exercises at the end of each chapter. This would help students to have a better understanding of the material shown in each chapter. A manual with the solutions of exercises accompanies this book.

Instructors can use this book as a text in a course of introduction to Bayesian Econometrics/Statistics with a high emphasis on implementation and applications. This book is complementary, rather than substitute, of excellent books in the topic such as [37, 89, 45, 42, 58] and [56].

Acknowledgments

I started our GUI in the 2016 after being diagnosed with cervical dystonia. I used to work in this side project on weekends, I named this time "nerd weekends", and it was a kind of release from my health condition. Once I got better, I invited Mateo Graciano, my former student, business partner and friend, to be part of the project, he helped me a lot developing our GUI, and I am enormously thankful to Mateo. I would also like to thank members

xiv Preface

of the BEsmarter research group from Universidad EAFIT, and NUMBATs members from Monash University for your comments and recommendations to improve our GUI.

This book is an extension of the paper A GUIded tour of Bayesian regression [87], which is a brief user guide of our GUI. So, I decided to write this book to show the underlying theory and codes in our GUI, and use it as a text book in my course in Bayesian econometrics/statistics. I acknowledge and offer my gratitude to my students in this subject, their insight and thoughtful questions have helped me to get a better understanding of this material.

I also thank Chris Parmeter for your suggestions about how to present our user guide, Professor Raul Pericchi and Juan Carlos Correa who introduced me to Bayesian statistics, Liana Jacobi and Chun Fung Kwok (Jackson) from the University of Melbourne and David Frazier from Monash University for nice talks and amazing collaborations in Bayesian Econometrics/statistics, Professor Peter Diggle to support my career, and particularly, Professor Gael Martin, who gave me a chance to work with her, she is an inspiring intellectual figure. Finally, my colleagues and staff from Universidad EAFIT have always given me their support.

To my parents, Orlando and Nancy, who have given me their unconditional support. They have taught me that the primary aspect of the human being's spiritual evolution is humility. I am in my way to learn this.

Symbols

Symbol Description

\neg	Negation symbol	argmax	Argument of the maximum
\propto	Proportional symbol	argmin	Argument of the minimum
\perp	Independence symbol	tr	Trace operator
$\mathcal R$	The Real set	vec	Vectorization operator
Ø	Empty set	\lim	Limit
1	Indicator function	\otimes	Kronecker product
P	Probability measure	$\operatorname{diag}\{\cdot\}$	Diagonal matrix
:=	Is defined as	$dim\{\cdot\}$	Dimension of an object

Part I

Foundations: Theory, simulation methods and programming

Basic formal concepts

We introduce formal concepts in Bayesian inference starting with the Bayes' rule, all its components with their formal definitions and basic examples. In addition, we present some nice features of Bayesian inference such as Bayesian updating, and asymptotic sampling properties, and the basics of Bayesian inference based on decision theory under uncertainty, presenting important concepts like loss function, risk function and optimal rules.

1.1 The Bayes' rule

As expected the point of departure to perform Bayesian inference is the Bayes' rule, which is the Bayes' solution to the inverse probability of causes, this rule combines prior beliefs with objective probabilities based on repeatable experiments. In this way, we can 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 times the marginal probability of A_i over the marginal probability of B,

$$P(A_i|B) = \frac{P(A_i, B)}{P(B)}$$

$$= \frac{P(B|A_i) \times P(A_i)}{P(B)},$$
(1.1)

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

In the Bayesian framework, B is sample information that updates a probabilistic statement about an unknown object A_i following probability rules. This is done by means of the Bayes' rule using prior "beliefs" about A_i , that is, $P(A_i)$, sample information relating B to the particular state of the nature A_i through a probabilistic statement, $P(B|A_i)$, and the probability of observing that specific sample information P(B).

¹Observe that I use the term "Bayes' rule" rather than "Bayes' theorem". It was Laplace [60] who actually generalized the Bayes' theorem [7]. His generalization is named the Bayes' rule.

Let's see 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%, P(+|disease) = 0.98. On the other hand, the prior information regarding being infected with this disease comes from a base incidence rate that is equal to 0.002, that is P(disease) = 0.002. Then, what is the probability of being actually infected?

This is an example of the base rate fallacy, where having a positive test result from a disease whose base incidence rate is tiny gives a low probability of actually having the disease.

The key to answer the question is based on understanding the difference between the probability of having the disease given a positive result, P(disease|+), versus the probability of a positive result given the disease, P(+|disease). The former is the important result, and the Bayes' rule help us to get the answer. Using the Bayes' rule (equation 1.1):

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

where $P(+) = P(+|\text{disease}) \times P(\text{disease}) + P(+|\neg \text{disease}) \times P(\neg \text{disease})$.

R code. The base rate fallacy

```
PD <- 0.002 # Probability of disease
PPD <- 0.98 # True positive (Sensitivity)
PDP <- PD * PPD / (PD * PPD + (1 - PD)*(1 - PPD))

paste("Probability of disease given a positive test is", sep = " ", round(PDP, 2))

"Probability of disease given a positive test is 0.09"
```

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

Another interesting example, which is at the heart of the origin of the Bayes' theorem [7], is related to the existence of God [99]. The Section X of David Hume's "An Inquiry concerning Human Understanding, 1748" is named Of Miracles. There, Hume argues that when someone claims to have seen a

 $^{^2\}neg$ is the negation symbol. In addition, we have that $P(B|A)=1-P(B|A^c)$ in this example, where A^c is the complement of A. However, it is not always the case that $P(B|A)\neq 1-P(B|A^c)$.

The Bayes' rule 5

miracle, this is poor evidence it actually happened, since it goes against what we see every day. Then, Richard Price, who actually finished and published "An essay towards solving a problem in the doctrine of chances" in 1763 after Bayes died in 1761, argues against Hume saying that there is a huge difference between *impossibility* as used commonly in conversation and *physical impossibility*. Price used an example of a dice with a million sides, where *impossibility* is getting a particular side when throwing this dice, and *physical impossibility* is getting a side that does not exist. In millions throws, the latter case never would occur, but the former eventually would.

Let's say that there are two cases of resurrection (Res), Jesus Christ and Elvis, and the total number of people who have ever lived is 108.5 billion,³ then the prior base rate is $2/(108.5 \times 10^9)$. On the other hand, let's say that the sample information comes from a very reliable witness whose true positive rate is 0.9999999. Then, what is the probability of this miracle?⁴

Using the Bayes' rule:

```
\begin{split} P(\text{Res}|\text{Witness}) &= \frac{P(\text{Witness}|\text{Res}) \times P(\text{Res})}{P(\text{Witness})} \\ &= \frac{2/(108.5*10^9) \times 0.9999999}{2/(108.5*10^9) \times 0.9999999 + (1-2/(108.5*10^9)) \times (1-0.9999999)} \\ &= 0.000184297806959661 \end{split}
```

where $P(\text{Witness}) = P(\text{Witness}|\text{Res}) \times P(\text{Res}) + (1 - P(\text{Witness}|\text{Res})) \times (1 - P(\text{Res}))$.

Thus, 1.843×10^{-4} is the probability of a resurrection given a very reliable witness.

Observe that we can condition on many events in the Bayes' rule. Let's have two conditioning events B and C, then equation 1.1 becomes

 $^{^3} https://www.wolframalpha.com/input/?i=number+of+people+who+have+ever+lived+on+Earth\ ^4 https://www.r-bloggers.com/2019/04/base-rate-fallacy-or-why-no-one-is-justified-to-believe-that-jesus-rose/$



FIGURE 1.1
The Monty Hall problem.

$$P(A_i|B,C) = \frac{P(A_i,B,C)}{P(B,C)} = \frac{P(B|A_i,C) \times P(A_i|C) \times P(C)}{P(B|C)P(C)}.$$
 (1.2)

Let's use this rule in one of the most intriguing statistical puzzles, the Monty Hall problem, to illustrate how to use equation 1.2 [94, 95]. This was the situation faced by a contestant in the American television game show Let's Make a Deal. There, the contestant was asked to choose a door where behind one door there is a car, and behind the others, goats. Let's say that the contestant picks door No. 1, and the host (Monty Hall), who knows what is behind each door, opens door No. 3, where there is a goat (see Figure 1.1). Then, the host asks the tricky question to the contestant, do you want to pick door No. 2?

Let's name P_i the event **contestant picks door No.** i, which stays close, H_i the event **host picks door No.** i, which is open, and there is a goat, and 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 $P(C_i) = 1/3$, i = 1, 2, 3, and now $P(C_i|H_3) = 1/2$, i = 1, 2 as the host opened door No. 3. So, why bothering 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 picks a door where there is a goat given contestant choice. In this particular

The Bayes' rule 7

setting, $P(H_3|C_3, P_1) = 0$, $P(H_3|C_2, P_1) = 1$ and $P(H_3|C_1, P_1) = 1/2$. Then, using equation 1.2

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

where the third equation uses the fact that C_i and P_i are independent events, and $P(H_3|P_1) = 1/2$ due to this depending just on P_1 (not on C_2).

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

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

R code. The Monty Hall problem

```
set.seed(0101) # Set simulation seed
2 S <- 100000 # Simulations
  Game <- function(switch = 0){</pre>
     # switch = 0 is not change
     # switch = 1 is to change
    opts <- 1:3
     car <- sample(opts, 1) # car location</pre>
     guess1 <- sample(opts, 1) # Initial guess</pre>
     if(car != guess1) {
     host <- opts[-c(car, guess1)]
    } else {
12
     host <- sample(opts[-c(car, guess1)], 1)
    win1 <- guess1 == car # Win no change
15
     guess2 <- opts[-c(host, guess1)]</pre>
     win2 <- guess2 == car # Win change
     if(switch == 0){
      win <- win1
    } else {
      win <- win2
21
22
23
     return(win)
24 }
26 #Win probabilities not changing
        <- mean(replicate(S, Game(switch = 0)))</pre>
28 Prob
29 0.3334
31 #Win probabilities changing
32 Prob <- mean(replicate(S, Game(switch = 1)))</pre>
33 Prob
34 0.6654
```

1.2 Bayesian framework: A brief summary of theory

For two random objects θ and y, the Bayes' rule may be analogously used,⁵

⁵From a Bayesian perspective θ is fixed, but unknown. Then, it is treated as a random object despite the lack of variability (see Chapter 2).

$$\pi(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})}{p(\mathbf{y})},$$
(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 (statistical model), and

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

is the marginal likelihood or prior predictive. Observe that for this expected value to be meaningful the prior should be a proper density, that is, integrates 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, $\int_{\Theta} \pi(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta} = \int_{\Theta} \frac{p(\mathbf{y}|\boldsymbol{\theta}) \times c}{\int_{\Theta} p(\mathbf{y}|\boldsymbol{\theta}) \times cd\boldsymbol{\theta}} d\boldsymbol{\theta} = 1$ 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, theory foundations, expert knowledge and/or mathematical convenience (see Chapters 3 and 4). 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 $\boldsymbol{\Theta}$, $p(\mathbf{y})$ 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} ? 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, let's say \mathcal{M}_m , $m = \{1, 2, ..., M\}$. Then, we should make explicit this in the Bayes' rule formulation,

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

The posterior model probability is

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

where $p(\mathbf{y}|\mathcal{M}_m) = \int_{\mathbf{\Theta}} p(\mathbf{y}|\boldsymbol{\theta}, \mathcal{M}_m) \times \pi(\boldsymbol{\theta}|\mathcal{M}_m) d\boldsymbol{\theta}$ 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 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 what 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 ,

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

Equation 1.7 is a very good shortcut to perform Bayesian inference about θ .

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

$$PO_{12} = \frac{\pi(\mathcal{M}_1|\mathbf{y})}{\pi(\mathcal{M}_2|\mathbf{y})}$$
$$= \frac{p(\mathbf{y}|\mathcal{M}_1)}{p(\mathbf{y}|\mathcal{M}_2)} \times \frac{\pi(\mathcal{M}_1)}{\pi(\mathcal{M}_2)}, \tag{1.8}$$

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 $p(\mathbf{y}|\mathcal{M}_m) = \int_{\Theta} p(\mathbf{y}|\boldsymbol{\theta}, \mathcal{M}_m) \pi(\boldsymbol{\theta}|\mathcal{M}_m) d\boldsymbol{\theta} = \mathbb{E}\left[p(\mathbf{y}|\boldsymbol{\theta}, \mathcal{M}_m)\right]$. For this integral to be meaningful, the prior must be proper. Using improper prior has unintended consequences when comparing models, for instance, parsimonious models are favored by posterior odds or Bayes factors depend on units of measure (see Chapter 4).

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 $\pi(\mathcal{M}_1|\mathbf{y}) + \pi(\mathcal{M}_2|\mathbf{y}) = 1$. Then, $\pi(\mathcal{M}_1|\mathbf{y}) = \frac{PO_{12}}{1+PO_{12}}$ and $\pi(\mathcal{M}_2|\mathbf{y}) = 1 - \pi(\mathcal{M}_1|\mathbf{y})$. In general, $\pi(\mathcal{M}_m|\mathbf{y}) = \frac{p(\mathbf{y}|\mathcal{M}_m) \times \pi(\mathcal{M}_m)}{\sum_{l=1}^{M} p(\mathbf{y}|\mathcal{M}_l) \times \pi(\mathcal{M}_l)}$. These posterior model probabilities can be used to perform Bayesian model averaging.

Table 1.1 shows guidelines for the interpretation of $2 \log(PO_{12})$ [55]. 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 in its own or use a non-informative prior. In this case, we can use the posterior predictive p-value [35, 36].⁶

 $^{^{6}}$ [4] 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.

TABLE 1.1Kass 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

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

This is done simulating pseudo-data from our sampling model $(\mathbf{y}^{(s)}, s = 1, 2, ..., S)$ 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 $(p(D_{\mathbf{y}}) \leq 0.05 \text{ or } p(D_{\mathbf{y}}) \geq 0.95)$ suggest potential discrepancy between the data and the model. [36] also suggest 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 $D(\mathbf{y}) = \mathbb{E}[D(\mathbf{y}, \boldsymbol{\theta})] = \int_{\boldsymbol{\Theta}} D(\mathbf{y}, \boldsymbol{\theta}) \pi(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta}$. These alternatives can be more computational demanding.

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

$$\pi(\mathbf{Y}_0|\mathbf{y}, \mathcal{M}_m) = \int_{\mathbf{\Theta}} \pi(\mathbf{Y}_0, \boldsymbol{\theta}|\mathbf{y}, \mathcal{M}_m) d\boldsymbol{\theta}$$
$$= \int_{\mathbf{\Theta}} \pi(\mathbf{Y}_0|\boldsymbol{\theta}, \mathbf{y}, \mathcal{M}_m) \pi(\boldsymbol{\theta}|\mathbf{y}, \mathcal{M}_m) d\boldsymbol{\theta}. \tag{1.9}$$

Observe that equation 1.9 is again an expectation $\mathbb{E}[\pi(\mathbf{Y}_0|\boldsymbol{\theta},\mathbf{y},\mathcal{M}_m)]$, this time using the posterior distribution. 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 5).

Bayesian model average (BMA) allows considering model uncertainty in prediction or any unknown probabilistic object. In the prediction case,

$$\pi(\mathbf{Y}_0|\mathbf{y}) = \sum_{m=1}^{M} \pi(\mathcal{M}_m|\mathbf{y})\pi(\mathbf{Y}_0|\mathbf{y}, \mathcal{M}_m), \qquad (1.10)$$

and parameters case,

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

where

$$\mathbb{E}(\boldsymbol{\theta}|\mathbf{y}) = \sum_{m=1}^{M} \hat{\boldsymbol{\theta}}_{m} \pi(\mathcal{M}_{m}|\mathbf{y}), \tag{1.12}$$

and

$$Var(\boldsymbol{\theta}|\mathbf{y}) = \sum_{m=1}^{M} \pi(\mathcal{M}_m|\mathbf{y}) \widehat{Var}(\boldsymbol{\theta}|\mathbf{y}, \mathcal{M}_m) + \sum_{m=1}^{M} \pi(M_m|\mathbf{y}) (\hat{\boldsymbol{\theta}}_m - \mathbb{E}[\boldsymbol{\theta}|\mathbf{y}])^2,$$
(1.13)

 $\hat{\boldsymbol{\theta}}_m$ and $\widehat{Var}(\boldsymbol{\theta}|\mathbf{y}, \mathcal{M}_m)$ are the posterior mean and variance under model m, respectively.

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

A nice advantage of the Bayesian approach, which is very useful in *state* space representations (see Chapter 9), is the way that the posterior distribution updates with new sample information. Given $\mathbf{y} = \mathbf{y}_{1:t+1}$ a sequence of observations from 1 to t+1, then

$$\pi(\boldsymbol{\theta}|\mathbf{y}_{1:t+1}) \propto p(\mathbf{y}_{1:t+1}|\boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})$$

$$= p(y_{t+1}|\mathbf{y}_{1:t},\boldsymbol{\theta}) \times p(\mathbf{y}_{1:t}|\boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})$$

$$\propto p(y_{t+1}|\mathbf{y}_{1:t},\boldsymbol{\theta}) \times \pi(\boldsymbol{\theta}|\mathbf{y}_{1:t}). \tag{1.14}$$

We observe that the new prior is just the posterior distribution using the previous observations. This is particular useful under the assumption of conditional independence, that is, $y_{t+1} \perp \mathbf{y}_{1:t}|\boldsymbol{\theta}$, then $p(y_{t+1}|\mathbf{y}_{1:t},\boldsymbol{\theta}) = p(y_{t+1}|\boldsymbol{\theta})$ such that the posterior can be recovered recursively [76]. This facilities online updating due to all information up to t being in $\boldsymbol{\theta}$. Then, $\pi(\boldsymbol{\theta}|\mathbf{y}_{1:t+1}) \propto p(y_{t+1}|\boldsymbol{\theta}) \times \pi(\boldsymbol{\theta}|\mathbf{y}_{1:t}) \propto \prod_{h=1}^{t+1} p(y_h|\boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})$. This recursive expression can be calculated faster at some specific point in time t compared to a batch mode algorithm, which requires processing simultaneously all information up to t.

It is also important to wonder about the sampling properties of "Bayesian estimators". This topic has attracted attention of statisticians and econometricians long time ago. For instance, asymptotic posterior concentration at the population parameter vector is discussed by [13]. Convergence of posterior

distributions is stated by the Bernstein-von Mises theorem [61, 103], which creates 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 α level converge asymptotically to confidence intervals at α level. This suggests that Bayesian inference is asymptotically correct from a sampling perspective in these settings.

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

$$\pi(\theta|\mathbf{y}) \propto p(\mathbf{y}|\theta) \times \pi(\theta)$$

$$= \exp\left\{N \times \bar{l}(\mathbf{y}|\theta)\right\} \times \pi(\theta). \tag{1.15}$$

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

Maximum likelihood theory shows that $\lim_{N\to\infty} \bar{l}(\mathbf{y}|\theta) \to \bar{l}(\mathbf{y}|\theta_0)$ where θ_0 is the population parameter of the data generating process. In addition, doing a second order Taylor expansion of the log likelihood at the Maximum likelihood estimator,

$$\begin{split} l(\mathbf{y}|\theta) &\approx l(\mathbf{y}|\hat{\theta}) + \frac{dl(\mathbf{y}|\theta)}{d\theta} \bigg|_{\hat{\theta}} (\theta - \hat{\theta}) + \frac{1}{2} \left. \frac{d^2 l(\mathbf{y}|\theta)}{d\theta^2} \right|_{\hat{\theta}} (\theta - \hat{\theta})^2 \\ &= l(\mathbf{y}|\hat{\theta}) + \frac{1}{2} \sum_{i=1}^{N} \frac{d^2 l(y_i|\theta)}{d\theta^2} \bigg|_{\hat{\theta}} (\theta - \hat{\theta})^2 \\ &= l(\mathbf{y}|\hat{\theta}) - \frac{1}{2} N \left[-\bar{l}'' \big|_{\hat{\theta}} \right] (\theta - \hat{\theta})^2 \\ &= l(\mathbf{y}|\hat{\theta}) - \frac{N}{2\sigma^2} (\theta - \hat{\theta})^2 \end{split}$$

where
$$\frac{dl(\mathbf{y}|\theta)}{d\theta}\Big|_{\hat{\theta}} = 0$$
, $\bar{l}'' \equiv \frac{1}{N} \sum_{i=1}^{N} \frac{d^2l(y_i|\theta)}{d\theta^2}\Big|_{\hat{\theta}}$ and $\sigma^2 := \left[-\bar{l}''\Big|_{\hat{\theta}}\right]^{-1}$. 8 Then,

⁷Take into account that in the likelihood function the argument is θ . However, we keep the notation for facility in exposition.

⁸The last definition follows from standard theory in maximum likelihood estimation (see [17, Chap. 10] and [106, Chap. 13]).

$$\pi(\theta|\mathbf{y}) \propto \exp\left\{l(\mathbf{y}|\theta)\right\} \times \pi(\theta)$$

$$\approx \exp\left\{l(\mathbf{y}|\hat{\theta}) - \frac{N}{2\sigma^2}(\theta - \hat{\theta})^2\right\} \times \pi(\theta)$$

$$\propto \exp\left\{-\frac{N}{2\sigma^2}(\theta - \hat{\theta})^2\right\} \times \pi(\theta)$$

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

1.2.1 Example: Health insurance

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

Assuming that this is a random sample⁹ from a data generating process (statistical model) that is Poisson, that is, $Y_i \sim P(\lambda)$, and your probabilistic prior beliefs about λ are well described by a Gamma distribution with shape and scale parameters α_0 and β_0 , $\lambda \sim G(\alpha_0, \beta_0)$, then, you are interested in calculating the probability $P(Y_0 > 0|\mathbf{y})$. You need to calculate the posterior predictive density $\pi(Y_0|\mathbf{y})$ to answer this question in a Bayesian way.

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

$$\pi(Y_0|\mathbf{y}) = \int_0^\infty \frac{\lambda^{y_0} \exp\left\{-\lambda\right\}}{y_0!} \times \pi(\lambda|\mathbf{y}) d\lambda,$$

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 $\int_0^\infty \pi(\lambda|\mathbf{y}) d\lambda = 1$, then 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

⁹Independent and identically distributed draws.

$$\mathbb{E}[\lambda|\mathbf{y}] = \alpha_n \beta_n$$

$$= \left(\sum_{i=1}^N y_i + \alpha_0\right) \left(\frac{\beta_0}{\beta_0 N + 1}\right)$$

$$= \bar{y} \left(\frac{N\beta_0}{N\beta_0 + 1}\right) + \alpha_0 \beta_0 \left(\frac{1}{N\beta_0 + 1}\right)$$

$$= w\bar{y} + (1 - w)\mathbb{E}[\lambda],$$

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

The predictive distribution is

$$\begin{split} \pi(Y_0|\mathbf{y}) &= \int_0^\infty \frac{\lambda^{y_0} \exp\left\{-\lambda\right\}}{y_0!} \times \frac{1}{\Gamma(\alpha_n)\beta_n^{\alpha_n}} \lambda^{\alpha_n - 1} \exp\left\{-\lambda/\beta_n\right\} d\lambda \\ &= \frac{1}{y_0! \Gamma(\alpha_n)\beta_n^{\alpha_n}} \int_0^\infty \lambda^{y_0 + \alpha_n - 1} \exp\left\{-\lambda \left(\frac{1 + \beta_n}{\beta_n}\right)\right\} d\lambda \\ &= \frac{\Gamma(y_0 + \alpha_n) \left(\frac{\beta_n}{\beta_n + 1}\right)^{y_0 + \alpha_n}}{y_0! \Gamma(\alpha_n)\beta_n^{\alpha_n}} \\ &= \left(\frac{y_0 + \alpha_n - 1}{y_0}\right) \left(\frac{\beta_n}{\beta_n + 1}\right)^{y_0} \left(\frac{1}{\beta_n + 1}\right)^{\alpha_n}. \end{split}$$

The third equality follows from the kernel of a Gamma density, and the fourth from $\binom{y_0+\alpha_n-1}{y_0}=\frac{(y_0+\alpha_n-1)(y_0+\alpha_n-2)...\alpha_n}{y_0!}=\frac{\Gamma(y_0+\alpha_n)}{\Gamma(\alpha_n)y_0!}$ using a property of the Gamma function.

Observe that this is a Negative Binomial density, that is $Y_0|\mathbf{y} \sim NB(\alpha_n, p_n)$ 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 $\beta_0 = 1/0.001$ which imply vague prior information about λ due to having a large degree of variability compared to the mean information.¹⁰ In particular, $\mathbb{E}[\lambda] = 1$ and $\mathbb{V}ar[\lambda] = 1000$.

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.

¹⁰We should be aware that there may be technical problems using this king of hyperparameters in this setting [38].

Another approach is using *Empirical Bayes*, where we set the hyperparameters maximizing the logarithm of the marginal likelihood, that is, $\left[\hat{\alpha}_0 \ \hat{\beta}_0\right]^{\top} = \underset{\alpha_0,\beta_0}{\operatorname{argmax}} \, \ln p(\mathbf{y})$ where

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

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

$$= \frac{\Gamma(\sum_{i=1}^N y_i + \alpha_0) \left(\frac{\beta_0}{N\beta_0 + 1}\right)^{\sum_{i=1}^N y_i} \left(\frac{1}{N\beta_0 + 1}\right)^{\alpha_0}}{\Gamma(\alpha_0) \prod_{i=1}^N y_i}$$

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

$$PO_{12} = \frac{p(\mathbf{y}|\text{Empirical Bayes})}{p(\mathbf{y}|\text{Vague prior})}$$

$$= \frac{\frac{\Gamma(\sum_{i=1}^{N} y_i + 51.807) \left(\frac{0.023}{N \times 0.023 + 1}\right)^{\sum_{i=1}^{N} y_i} \left(\frac{1}{N \times 0.023 + 1}\right)^{51.807}}{\Gamma(\sum_{i=1}^{N} y_i + 0.001) \left(\frac{1/0.001}{N/0.001 + 1}\right)^{\sum_{i=1}^{N} y_i} \left(\frac{1}{N/0.001 + 1}\right)^{0.001}}{\Gamma(0.001)}$$

$$\approx 919.$$

Then, $2 \times \log(PO_{12}) = 13.64$, there is 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 average (BMA). In particular,

$$\mathbb{E}(\lambda|\mathbf{y}) = 1.2 \times 0.999 + 1.2 \times 0.001 = 1.2$$

$$Var(\lambda|\mathbf{y}) = 0.025 \times 0.999 + 0.24 \times 0.001$$

$$+ (1.2 - 1.2)^2 \times 0.999 + (1.2 - 1.2)^2 \times 0.001 = 0.025.$$

The BMA predictive distribution is a mix of negative binomial distributions, that is, $y_0|\mathbf{y} \sim 0.999 \times NB(57.8, 0.02) + 0.001 \times NB(6.001, 0.17)$.

R code. Health insurance, predictive distribution using vague hyperparameters

```
set.seed(010101)
2 y <- c(0, 3, 2, 1, 0) # Data
3 N <- length(y)
4 ProbBo <- function(y, a0, b0){
    N <- length(y)
    #sample size
    an \leftarrow a0 + sum(y)
     # Posterior shape parameter
    bn <- b0 / ((b0 * N) + 1)
# Posterior scale parameter
    p <- bn / (bn + 1)
     # Probability negative binomial density
     Pr <- 1 - pnbinom(0, size=an, prob=(1 - p))</pre>
     # Probability of visiting the Doctor at least once next
     # Observe that in R there is a slightly different
       parametrization.
     return(Pr)
16
17 }
18 # Using a vague prior:
19 a0 <- 0.001 # Prior shape parameter
_{20} b0 <- 1 / 0.001 # Prior scale parameter
PriMeanV <- a0 * b0 # Prior mean
PriVarV <- a0 * b0^2 # Prior variance
^{23} Pp <- ProbBo(y, a0 = 0.001, b0 = 1 / 0.001)
24 # This setting is vague prior information.
25 Pp
26 0.67
```

R code. Health insurance, predictive distribution using empirical Bayes

```
1 # Using Empirical Bayes
2 LogMgLik <- function(theta, y){</pre>
3 N <- length(y)
#sample size
a0 <- theta[1]</pre>
6 # prior shape hyperparameter
7 b0 <- theta[2]
_{8} # prior scale hyperparameter
   an <- sum(y) + a0
   # posterior shape parameter
   if(a0 <= 0 || b0 <= 0){
    #Avoiding negative values
    lnp <- -Inf
13
14
    }else{
    lnp <- lgamma(an) + sum(y)*log(b0/(N*b0+1)) - a0*log(N*b0
       +1) - lgamma(a0)
16 }
17 # log marginal likelihood
   return(-lnp)
19 }
20 theta0 \leftarrow c(0.01, 1/0.1)
_{21} # Initial values
22 control <- list(maxit = 1000)</pre>
23 # Number of iterations in optimization
24 EmpBay <- optim(theta0, LogMgLik, method = "BFGS", control =
       control , hessian = TRUE , y = y)
25 # Optimization
26 EmpBay $convergence
28 aOEB <- EmpBay par [1]
29 # Prior shape using empirical Bayes
30 a0EB
31 51.81
32 bOEB <- EmpBay par [2]
33 # Prior scale using empirical Bayes
34 b0EB
35 0.023
36 PriMeanEB <- a0EB * b0EB
37 # Prior mean
38 PriVarEB <- a0EB * b0EB^2
39 # Prior variance
40 PpEB \leftarrow ProbBo(y, a0 = a0EB, b0 = b0EB)
41 # This setting is using emprical Bayes.
42 PpEB
43 0.70
```

R code. Health insurance, density plots

```
1 # Density figures:
2 # This code helps plotting densities
_3 lambda <- seq(0.01, 10, 0.01)
4 # Values of lambda
5 VaguePrior <- dgamma(lambda,shape=a0,scale = b0)</pre>
6 EBPrior <- dgamma(lambda, shape=a0EB, scale = b0EB)
7 PosteriorV <- dgamma(lambda, shape = a0 + sum(y), scale = b0</pre>
       / ((b0 * N) + 1))
8 PosteriorEB <- dgamma(lambda, shape = a0EB+sum(y), scale =</pre>
      bOEB / ((bOEB * N) + 1))
9 # Likelihood function
10 Likelihood <- function(theta, y){</pre>
11 LogL <- dpois(y, theta, log = TRUE)</pre>
12 Lik <- prod(exp(LogL))</pre>
13 return(Lik)
14 }
15 Liks <- sapply(lambda, function(par) {Likelihood(par, y = y)
      })
16 Sc <- max(PosteriorEB)/max(Liks)</pre>
17 #Scale for displaying in figure
18 LiksScale <- Liks * Sc
19 data <- data.frame(cbind(lambda, VaguePrior, EBPrior,</pre>
      PosteriorV, PosteriorEB, LiksScale)) #Data frame
20 require(ggplot2) # Cool figures
21 require(latex2exp) # LaTeX equations in figures
require(ggpubr) # Multiple figures in one page
23 fig1 <- ggplot(data = data, aes(lambda, VaguePrior)) + geom_
      line() + xlab(TeX("$\\lambda$")) + ylab("Density") +
      ggtitle("Prior: Vague Gamma")
24 fig2 <- ggplot(data = data, aes(lambda, EBPrior)) + geom_
      line() + xlab(TeX("$\\lambda$")) + ylab("Density") +
      ggtitle("Prior: Empirical Bayes Gamma")
25 fig3 <- ggplot(data = data, aes(lambda, PosteriorV)) + geom_
      line() + xlab(TeX("$\\lambda$")) + ylab("Density") +
      ggtitle("Posterior: Vague Gamma")
26 fig4 <- ggplot(data = data, aes(lambda, PosteriorEB)) + geom
       _{	t line()} + {	t xlab(TeX("$ \ \ )} + {	t ylab("Density")} +
      ggtitle("Posterior: Empirical Bayes Gamma")
27 FIG <- ggarrange(fig1, fig2, fig3, fig4, ncol = 2, nrow = 2)
28 annotate_figure(FIG, top = text_grob("Vague versus Empirical
       Bayes: Poisson-Gamma model", color = "black", face = "
      bold", size = 14))
29 dataNew <- data.frame(cbind(rep(lambda, 3), c(EBPrior,
      PosteriorEB, LiksScale), rep(1:3, each = 1000))) #Data
30 colnames(dataNew) <- c("Lambda", "Density", "Factor")
31 dataNew$Factor <- factor(dataNew$Factor, levels=c("1", "3",</pre>
      "2"),
32 labels=c("Prior", "Likelihood", "Posterior"))
33 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 prior and posterior densities based on vague and Empirical Bayes hyperparameters. We see that prior and posterior densities using the latter are more informative as expected.

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

R code. Health insurance, Predictive density

```
# Predictive distributions
  PredDen <- function(y, y0, a0, b0){</pre>
    N <- length(y)
    #sample size
    an \leftarrow a0 + sum(y)
    # Posterior shape parameter
    bn \leftarrow b0 / ((b0 * N) + 1)
    # Posterior scale parameter
    p <- bn / (bn + 1)
    # Probability negative binomial density
    Pr <- dnbinom(y0, size=an, prob=(1 - p))
    # Predictive density
12
    # Observe that in R there is a slightly different
      parametrization.
     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=a0EB, b0=b0EB)
19 dataPred <- as.data.frame(cbind(y0, PredVague, PredEB))</pre>
20 colnames(dataPred) <- c("y0", "PredictiveVague",</pre>
       PredictiveEB")
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")) +
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 the next year, having one, two, and so on using Empirical Bayes and vague hyperparameters. The predictive probability of not having any visits are approximately equal to 30% and 33% based on the Empirical Bayes and vague hyperparameters.



FIGURE 1.4
Predictive density: Vague and Empirical Bayes.

R code. Health insurance, Bayesian model average 1 # Posterior odds: Vague vs Empirical Bayes 2 PO12 <- exp(-LogMgLik(c(aOEB, bOEB), y = y))/exp(-LogMgLik(c (a0, b0), y = y))3 PO12 4 919 5 PostProMEM <- P012/(1 + P012)</pre> 6 PostProMEM 8 # Posterior model probability Empirical Bayes 9 PostProbMV <- 1 - PostProMEM 10 PostProbMV 11 0.002 12 # Posterior model probability vague prior 13 # Bayesian model average (BMA) 14 PostMeanEB \leftarrow (a0EB + sum(y)) * (b0EB / (b0EB * N + 1)) 15 # Posterior mean Empirical Bayes 16 PostMeanV \leftarrow (a0 + sum(y)) * (b0 / (b0 * N + 1)) 17 # Posterior mean vague priors 18 BMAmean <- PostProMEM * PostMeanEB + PostProbMV * PostMeanV 19 BMAmean 20 1.2 21 # BMA posterior mean 22 PostVarEB <- (a0EB + sum(y)) * (b0EB/(b0EB * N + 1))^2 23 # Posterior variance Empirical Bayes 24 PostVarV \leftarrow (a0 + sum(y)) * (b0 / (b0 * N + 1))^2 25 # Posterior variance vague prior 26 BMAVar <- PostProMEM * PostVarEB + PostProbMV*PostVarV + PostProMEM * (PostMeanEB - BMAmean)^2 + PostProbMV * (PostMeanV - BMAmean)^2 27 # BMA posterior variance 28 BMAVar 29 0.025

$ig| R \ code. \ Health \ insurance, \ Bayesian \ model \ average$

Predictive density: BMA



FIGURE 1.5 Bayesian model average: Predictive density.

Figure 1.5 displays the predictive density using Bayesian model average based on the vague and Empirical Bayes hyperparameters. This figure essentially 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 displays how the posterior distribution updates given new sample information based on an initial non-informative prior (iteration 1). We



FIGURE 1.6 Bayesian updating: Posterior densities.

see that iteration 5 is based on all the sample information in our example, as a consequence, the posterior density in iteration 5 is equal to the posterior density in Figure 1.3.

R code. Health insurance, Bayes updating

```
1 # Bayesian updating
2 BayUp <- function(y, lambda, a0, b0){</pre>
    N <- length(y)
    #sample size
    an \leftarrow a0 + sum(y)
    # Posterior shape parameter
    bn <- b0 / ((b0 * N) + 1)
     # Posterior scale parameter
    p <- dgamma(lambda, shape = an, scale = bn)
    # Posterior density
    return(list(Post = p, a0New = an, b0New = bn))
12 }
13 PostUp <- NULL
14 for(i in 1:N){
   if(i == 1){
15
       PostUpi \leftarrow BayUp(y[i], lambda, a0 = 0.001, b0 = 1/0.001)
     else{
      PostUpi <- BayUp(y[i], lambda, a0 = PostUpi$a0New, b0 =
18
       PostUpi$b0New)
19
    PostUp <- cbind(PostUp, PostUpi$Post)</pre>
20
21 }
22 DataUp <- data.frame(cbind(rep(lambda, 5), c(PostUp), rep</pre>
       (1:5, each = 1000))) #Data frame
23 colnames(DataUp) <- c("Lambda", "Density", "Factor")</pre>
24 DataUp$Factor <- factor(DataUp$Factor, levels=c("1", "2", "3
       ", "4", "5"),
25 labels=c("Iter 1", "Iter 2", "Iter 3", "Iter 4", "Iter 5"))
26 ggplot(data = DataUp, aes_string(x = "Lambda", y = "Density"
       , group = "Factor")) + geom_line(aes(color = Factor)) +
       xlab(TeX("$\\lambda$")) + ylab("Density") + ggtitle("
       Bayesian updating: Poisson-Gamma model with vague prior"
       ) + guides(color=guide_legend(title="Update")) + scale_
       color_manual(values = c("red", "purple", "blue", "yellow
       ", "black"))
27 S <- 100000 # Posterior draws
28 PostMeanLambdaUps <- sapply(1:N, function(i) {mean(sample())</pre>
       lambda, S, replace = TRUE, prob = PostUp[ , i])))) #
       Posterior mean update i
29 paste("Posterior means using all information and sequential
       updating are:", round(PostMeanV, 2), "and", round(
       PostMeanLambdaUps[5], 2), sep = " ")
30 Posterior means using all information and sequential
       updating are: 1.2 and 1.2
31 PostVarLambdaUps <- sapply(1:N, function(i) {var(sample(</pre>
       lambda, S, replace = TRUE, prob = PostUp[ , i])))) #
       Posterior variance update i
32 paste("Posterior variances using all information and
       sequential updating are:", round(PostVarV, 2), "and",
round(PostVarLambdaUps[5], 2), sep = " ")
33 Posterior variances using all information and sequential
       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 demand to report 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 a decision regarding what are the optimal posterior values to report [11].

The point of departure is a loss function, which is a non-negative real value function whose arguments are the unknown state of nature (Θ) , and a set of actions to be made (A), that is,

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

This function is a mathematical expression of the loss of making mistakes. In particular, selecting action $a \in \mathcal{A}$ when $\theta \in \Theta$ is the true. In our case, the unknown state of nature can be 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 $(\mathbb{E}[L(\boldsymbol{\theta}, a)|\mathbf{y}])$,

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

where $\mathbb{E}[L(\boldsymbol{\theta}, a)|\mathbf{y}] = \int_{\boldsymbol{\Theta}} L(\boldsymbol{\theta}, a) \pi(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta}$. 11

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

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

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

To get this results, let's use the first condition order, differentiate the risk function with respect to a, interchange differential and integral order, and set this equal to zero, $-2\int_{\Theta}[\theta-a^*]\pi(\theta|\mathbf{y})d\theta=0$ 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 θ when facing the quadratic loss function.

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

- The generalized quadratic loss function, $L(\theta, a) = w(\theta)[\theta a]^2$, where $w(\theta) > 0$ is a weighting function, gives as optimal decision rule the weighted mean. We should follow same steps as the previous result to get $a^*(\mathbf{y}) = \frac{\mathbb{E}[w(\theta) \times \theta|\mathbf{y}]}{\mathbb{E}[w(\theta)|\mathbf{y}]}$. 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 optimal action the posterior median (Exercise 5).
- The generalized absolute error function,

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

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,

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

then, $\int_{-\infty}^{a^*} \pi(\theta|\mathbf{y}) d\theta = \frac{K_0}{K_0 + K_1}$, that is, any $K_0/(K_0 + K_1)$ -percentile of $\pi(\theta|\mathbf{y})$ is an optimal Bayesian estimate of θ .

We can also use decision theory under uncertainty in hypothesis testing. In particular, testing $H_0: \theta \in \Theta_0$ versus $H_1: \theta \in \Theta_1$, $\Theta = \Theta_0 \cup \Theta_1$ and $\emptyset = \Theta_0 \cap \Theta_1$, there are two actions of interest, a_0 and a_1 , where a_j denotes no rejecting H_j , $j = \{0, 1\}$.

Given the $0 - K_i$ loss function,

$$L(\theta, a_j) = \left\{ \begin{matrix} 0, & \theta \in \Theta_j \\ K_j, & \theta \in \Theta_j, j \neq i \end{matrix} \right\},$$

where there is no loss if the right decision is made, for instance, no rejecting H_0 when $\theta \in \Theta_0$, and the loss is K_j when an error is made, for instance, type I error, rejecting the null hypothesis (H_0) when it is true $(\theta \in \Theta_0)$, implies a loss equal to K_1 due to picking a_1 , no rejecting H_1 .

The posterior expected loss associated with decision a_j , that is, no rejecting H_j , is $\mathbb{E}[L(\theta, a_j)|\mathbf{y}] = 0 \times P(\Theta_j|\mathbf{y}) + K_jP(\Theta_i|\mathbf{y}) = K_jP(\Theta_i|\mathbf{y}), j \neq i$. Therefore, the Bayes optimal decision is the one that gives the smallest posterior expected loss, that is, the null hypothesis is rejected $(a_1 \text{ is not rejected})$, when

 $K_0P(\Theta_1|\mathbf{y}) > K_1P(\Theta_0|\mathbf{y})$. Given our framework $(\Theta = \Theta_0 \cup \Theta_1, \emptyset = \Theta_0 \cap \Theta_1)$, then $P(\Theta_0|\mathbf{y}) = 1 - P(\Theta_1|\mathbf{y})$, and as a consequence, $P(\Theta_1|\mathbf{y}) > \frac{K_1}{K_1 + K_0}$, that is, 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 $\Theta_{C(\mathbf{y})} \subset \Theta$ a credible set for θ , and $L(\theta, \Theta_{C(\mathbf{y})}) = 1 - \mathbb{1} \{ \theta \in \Theta_{C(\mathbf{y})} \}$, where

$$\mathbb{1}\left\{\theta \in \Theta_{C(\mathbf{y})}\right\} = \left\{\begin{matrix} 1, & \theta \in \Theta_{C(\mathbf{y})} \\ 0, & \theta \notin \Theta_{C(\mathbf{y})} \end{matrix}\right\}.$$

Then,

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

where the 0–1 loss function is equal to zero if $\theta \in \Theta_{C(\mathbf{y})}$, and one if $\theta \notin \Theta_{C(\mathbf{y})}$. Then, 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 $\theta \in \Theta_{C(\mathbf{y})}$; then, we can measure the accuracy of the report by $L(\theta, \alpha(\mathbf{y})) = [\mathbbm{1}\{\theta \in \Theta_{C(\mathbf{y})}\} - \alpha(\mathbf{y})]^2$. 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 $\mathbb{E}[\mathbbm{1}\{\theta \in \Theta_{C(\mathbf{y})}\}|\mathbf{y}] = P(\theta \in \Theta_{C(\mathbf{y})}|\mathbf{y})$. This probability can be calculated given the posterior distribution, that is, $P(\theta \in \Theta_{C(\mathbf{y})}|\mathbf{y}) = \int_{\Theta_{C(\mathbf{y})}} \pi(\theta|\mathbf{y})d\theta$. This is a measure of the belief that $\theta \in \Theta_{C(\mathbf{y})}$ given the prior beliefs and sample information.

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

Two alternatives to report credible sets are the symmetric credible set and the highest posterior density set (HPD). The former is based on $\frac{\alpha}{2}\%$ and $(1-\frac{\alpha}{2})\%$ percentiles of the posterior distribution, and the latter is a $100(1-\alpha)\%$ credible interval for θ with the property that it has the smallest distance compared to any other $100(1-\alpha)\%$ credible interval for θ based on the posterior distribution. That is, $C(\mathbf{y}) = \{\theta : \pi(\theta|\mathbf{y}) \ge k(\alpha)\}$, where $k(\alpha)$ is the largest number such that $\int_{\theta:\pi(\theta|\mathbf{y})\ge k(\alpha)} \pi(\theta|\mathbf{y})d\theta = 1-\alpha$. The HPD set can be a collection of disjoint intervals when working with multimodal posterior densities. In addition, they have the limitation of not necessary 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, $\mathbb{E}_{Y_0}[L(Y_0, a)] = \int_{\mathcal{Y}_0} L(Y_0, a)\pi(Y_0|\mathbf{y})dY_0$, 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 BMA allows incorporating model uncertainty in a regression framework, sometimes it is desirable to select just one model. A compelling alternative is 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 [23].

1.3.1 Example: Health insurance continues

We show some optimal rules in the health insurance example. In particular, the best point estimates of λ given the quadratic, absolute and generalized absolute loss functions. For the third, we assume that underestimating λ is twice as costly as overestimating it, that is, $K_0 = 2$ and $K_1 = 1$.

Taking into account that the posterior distribution of λ is $G(\alpha_0 + \sum_{i=1}^{N} y_i, \beta_0/(\beta_0 N + 1))$, using the hyperparameters from empirical Bayes, we have that $\mathbb{E}[\lambda|\mathbf{y}] = \alpha_n \beta_n = 1.2$, the median is 1.19, and the 2/3-th quantile is 1.26. Those are the optimal point estimates for the quadratic, absolute and generalized absolute loss functions.

In addition, we test the null hypothesis H_0 . $\lambda \in [0,1)$ versus H_1 . $\lambda \in [1,\infty)$ setting $K_0 = K_1 = 1$ we should reject the null hypothesis due to $P(\lambda \in [0,1)) = 0.9 > K_1/(K_0 + K_1) = 0.5$.

We get that the 95% symmetric credible interval is (0.91, 1.53), and the highest posterior density interval is (0.9, 1.51). Finally, the optimal point prediction under a quadratic loss function is 1.2, which is the mean value of the posterior predictive distribution, and the optimal model assuming a 0-1 loss function is the model using the hyperparameters from the empirical Bayes procedure due to the posterior model probability of this model being approximately 1, whereas the posterior model probability of the model using vague hyperparameters is approximately 0.

R code. Health insurance, Bayesian reports

```
1 an \leftarrow sum(y) + a0EB
2 # Posterior shape parameter
3 bn <- b0EB / (N*b0EB + 1)
4 # Posterior scale parameter
5 S <- 1000000
6 # Number of posterior draws
7 Draws <- rgamma (1000000, shape = an, scale = bn)
8 # Posterior draws
9 ##### Point estimation #######
10 OptQua <- an*bn
11 # Mean: Optimal choice quadratic loss function
12 OptQua
13 1.200952
14 OptAbs <- qgamma(0.5, shape = an, scale = bn)
# Median: Optimal choice absolute loss function
16 OptAbs
17 1.194034
_{\rm 18} # Setting KO = 2 and K1 = 1, that is, to underestimate
      lambda is twice as costly as to overestimate it.
19 KO <- 2; K1 <- 1
20 OptGenAbs <- quantile(Draws, KO/(KO + K1))</pre>
21 # Median: Optimal choice generalized absolute loss function
22 OptGenAbs
23 66.66667%
24 1.262986
25 ##### Hypothesis test #######
26 # HO: lambda in [0,1) vs H1: lambda in [1, Inf]
27 KO <- 1; K1 <- 1
28 ProbHO <- pgamma(1, shape = an, scale = bn)
29 ProbHO # Posterior probability HO
30 0.09569011
31 ProbH1 <- 1 -ProbH0
32 ProbH1 # Posterior probability H1
33 0.9043099
34 # We should reject HO given ProbH1 > K1 / (KO + K1)
35 ##### Credible intervals #######
36 LimInf <- qgamma(0.025, shape = an, scale = bn) # Lower
      bound
37 LimInf
38 0.9114851
39 LimSup <- qgamma(0.975, shape = an, scale = bn) # Upper
      bound
40 LimSup
41 1.529724
42 HDI <- HDInterval::hdi(Draws, credMass = 0.95) # Highest
      posterior density credible interval
43 HDI
      lower
44
45 0.8971505 1.5125911
46 attr(,"credMass")
47 [1] 0.95
48 ##### Predictive optimal choices #######
_{
m 49} p <- bn / (bn + 1) # Probability negative binomial density
50 OptPred <- p/(1-p)*an # Optimal point prediction given a
      quadratic loss function in prediction
51 OptPred
52 1.200952
```

1.4 Summary

We introduce the Bayes' rule to update probabilistic statements using funny examples. Then, we study the three probabilistic objects of main relevance in Bayesian inference: the posterior distribution, the marginal likelihood and the predictive density. The first allows performing inference regarding parameters, the second is required to perform hypothesis test for model selection using the Bayes factor, and the third to perform probabilistic predictions. We also review some sampling properties of Bayesian estimators, and Bayes update. All those concepts were developed using a simple example in R software. Finally, we introduce some concepts of decision theory that can be used to report summary statistics minimizing posterior expected losses.

1.5 Exercises

1. The court case: the blue or green cap

A cab was involved in a hit and run accident at night. There are two cab companies in the town: blue and green. The former has 150 cabs, and the latter 850 cabs. A witness said that a blue cab was involved in the accident; the court tested his/her reliability under the same circumstances, and got that 80% of the times the witness correctly identified the color of the cab. What is the probability that the color of the cab involved in the accident was 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* switching the decision if there are four doors, where there are three goats and one car? Solve this problem 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 that you are analyzing to buy a car insurance next year. To make a better decision you want to know what is the probability that you have a car claim next year? You have the records of your car claims in 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, $Y_i \sim Ber(p)$, and your probabilistic prior beliefs about p are well described by a beta distribution with parameters α_0 and β_0 , $p \sim B(\alpha_0, \beta_0)$, then, you are

Exercises 33

interested in calculating the probability of a claim the next year $P(Y_0 = 1|\mathbf{y})$.

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

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

Conceptual differences of the Bayesian and Frequentist approaches

We give some of the conceptual differences between the Bayesian and Frequentist inferential approaches. We emphasize in the Bayesian concepts as most of the readers can be familiarized with the Frequentist statistical framework.

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:
 - watched this particular contestant participating in this show many times.
 - 2. seen him asked this same question each time,
 - 3. and computed the relative frequency with which he gives the correct answer,

you need to answer this question as a Bayesian!

Uncertainty about the event answer this question needs to be expressed as a "degree of belief" informed both by information coming from data on the skill of the particular participant, and how much he knows about inventors, and possibly prior knowledge on his performance in other game shows. Of course, your prior knowledge of the contestant may be minimal, or it may be very informed. Either way, your final answer remains a degree of belief held 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 simply 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 of nature occurs. To the Frequentist, probability is linked intrinsically 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 as a physical phenomenon, like mass or wavelength, whereas Bayesians stipulate that probability lives in the mind of scientists as any scientific construct [74].

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?", and 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 say that psychological experiments suggest that human beings suffer from *anchoring*, that is, a cognitive bias that causes us to rely too heavily on the previous information (prior) such that the updating process (posterior) due to new information (likelihood) being low compared to the Bayes' rule [52].

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 [86, 26, 91, 63], but others who adopt *objective* prior distributions such as Jeffreys', reference, empirical or robust [6, 59, 50, 10] to operationalize Bayes' rule, and thereby weight quantitative (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.¹ In addition, both Frequentist and Bayesian 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." [14].

We also know that "All models are wrong, but some are useful" [15], that 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 [35, 4] or posterior odds ratios for relative statements [49, 55]. This is because the marginal likelihood, conditional on data, is interpreted as the evidence of the prior distribution [9].

In addition, what is objectivity 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 [63]. In fact "Student" (William Gosset) saw statistical significance at any level as being "nearly valueless" in itself [109]. But, this is not just a situation in the Frequentist approach. The cut-offs given to "establish" scientific evidence against a null hypothesis in terms of log_{10} scale [50] or log_e scale [55] 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 what is required to perform estimation, hypothesis testing (model selection) and prediction in the Bayesian approach is to apply the Bayes' rule. This means coherence under a probabilistic view. But, there is no free lunch, coherence reduces flexibility. On the other hand, the Frequestist approach may be not coherent from a probabilistic point of view, but it is very flexible. This approach can be seen as a tool kit that offers inferential solutions under the umbrella of understanding probability as relative frequency. For instance, a point estimator in a Frequentist approach is found such that satisfies good

¹Type I error is rejecting the null hypothesis when this is true, and the Type II error is not rejecting the null hypothesis when this is false.

sampling properties like unbiasness, efficiency, or a large sample property as consistency.

A remarkable difference is that optimal Bayesian decisions are calculated minimizing the expected value of the loss function with respect to the posterior distribution, that is, it is conditional on observed data. On the other hand, Frequentist "optimal" actions are base on the expected values over the distribution of the estimator (a function of data) conditional on the unknown parameters, that is, it considers sampling variability.

The Bayesian approach allows to obtain the posterior distribution of any unknown object such as parameters, latent variables, future or unobserved variables or models. A nice advantage is that prediction can take into account estimation error, and predictive distributions (probabilistic forecasts) can be easily recovered.

Hypothesis testing (model selection) is based on *inductive logic* reasoning (*inverse probability*); on the basis of what we see, we evaluate what hypothesis is most tenable, and is performed using posterior odds, which in turn are based on Bayes factors that evaluate evidence in favor of a null hypothesis taking explicitly the alternative [55], following the rules of probability [63], comparing how well the hypothesis predicts data [44], minimizing the weighted sum of type I and type II error probabilities [27, 75], and taking the implicit balance of losses [50, 12] into account. Posterior odds allows to use the same framework to analyze nested and non-nested models and perform model average. However, Bayes factors cannot be based on improper or vague priors [56], the practical interplay between model selection and posterior distributions is not as easy as it maybe in the Frequentist approach, and the computational burden can be more demanding due to solving potentially difficult integrals.

On the other hand, the Frequentist approach establishes most of its estimators as the solution of a system of equations. Observe that optimization problems reduce to solve systems. We can potentially get the distribution of these estimators, but most of the time it is needed asymptotic arguments or resampling techniques. Hypothesis testing requires pivotal quantities and/or also resampling, and prediction most of the time is based on a plug-in approach, which means not taking estimation error into account. In addition, ancillary statistics can be used to build prediction intervals. 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 instance in the normal population mean hypothesis framework you cannot reject at α significance level (Type I error) any null hypothesis H_0 . $\mu = \mu^0$ if μ^0 is in the $1 - \alpha$ confidence interval $P(\mu \in [\hat{\mu} - |t_{N-1}^{\alpha/2}| \times \hat{\sigma}_{\hat{\mu}}, \hat{\mu} + |t_{N-1}|^{\alpha/2} \times \hat{\sigma}_{\hat{\mu}}]) = 1 - \alpha$, where $\hat{\mu}$ and $\hat{\sigma}_{\hat{\mu}}$ are the

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

³An ancillary statistic is a pivotal quantity that is also a statistic.

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

A remarkable difference between the Bayesian and the 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 0 or 1 under this approach, as μ can be there or not, the problem is that we will never know in applied settings. This due to that $P(\mu \in [\hat{\mu} - |t_{N-1}^{0.025}| \hat{\times} \sigma_{\hat{\mu}}, \hat{\mu} + |t_{N-1}^{0.025}| \times \hat{\sigma}_{\hat{\mu}}]) = 0.95$ being in the sense of repeated sampling. On the other hand, once we have the posterior distribution, we can say that $P(\mu \in [0.2, 0.4]) = 0.95$ under the Bayesian framework.

Following common practice, most of researchers and practitioners do hypothesis testing based on the p-value in the Frequentist framework. But, **what** is a p-value? Most of the users do not know the answer due to many time statistical inference is not performed by statisticians [10].⁴ 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. So,

"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." [50]

It seems that common Frequentist inferential practice intertwined two different logic reasoning arguments: the p-value [32] and significance level [72]. The former is an informal short-run criterion, whose philosophical foundation is reduction to absurdity, which measures the discrepancy between the data and the null hypothesis. So, the p-value is not a direct measure of the probability that the null hypothesis is false. The latter, whose philosophical foundations is deduction, is based on a long-run performance such that controls the overall number of incorrect inferences in the repeated sampling without care of individual cases. The p-value fallacy consists in interpreting the p-value as the strength of evidence against the null hypothesis, and using it simultaneously with the frequency of type I error under the null hypothesis [44].

The American Statistical Association has several concerns regarding the use of the p-value as a cornerstone to perform hypothesis testing in science. This concern motivates the ASA's statement on p-values [104], which can be summarized in the following principles:

• "P-values can indicate how incompatible the data are with a specified statistical model."

⁴https://fivethirtyeight.com/features/not-even-scientists-can-easily-explain-p-values/

- "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 only 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 difference between the Frequentists and the Bayesians is the way how scientific hypothesis are tested. The former use the p-value, whereas the latter use the Bayes factor. Observe that the p-value is associated with the probability of the data given the hypothesis, whereas the Bayes factor is associated with the probability of the hypothesis given the data. However, there is an approximate link between the t statistic and the Bayes factor for regression coefficients [80]. In particular, $|t| > (log(N)+6)^{1/2}$, corresponds to strong evidence in favor of rejecting the not relevance of a control in a regression. Observe that in this setting the threshold of the t statistic, and as a consequence the significant level, depends on the sample size. Observe that 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 [93] and [8] for nice exercises to reveal potential flaws of the p-value (p) due to $p \sim U[0,1]$ under the null hypothesis,⁵ and calibrations of the p-value to interpret them 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 $\alpha(p) = (1 + [-e \times p \times \log(p)]^{-1})^{-1}$ as the error probability α in rejecting H_0 . Take into account that B(p) and $\alpha(p)$ are lower bounds.

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

 $^{^5}$ https://joyeuserrance.wordpress.com/2011/04/22/proof-that-p-values-under-the-null-are-uniformly-distributed/ for a simple proof.

full assurance, of the true of the conclusion, that is why we get probabilistic statements.

Here, there is a difference between effects of causes (forward causal inference) and causes of effects (reverse causal inference) [39, 25]. To illustrate this point, imagine that a firm increases the price of a specific good, then economic theory would say that its demand decreases. The premise (null hypothesis) is a price increase, and the consequence is a demand reduction. Another view would be to observe a demand reduction, and try to identify which cause is more tenable. For instance, demand reduction can be caused by any positive supply shocks or any negative demand shocks. The Frequentist logic sees the first view, and the Bayesian reasoning gives the probability associated with possible causes.

2.4 The likelihood principle

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

We follow [9], who in turns followed [64], to illustrate the likelihood principle. We are given a coin such that we are interested in the probability, θ , of having it come up heads when flipped. It is desired 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, the results of which is the observation of 9 heads and 3 tails.

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

- 1. The experiment consisted of a predetermine 12 flips, so that Y = [Heads] would be $B(12, \theta)$, then $p_1(y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y} = 220 \times \theta^9 (1-\theta)^3$.
- 2. The experiment consisted of flipping the coin until 3 tails were observed (r=3). Then, Y, the number of failures (heads) until getting 3 tails, is $NB(3, 1-\theta)$. Then, $p_2(y|\theta) = {y+r-1 \choose r-1}(1-(1-\theta)^y(1-\theta)^r = 55 \times \theta^9(1-\theta)^3$.

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

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

R code. The likelihood principle: Binomial model

```
1 success <- 9
2 # Number of observed success in n trials
3 n <- 12
4 # Number of trials
5 siglevel <- sum(sapply(9:n,function(y)dbinom(y,n,0.5)))
6 siglevel
7 0.073</pre>
```

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

```
\alpha_2 = P_{1/2}(Y \ge 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
2 # Number of target success (tails)
3 failures <- 9
4 # Number of failures
5 siglevel <- 1 - pnbinom((failures - 1), success, 0.5)
6 siglevel
7 0.0327</pre>
```

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

2.5 Why is not the Bayesian approach that popular?

At this stage, we may wonder why the Bayesian statistical framework is not the dominant inferential approach despite that it has its historical origin in 1763 [7], whereas the Frequentist statistical framework was largely developed in the early 20th century. The scientific battle over the Bayesian inferential approach lasted for 150 years, and this maybe explained by some of the following facts.

There is an issue regarding apparent subjectivity as the Bayesian inferential approach runs counter the strong conviction that science demands objectivity, and Bayesian probability is a measure of degrees of belief, where the initial prior maybe just a guess; this was not accepted as objective and rigorous science. Initial critics said that Bayes was quantifying ignorance as he set equal probabilities to any potential result. As a consequence, prior distributions were damned [70].

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

However, it was Laplace who independently generalized Bayes' theorem in 1781. He used it initially in gambling problems, and soon after in astronomy, mixing different sources of information in order to leverage research in specific situations where data was scarce. Then, he wanted to use his discovery to find the probability of causes, and thought that this required large data sets, and turned into demography. In this field, he had to perform large calculations that demanded to develop smart approximations, creating the Laplace's approximation and the central limit theorem [59]; although, at the cost of apparently leaving his research on Bayesian inference.

Once Laplace was gone in 1827, the Bayes' rule disappeared from the scientific spectrum for almost a century. In part, personal attacks against Laplace made the rule be forgotten, and also, the old fashion thought that statistics does not have to say anything about causation, and that the prior is very subjective to be compatible with science. Although, practitioners used it to solve problems in astronomy, communication, medicine, military and social issues with remarkable results.

Thus, the concept of degrees of belief to operationalize probability was gone in name of scientific objectivity, and probability as the frequency an event occurs in many repeatable trials became the rule. Laplace critics argued that those concepts were diametric opposites, although, Laplace considered them as basically equivalent when large sample sizes are involved [70].

The era of the Frequentists or sampling theorists began, lead by Karl Pearson, and his nemesis, Ronald Fisher, both brilliant, against the inverse probability approach, persuasive and dominant characters that made impossible to argue against their ideas. Karl Pearson legacy was taken by his son Egon, and Egon's friend Neyman, both inherited the anti-Bayesian and anti-Fisher legacy.

Despite the anti-Bayesian campaign among statisticians, there were some independent characters developing Bayesian ideas, Borel, Ramsey and de Fineti, all of them isolated in different countries, France, England and Italy. However, the anti-Bayesian trio of Fisher, Neyman and Egon Person got all the attention during the 1920s and 1930s. Only, a geophysicist, Harold Jeffreys, kept alive Bayesian inference in the 1930s and 1940s. Jeffreys was a

very quiet, shy, uncommunicative gentleman working at Cambridge in the astronomy department. He was Fisher's friend thanks to his character, although they were diametric opposites regarding the Bayesian inferential approach, facing very high intellectual battles. Unfortunately for the Bayesian approach, Jeffreys lost, he was very technical using confusing high level mathematics, worried about inference from scientific evidence, not guiding future actions based on decision theory, which was very important in that era for mathematical statistics due to the Second World War. On the other hand, Fisher was a very dominant character, persuasive in public and a master of practice, his techniques were written in a popular style with minimum mathematics.

However, Bayes' rule achieved remarkable results in applied settings like the AT&T company or the social security system in USA. Bayesian inference also had a relevant role during the second World War and the Cold War. Alan Turing used inverse probability at Bletchley Park to crack German messages called Enigma code used by U-boats, Andrei Kolmogorov used it to improved firing tables of Russia's artillery, Bernard Koopman applied it for searching targets in the open sea and the RAND Corporation used it in the Cold War. Unfortunately, these Bayesian developments were top secrets for almost 40 years that keep classified the contribution of inverse probability in modern human history.

During 1950s and 1960s three mathematicians lead the rebirth of the Bayesian approach, Good, Savage and Lindley. However, it seems that they were unwilling to apply their theories to real problems, and despite that the Bayesian approach proved its worth, for instance, in business decisions, navy search, lung cancer, etc, it was applied to simple models due to its mathematical complexity and requirement of large computations. But, there were some breakthrough that change this. First, hierarchical models introduced by Lindley and Smith, where a complex model is decomposed into many easy to solve models, and second, Markov chain Monte Carlo methods developed by Hastings in the 1970s [47] and the Geman brothers in the 1980s [40]. These methods were introduced into the Bayesian inferential framework in the 1990s by Gelfand and Smith [33], and Tierney [101], when desktop computers got enough computational power to solve complex models. Since then, the Bayesian inferential framework has gained increasing popularity among practitioners and scientists.

A simple working example

We will illustrate some conceptual differences between the Bayesian and Frequentist statistical approaches performing inference given a random sample $\mathbf{y} = [y_1, y_2, \dots, y_N], \text{ where } y_i \stackrel{iid}{\sim} N(\mu, \sigma^2), i = 1, 2, \dots, N.$ 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 4), that is, this prior is perfectly compatible with sample information. In addition, we are assuming independent priors for μ and σ . Then,

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

$$= \frac{1}{\sigma} \times (\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^N ((y_i - \bar{y}) - (\mu - \bar{y}))^2\right\}$$

$$= \frac{1}{\sigma} \exp\left\{-\frac{N}{2\sigma^2} (\mu - \bar{y})^2\right\} \times (\sigma)^{-N} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \bar{y})^2\right\}$$

$$= \frac{1}{\sigma} \exp\left\{-\frac{N}{2\sigma^2} (\mu - \bar{y})^2\right\} \times (\sigma)^{-(\alpha_n + 1)} \exp\left\{-\frac{\alpha_n \hat{\sigma}^2}{2\sigma^2}\right\},$$

where $\bar{y} = \frac{\sum_{i=1}^{N}}{N}$, $\alpha_n = N - 1$ 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 den-

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 [108], $\sigma|\mathbf{y} \sim IG(\alpha_n, \hat{\sigma}^2)$. 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 μ 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 totally different. The confidence interval has a probabilistic interpretation under sampling variability of \bar{Y} , that is, in repeated sampling $(1-\alpha)\%$ of the intervals $\bar{Y}\pm|z_{\alpha/2}|\frac{\sigma}{\sqrt{N}}$ would include μ , but given an observed realization of \bar{Y} , say \bar{y} , the probability of $\bar{y}\pm|z_{\alpha/2}|\frac{\sigma}{\sqrt{N}}$ including μ is 1 or 0, that is why we say a $(1-\alpha)\%$ confidence interval. On the other hand, $\bar{y}\pm|z_{\alpha/2}|\frac{\sigma}{\sqrt{N}}$ has a simple probabilistic interpretation in the Bayesian framework, there is a $(1-\alpha)\%$ probability that μ lies in this interval.

If we want to get the marginal posterior density of μ ,

$$\pi(\mu|\mathbf{y}) = \int_0^\infty \pi(\mu, \sigma|\mathbf{y}) d\sigma$$

$$\propto \int_0^\infty \frac{1}{\sigma} \times (\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2\right\} d\sigma$$

$$= \int_0^\infty \left(\frac{1}{\sigma}\right)^{N+1} \exp\left\{-\frac{N}{2\sigma^2} \frac{\sum_{i=1}^N (y_i - \mu)^2}{N}\right\} d\sigma$$

$$= \left[\frac{2}{\Gamma(N/2)} \left(\frac{N \sum_{i=1}^N (y_i - \mu)^2}{2N}\right)^{N/2}\right]^{-1}$$

$$\propto \left[\sum_{i=1}^N (y_i - \mu)^2\right]^{-N/2}$$

$$= \left[\sum_{i=1}^N ((y_i - \bar{y}) - (\mu - \bar{y}))^2\right]^{-N/2}$$

$$= \left[\alpha_n \hat{\sigma}^2 + N(\mu - \bar{y})^2\right]^{-N/2}$$

$$\propto \left[1 + \frac{1}{\alpha_n} \left(\frac{\mu - \bar{y}}{\hat{\sigma}/\sqrt{N}}\right)^2\right]^{-(\alpha_n + 1)/2}.$$

The fourth line is due to having the kernel of a inverted gamma density with N degrees of freedom in the integral [108].

The last expression is 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 $\frac{\hat{\sigma}^2}{N}\left(\frac{\alpha_n}{\alpha_n-2}\right)$. Then, $\mu|\mathbf{y}\sim t\left(\bar{y},\frac{\hat{\sigma}^2}{N}\left(\frac{\alpha_n}{\alpha_n-2}\right),\alpha_n\right)$.

Observe that a $(1-\alpha)\%$ confidence interval and $(1-\alpha)\%$ credible interval have exactly the same expression, $\bar{y}\pm|t_{\alpha/2}^{\alpha_n}|\frac{\hat{\sigma}}{\sqrt{N}}$, where $t_{\alpha/2}^{\alpha_n}$ is the $\alpha/2$ per-

centile of a Student's t distribution. But again, the interpretations are totally different.

The mathematical similarity between the Frequentist and Bayesian expressions in this example is due to using 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, $Y_i \sim N(\mu, \sigma)$. The sample mean and variance are equal to 102 and 10, respectively. Assuming an improper prior equal to $1/\sigma$,

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

Summary 47

Using $\mu|\mathbf{y} \sim t\left(\bar{y}, \frac{\hat{\sigma}^2}{N}\left(\frac{\alpha_n}{\alpha_{n-2}}\right), \alpha_n\right)$, which implies that $\bar{y} \pm |t_{\alpha/2}^{\alpha_n}| \frac{\hat{\sigma}}{\sqrt{N}}$, where $\bar{y} = 102$, $\hat{\sigma}^2 = 10$ and $\alpha_n = 49$, the 95% confidence and credible intervals for μ are the same (101.1, 102.9), and $P(\mu > 103) = 1.49\%$ 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
7 LimInf
8 101.101
9 # Lower bound
10 LimSup <- y_bar + abs(qt(0.025, alpha)) * serror
11 LimSup
12 102.898
13 # Upper bound
14 y.cut <- 103
15 P <- 1-metRology::pt.scaled(y.cut, df = alpha, mean = y_bar, sd = serror)
16 P
17 0.0149
18 # Probability of mu greater than y.cut</pre>
```

2.7 Summary

The differences between the Bayesian and Frequentist inferential approaches are philosophical, including as pertains to the role of probability; pedagogical, in particular as relates to the use of inference to inform decision making; and methodological, as having differences in their mathematical and computational frameworks. Although at methodological level, the debate has become considerably muted, except for some aspects of inference, with the recognition that each approach has a great deal to contribute to statistical practice [43, 5, 54]. As Bradley Efron said "Computer-age statistical inference at its most successful **combines** elements of the two philosophies" [31].

2.8 Exercises

1. Jeffreys-Lindley's paradox

The **Jeffreys-Lindley's paradox** [50, 65] is an apparent disagreement between the Bayesian and Frequentist frameworks to a hypothesis testing situation.

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

- •Show that the posterior model probability for the model under the null is approximately 0.95. Assume $\pi(H_0) = \pi(H_1) = 0.5$, and $\pi(\theta)$ equal to U(0,1) under H_1 .
- •Show that the *p*-value for this hypothesis test is equal to 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 . $\mu = \mu_0$ vs 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 $\pi(\mu, \sigma) \propto 1/\sigma$ under the alternative hypothesis.

Show that

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

$$PO_{01} = \frac{p(\mathbf{y}|\mathcal{M}_0)}{p(\mathbf{y}|\mathcal{M}_1)}$$

$$= \frac{\Gamma((\alpha_n + 1)/2)}{\Gamma(1/2)\Gamma(\alpha_n/2)} (\alpha_n \hat{\sigma}^2/N)^{-1/2} \left[1 + \frac{(\mu_0 - \bar{y})^2}{\alpha_n \hat{\sigma}^2/N} \right]^{-\left(\frac{\alpha_n + 1}{2}\right)}$$

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 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 . $\mu = \mu_0$ vs 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?

Exercises 49

•Find the posterior model probability of the null model for each μ_0 .

Objective and subjective Bayesian approaches

Cornerstone models: Conjugate families

We will introduce conjugate families in basic statistical models with examples, solving them analytically and computationally using R. We will have some mathematical, and computational exercises in R.

4.1 Motivation of conjugate families

Observing the three fundamental pieces of Bayesian analysis: the posterior distribution (parameter inference), the marginal likelihood (hypothesis testing), and the predictive distribution (prediction), equations 4.1, 4.2 and 4.3, respectively,

$$\pi(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})}{p(\mathbf{y})},$$
(4.1)

$$p(\mathbf{y}) = \int_{\mathbf{\Theta}} p(\mathbf{y}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta},\tag{4.2}$$

and

$$p(\mathbf{Y}_0|\mathbf{y}) = \int_{\mathbf{\Theta}} p(\mathbf{Y}_0|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta}, \tag{4.3}$$

we can understand that some of the initial limitations of the application of the Bayesian analysis were associated with the absence of algorithms to draw from non-standard posterior distributions (equation 4.1), and the lack of analytical solutions of the marginal likelihood (equation 4.2) and the predictive distribution (equation 4.3). Both issues requiring computational power.

Although there were algorithms to sample from non-standard posterior distributions since the second half of the last century [71, 47, 40], their particular application in the Bayesian framework emerged later [33, 101], maybe until the increasing computational power of desktop computers. However, it is also common practice nowadays to use models that have standard conditional posterior distributions to mitigate computational requirements. In addition, nice mathematical tricks plus computational algorithms [34, 21, 22]

and approximations [102, 51] are used to obtain the marginal likelihood (prior predictive).

Despite these advances, there are two potentially conflicting desirable model specification features that we can see from equations 4.1, 4.2 and 4.3: (1) analytical solutions and (2) the posterior distribution in the same family as the prior distribution for a given likelihood. The latter is called *conjugate priors*, a family of priors that is closed under sampling [92].

These features are desirable as the former implies facility to perform hypothesis testing and predictive analysis, and the latter means invariance of the prior-to-posterior updating. Both features imply less computational burden.

We can easily achieve each of these features independently, for instance using improper priors for analytical tractability, and defining in a broad sense the family of prior distributions for prior conjugacy. However, these features are in conflict.

Fortunately, we can achieve these two nice characteristics if we assume that the data generating process is given by a distribution function in the exponential family. That is, given a random sample $\mathbf{y} = [y_1, y_2, \dots, y_N]^{\mathsf{T}}$, a probability density function $p(\mathbf{y}|\boldsymbol{\theta})$ belongs to the exponential family if it has the form

$$p(\mathbf{y}|\boldsymbol{\theta}) = \prod_{i=1}^{N} h(y_i) C(\boldsymbol{\theta}) \exp \left\{ \eta(\boldsymbol{\theta})^{\top} \mathbf{T}(y_i) \right\}$$

$$= h(\mathbf{y}) C(\boldsymbol{\theta})^{N} \exp \left\{ \eta(\boldsymbol{\theta})^{\top} \mathbf{T}(\mathbf{y}) \right\}$$

$$= h(\mathbf{y}) \exp \left\{ \eta(\boldsymbol{\theta})^{\top} \mathbf{T}(\mathbf{y}) - A(\boldsymbol{\theta}) \right\},$$
(4.4)

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

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

$$p(\mathbf{y}|\boldsymbol{\theta}) = h(\mathbf{y})D(\boldsymbol{\eta})^N \exp\left\{\boldsymbol{\eta}^\top \mathbf{T}(\mathbf{y})\right\}$$
$$= h(\mathbf{y}) \exp\left\{\boldsymbol{\eta}^\top \mathbf{T}(\mathbf{y}) - B(\boldsymbol{\eta})\right\}.$$

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

4.1.1 Examples of exponential family distributions

1. Discrete distributions

Let's show that some of the most common distributions for random variables that 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.

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

$$= \frac{\lambda^{\sum_{i=1}^{N} y_i} \exp(-N\lambda)}{\prod_{i=1}^{N} y_i!}$$

$$= \frac{\exp(-N\lambda) \exp(\sum_{i=1}^{N} y_i \log(\lambda))}{\prod_{i=1}^{N} y_i!},$$

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

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.

$$p(\mathbf{y}|\theta) = \prod_{i=1}^{N} \theta^{y_i} (1 - \theta)^{1 - y_i}$$

= $\theta^{\sum_{i=1}^{N} y_i} (1 - \theta)^{N - \sum_{i=1}^{N} y_i}$
= $(1 - \theta)^N \exp\left\{\sum_{i=1}^{N} y_i \log\left(\frac{\theta}{1 - \theta}\right)\right\}$,

then
$$h(\mathbf{y}) = \mathbb{I}[y_i \in \{0, 1\}]$$
 (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 $\mathbf{y}_i = [y_{i1}, \dots, y_{im}], \sum_{l=1}^m y_{il} = n, n$ 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.

$$p(\mathbf{y}|\boldsymbol{\theta}) = \prod_{i=1}^{N} \frac{n!}{\prod_{l=1}^{m} y_{il}!} \prod_{l=1}^{m} \theta_{l}^{y_{il}}$$

$$= \frac{(n!)^{N}}{\prod_{i=1}^{N} \prod_{l=1}^{m} y_{il}!} \exp\left\{\sum_{i=1}^{N} \sum_{l=1}^{m} y_{il} \log(\theta_{l})\right\}$$

$$= \frac{(n!)^{N}}{\prod_{i=1}^{N} \prod_{l=1}^{m} y_{il}!} \exp\left\{\left(N \times n - \sum_{i=1}^{N} \sum_{l=1}^{m-1} y_{il}\right) \log(\theta_{m}) + \sum_{i=1}^{N} \sum_{l=1}^{m-1} y_{il} \log(\theta_{l})\right\}$$

$$= \frac{(n!)^{N}}{\prod_{i=1}^{N} \prod_{l=1}^{m} y_{il}!} \theta_{m}^{N \times n} \exp\left\{\sum_{i=1}^{N} \sum_{l=1}^{m-1} y_{il} \log(\theta_{l}/\theta_{m})\right\},$$
then $h(\mathbf{y}) = \frac{(n!)^{N}}{\prod_{i=1}^{N} \prod_{l=1}^{m} y_{il}!}, \quad \eta(\boldsymbol{\theta}) = \left[\log\left(\frac{\theta_{1}}{\theta_{m}}\right) \dots \log\left(\frac{\theta_{m-1}}{\theta_{m}}\right)\right],$

$$T(\mathbf{y}) = \left[\sum_{i=1}^{N} y_{i1} \dots \sum_{i=1}^{N} y_{im-1}\right] \text{ and } C(\boldsymbol{\theta}) = \theta_{m}^{n}.$$

2. Continuous distributions

Let's show that some of the most common distributions for random variables that can take any value within a certain range or interval, often 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.

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

then
$$h(\mathbf{y}) = (2\pi)^{-N/2}$$
, $\eta(\mu, \sigma^2) = \left[\frac{\mu}{\sigma^2} \frac{-1}{2\sigma^2}\right]$, $T(\mathbf{y}) = \left[\sum_{i=1}^{N} y_i \sum_{i=1}^{N} y_i^2\right]$ 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
$$B(\eta) = \frac{N}{2} \log(-2\eta_2) - \frac{N}{4} \frac{\eta_1^2}{\eta_2}$$
. Then,

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

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.

$$p(\mathbf{y}_{1}, \dots, \mathbf{y}_{N} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^{N} (2\pi)^{-p/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left\{-\frac{1}{2} \left(\mathbf{y}_{i} - \boldsymbol{\mu}\right)^{\top} \boldsymbol{\Sigma}^{-1} \left(\mathbf{y}_{i} - \boldsymbol{\mu}\right)\right\}$$

$$= (2\pi)^{-pN/2} |\boldsymbol{\Sigma}|^{-N/2} \exp\left\{-\frac{1}{2} tr \left[\sum_{i=1}^{N} \left(\mathbf{y}_{i} - \boldsymbol{\mu}\right)^{\top} \boldsymbol{\Sigma}^{-1} \left(\mathbf{y}_{i} - \boldsymbol{\mu}\right)\right]\right\}$$

$$= (2\pi)^{-pN/2} |\boldsymbol{\Sigma}|^{-N/2} \exp\left\{-\frac{1}{2} tr \left[\left(\mathbf{S} + N \left(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}\right) \left(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}\right)^{\top}\right) \boldsymbol{\Sigma}^{-1}\right]\right\}$$

$$= (2\pi)^{-pN/2} \exp\left\{-\frac{1}{2} \left[\left(vec(\mathbf{S})^{\top} + Nvec\left(\hat{\boldsymbol{\mu}}\hat{\boldsymbol{\mu}}^{\top}\right)^{\top}\right) vec(\boldsymbol{\Sigma}^{-1})\right]$$

$$-2N\hat{\boldsymbol{\mu}}^{\top} vec\left(\boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1}\right) + Ntr\left(\boldsymbol{\mu}\boldsymbol{\mu}^{\top} \boldsymbol{\Sigma}^{-1}\right) + N\log(|\boldsymbol{\Sigma}|)\right\},$$

where the second line uses the trace operator (tr), and its invariability under cyclic permutation is used in the third line. In addition, we add and subtract $\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{y}_{i}$ in each parenthesis such that we get $\mathbf{S} = \sum_{i=1}^{N} (\mathbf{y}_{i} - \hat{\boldsymbol{\mu}}) (\mathbf{y}_{i} - \hat{\boldsymbol{\mu}})^{\top}$. We get the fourth line after collecting terms, and using some properties of the trace operator to introduce the vectorization operator (vec), that is, $tr(\mathbf{AB}) = vec(\mathbf{A}^{\top})^{\top} vec(\mathbf{B})$, and $vec(\mathbf{A} + \mathbf{B}) = vec(\mathbf{A}) + vec(\mathbf{B})$.

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

4.2 Conjugate prior to exponential family

Theorem 4.2.1

The prior distribution $\pi(\boldsymbol{\theta}) \propto C(\boldsymbol{\theta})^{b_0} \exp \left\{ \eta(\boldsymbol{\theta})^{\top} \mathbf{a}_0 \right\}$ is conjugate to the exponential family (equation 4.4).

Proof

$$\pi(\boldsymbol{\theta}|\mathbf{y}) \propto C(\boldsymbol{\theta})^{b_0} \exp\left\{\eta(\boldsymbol{\theta})^{\top} \mathbf{a}_0\right\} \times h(\mathbf{y}) C(\boldsymbol{\theta})^N \exp\left\{\eta(\boldsymbol{\theta})^{\top} \mathbf{T}(\mathbf{y})\right\}$$

$$\propto C(\boldsymbol{\theta})^{N+b_0} \exp\left\{\eta(\boldsymbol{\theta})^{\top} (\mathbf{T}(\mathbf{y}) + \mathbf{a}_0)\right\}.$$

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

Remarks

We see 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 view helps the elicitation process.

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

In the setting where there is a regular conjugate prior, [28] show that we obtain a posterior expectation of the sufficient statistics that is a weighted average between the prior expectation and the likelihood estimate.

4.2.1 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 λ has the form

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

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

Taking into account that $\sum_{i=1}^{N} y_i$ is a sufficient statistic for the Poisson distribution, then we can think about a_0 as the number of occurrences in b_0 experiments. Observe that

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

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

Observe that α_0/β_0 is the prior mean, and α_0/β_0^2 is the prior variance. Then, $\alpha_0 \to 0$ and $\beta_0 \to 0$ imply a non-informative prior such that the posterior mean converges to the maximum likelihood estimator $\bar{y} = \frac{\sum_{i=1}^{N} y_i}{N}$,

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

$$\mathbb{E} [\lambda | \mathbf{y}] = \frac{\alpha_n}{\beta_n}$$

$$= \frac{\sum_{i=1}^{N} y_i + \alpha_0}{N + \beta_0}$$

$$= \frac{N\bar{y}}{N + \beta_0} + \frac{\alpha_0}{N + \beta_0}.$$

The posterior mean is a weighted average between sample and prior information. This is a general result from regular conjugate priors [28]. Observe that $\mathbb{E}[\lambda|\mathbf{y}] = \bar{y}, \lim N \to \infty$.

In addition, $\alpha_0 \to 0$ and $\beta_0 \to 0$ corresponds to $\pi(\lambda) \propto \frac{1}{\lambda}$, which is an improper prior. Improper priors may have bad consequences on Bayes factors (hypothesis testing), see bellow a discussion of this in the linear regression framework. In this example, we can get 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 θ has the form

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

This is the kernel of a beta density, $B(\alpha_0, \beta_0)$, $\alpha_0 = a_0 + 1$ 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, which is the number of successes, then $b_0 - a_0$ is the number of failures. This implies that α_0 is the number of prior successes plus one, and β_0 is the number of prior failures plus one. Given that the mode of a beta distributed random variable is $\frac{\alpha_0 - 1}{\alpha_0 + \beta_0 - 2} = \frac{a_0}{b_0}$, then we have the prior probability of success. Setting $\alpha_0 = 1$ and $\beta_0 = 1$, which implies a 0-1 uniform distribution, corresponds to a setting with 0 successes (and 0 failures) in 0 experiments.

Observe that

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

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.

El 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

$$\begin{split} p(Y_0|\mathbf{y}) &= \int_0^1 \theta^{y_0} (1-\theta)^{1-y_0} \times \frac{\theta^{\alpha_n-1} (1-\theta)^{\beta_n-1}}{B(\alpha_n,\beta_n)} d\theta \\ &= \frac{B(\alpha_n+y_0,\beta_n+1-y_0)}{B(\alpha_n,\beta_n)} \\ &= \frac{\Gamma(\alpha_n+\beta_n)\Gamma(\alpha_n+y_0)\Gamma(\beta_n+1-y_0)}{\Gamma(\alpha_n+\beta_n+1)\Gamma(\alpha_n)\Gamma(\beta_n)} \\ &= \left\{ \frac{\frac{\alpha_n}{\alpha_n+\beta_n}, \quad y_0 = 1}{\frac{\beta_n}{\alpha_n+\beta_n}, \quad y_0 = 0} \right\}. \end{split}$$

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

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

where
$$\boldsymbol{\eta}(\boldsymbol{\theta}) = \left[\log\left(\frac{\theta_1}{\theta_m}\right), \dots, \log\left(\frac{\theta_{m-1}}{\theta_m}\right)\right], \ \mathbf{a}_0 = \left[a_{01}, \dots, a_{am-1}\right]^\top,$$

$$\boldsymbol{\alpha}_0 = \left[\alpha_{01}, \alpha_{02}, \dots, \alpha_{0m}\right], \ \alpha_{0l} = a_{0l} + 1, \ l = 1, 2, \dots, m-1 \text{ and }$$

$$\alpha_{0m} = b_0 - \sum_{l=1}^{m-1} a_{0l} + 1.$$

This is the kernel of a Dirichlet distribution, that is, the prior distribution is $D(\alpha_0)$.

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

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

$$\pi(\boldsymbol{\theta}|\mathbf{y}) \propto \prod_{l=1}^{m} \theta_{l}^{\alpha_{0l}-1} \times \prod_{l=1}^{m} \theta_{l}^{\sum_{i=1}^{N} y_{il}}$$
$$= \prod_{l=1}^{m} \theta_{l}^{\alpha_{0l} + \sum_{i=1}^{N} y_{il} - 1}.$$

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

$$\mathbb{E}[\theta_{j}|\mathbf{y}] = \frac{\alpha_{nj}}{\sum_{l=1}^{m} \left[\alpha_{0l} + \sum_{i=1}^{N} y_{il}\right]}$$

$$= \frac{\sum_{l=1}^{m} \alpha_{0l}}{\sum_{l=1}^{m} \left[\alpha_{0l} + \sum_{i=1}^{N} y_{il}\right]} \frac{\alpha_{0j}}{\sum_{l=1}^{m} \alpha_{0l}}$$

$$+ \frac{\sum_{l=1}^{m} \sum_{i=1}^{N} y_{il}}{\sum_{l=1}^{m} \sum_{i=1}^{N} y_{il}} \frac{\sum_{i=1}^{N} y_{ij}}{\sum_{l=1}^{m} \sum_{i=1}^{N} y_{il}}.$$

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

The marginal likelihood is

$$p(\mathbf{y}) = \int_{\mathbf{\Theta}} \frac{\prod_{l=1}^{m} \theta_{l}^{\alpha_{0l}-1}}{B(\boldsymbol{\alpha}_{0})} \times \prod_{i=1}^{N} \frac{n!}{\prod_{l=1}^{m} y_{il}} \prod_{l=1}^{m} \theta_{l}^{y_{il}} d\boldsymbol{\theta}$$

$$= \frac{N \times n!}{B(\boldsymbol{\alpha}_{0}) \prod_{i=1}^{N} \prod_{l=1}^{m} y_{il}!} \int_{\mathbf{\Theta}} \prod_{l=1}^{m} \theta_{l}^{\alpha_{0l} + \sum_{i=1}^{N} y_{il} - 1} d\boldsymbol{\theta}$$

$$= \frac{N \times n!}{B(\boldsymbol{\alpha}_{0}) \prod_{i=1}^{N} \prod_{l=1}^{m} y_{il}!} B(\boldsymbol{\alpha}_{n})$$

$$= \frac{N \times n! \Gamma(\sum_{l=1}^{m} \alpha_{0l})}{\Gamma(\sum_{l=1}^{m} \alpha_{0l} + N \times n)} \prod_{l=1}^{m} \frac{\Gamma(\alpha_{nl})}{\Gamma(\alpha_{0l}) \prod_{i=1}^{N} y_{il}!},$$

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

$$p(Y_0|\mathbf{y}) = \frac{n!\Gamma\left(\sum_{l=1}^{m} \alpha_{nl}\right)}{\Gamma\left(\sum_{l=1}^{m} \alpha_{nl} + n\right)} \prod_{l=1}^{m} \frac{\Gamma\left(\alpha_{nl} + y_{0l}\right)}{\Gamma\left(\alpha_{nl}\right) y_{0l}!}.$$

This is a Dirichlet-multinomial distribution with parameters α_n .

Example: English premier league, Liverpool vs Manchester city

Let's see an example based on data from the English Premier league. In particular, we want to get the probability that in the following five matches Liverpool versus Manchester city, the former wins two games, and the latter three game. This is done based on the historical records of the last five matches where Liverpool was local between January 14th, 2018 and April tenth, 2022. There were two wins for Liverpool, two draws, and one win for Manchester city.²

We use two strategies to get the hyperparameters. First, we estimate the hyperparameters of the Dirichlet distribution using betting odds from bookmakers at 19:05 hours October sixth, 2022 (Colombia time). We got information from 24 bookmakers (see file DataOddsLIVvsMAN.csv),³ and transform these odds in probabilities using a simple standardization approach, then we use maximum

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

 $^{^3} https://www.oddsportal.com/soccer/england/premier-league/liverpool-manchester-city-WrqgEz5S/$

likelihood to estimate the hyperparameters. Second, we use empirical Bayes, that is, we estimate the hyperparameters optimizing the marginal likelihood.

$R\ code.\ Multinomial\text{-}Dirichlet\ model:\ Liverpool\ vs}\ Manchester\ city$

```
1 # Multinomial-Dirichlet example: Liverpool vs Manchester
2 Data <-read.csv("DataApplications/DataOddsLIVvsMAN.csv", sep</pre>
      = ",", header = TRUE)
3 # Change path
4 attach (Data)
5 library(dplyr)
6 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)
11 # We get probabilities using simple standardization. There
      are more technical approaches to do this. See for
      instance Shin (1993) and Strumbelj (2014).
12 DirMLE <- sirt::dirichlet.mle(Probs)</pre>
13 # Use maximum likelihood to estimate parameters of the
14 # Dirichlet distribution
15 alphaOodds <- DirMLE$alpha
16 alpha0odds
        p1
                  p2
18 1599.122 1342.703 2483.129
20 y < -c(2, 2, 1)
21 # Historical records last five mathces
^{22} # Liverpool wins (2), draws (2) and Manchester
23 # city wins (1)
25 # Marginal likelihood
26 MarLik <- function(a0){
   n <- sum(y)
    Res1 <- sum(sapply(1:length(y),</pre>
    function(1){lgamma(a0[1]+y[1])-lgamma(a0[1])}))
    Res <- lgamma(sum(a0))-lgamma(sum(a0)+n)+Res1
    return(-Res)
32 }
33 EmpBay <- optim(alpha0odds, MarLik, method = "BFGS")
34 alpha0EB <- EmpBay$par
35 alpha0EB
36 p1
           p2
37 2362.622 2660.153 1279.510
38 # Bayes factor empirical Bayes vs betting odds.
39 # This is greather than 1 by construction
40 BF <- exp(-MarLik(alpha0EB))/exp(-MarLik(alpha0odds))
41 BF
42 2.085819
43 # Posterior distribution based on empirical Bayes
44 alphan <- alpha0EB + y
45 # Posterior parameters
46 S <- 100000
47 # Simulation draws from the Dirichlet distribution
48 thetas <- MCMCpack::rdirichlet(S, alphan)
49 colnames (thetas) <- c("Liverpool", "Draw", "Manchester")
```

$R\ code.\ Multinomial\text{-}Dirichlet\ model:\ Liverpool\ vs}\ Manchester\ city$

```
# Predictive distribution based on simulations
_{2} y0 <- c(2, 0, 3)
  # Liverpool two wins and Manchester city three wins in next
      five matches
  Pred <- apply(thetas, 1, function(p) {rmultinom(1, size =</pre>
      sum(y0), prob = p)
  ProYo <- sum(sapply(1:S,function(s){sum(Pred[,s]==y0)==3}))/</pre>
  ProY0
  0.0832
8 # Probability of y0
10 # Predictive distribution using analytical expression
  PredY0 <- function(y0){</pre>
    n <- sum(y0)
    Res1 <- sum(sapply(1:length(y), function(1){lgamma(alphan[</pre>
      1]+y0[1]) - lgamma(alphan[1])-lfactorial(y0[1])}))
    Res <- lfactorial(n) + lgamma(sum(alphan)) - lgamma(sum(
      alphan)+n) + Res1
    return(exp(Res))
16 }
17 PredYO(y0)
18 0.0833
```

We see that the Bayes factor gives evidence in favor of the hyperparameters based on empirical Bayes, this is by construction, as these hyperparameters maximize the marginal likelihood.

We observe that using the hyperparameters from empirical Bayes, the probability that in the next five games Liverpool wins two games and Manchester city wins three games is 8.33%. The result using the predictive distribution based on simulations is similar to the probability using the exact predictive.

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

$$\pi(\mu, \sigma^{2}) \propto \exp\left\{b_{0}\left(-\frac{\mu^{2}}{2\sigma^{2}} - \frac{\log \sigma^{2}}{2}\right)\right\} \exp\left\{a_{01}\frac{\mu}{\sigma^{2}} - a_{02}\frac{1}{\sigma^{2}}\right\}$$

$$= \exp\left\{b_{0}\left(-\frac{\mu^{2}}{2\sigma^{2}} - \frac{\log \sigma^{2}}{2}\right)\right\} \exp\left\{a_{01}\frac{\mu}{\sigma^{2}} - a_{02}\frac{1}{\sigma^{2}}\right\}$$

$$\times \exp\left\{-\frac{a_{01}^{2}}{2\sigma^{2}b_{0}}\right\} \exp\left\{\frac{a_{01}^{2}}{2\sigma^{2}b_{0}}\right\}$$

$$= \exp\left\{-\frac{b_{0}}{2\sigma^{2}}\left(\mu - \frac{a_{01}}{b_{0}}\right)^{2}\right\} \left(\frac{1}{\sigma^{2}}\right)^{\frac{b_{0}+1-1}{2}}$$

$$\times \exp\left\{\frac{1}{\sigma^{2}} - \frac{2b_{0}a_{02} + a_{01}^{2}}{2b_{0}}\right\}$$

$$= \underbrace{\left(\frac{1}{\sigma^{2}}\right)^{\frac{1}{2}} \exp\left\{-\frac{b_{0}}{2\sigma^{2}}\left(\mu - \frac{a_{01}}{b_{0}}\right)^{2}\right\}}_{1}$$

$$\times \underbrace{\left(\frac{1}{\sigma^{2}}\right)^{\frac{b_{0}-1}{2}} \exp\left\{-\frac{1}{\sigma^{2}} \frac{2b_{0}a_{02} - a_{01}^{2}}{2b_{0}}\right\}}_{2}.$$

The first part is the kernel of a normal density with mean $\mu_0 = a_{01}/\beta_0$ and variance σ^2/β_0 , $\beta_0 = b_0$ that is, $\mu|\sigma^2 \sim N(\mu_0, \sigma^2/\beta_0)$. The second part is the kernel of an inverse gamma density with shape parameter $\alpha_0/2 = \frac{\beta_0 - 3}{2}$, and scale parameter $\delta_0/2 = \frac{2\beta_0 a_{02} - a_{01}^2}{2\beta_0}$, $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$. Observe that $b_0 = \beta_0$ is the hypothetical sample size, and a_{01} is the hypothetical sum of prior observations, then, it makes sense that a_{01}/β_0 and σ^2/β_0 are the prior mean and variance, respectively.

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

distribution,

$$\pi(\mu, \sigma^{2}|\mathbf{y}) \propto \left(\frac{1}{\sigma^{2}}\right)^{1/2} \exp\left\{-\frac{\beta_{0}}{2\sigma^{2}}(\mu - \mu_{0})^{2}\right\} \left(\frac{1}{\sigma^{2}}\right)^{\alpha_{0}/2+1} \exp\left\{-\frac{\delta_{0}}{2\sigma^{2}}\right\}$$

$$\times (\sigma^{2})^{-N/2} \exp\left\{-\frac{1}{2\sigma^{2}}\sum_{i=1}^{N}(y_{i} - \mu)^{2}\right\}$$

$$= \left(\frac{1}{\sigma^{2}}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^{2}}\left(\beta_{0}(\mu - \mu_{0})^{2} + \sum_{i=1}^{N}(y_{i} - \bar{y})^{2} + N(\mu - \bar{y})^{2} + \delta_{0}\right)\right\}$$

$$\times \left(\frac{1}{\sigma^{2}}\right)^{\frac{\alpha_{0}+N}{2}+1} + \frac{(\beta_{0}\mu_{0} + N\bar{y})^{2}}{\beta_{0} + N} - \frac{(\beta_{0}\mu_{0} + N\bar{y})^{2}}{\beta_{0} + N}$$

$$= \left(\frac{1}{\sigma^{2}}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^{2}}\left((\beta_{0} + N)\left(\mu - \left(\frac{\beta_{0}\mu_{0} + N\bar{y}}{\beta_{0} + N}\right)\right)^{2}\right)\right\}$$

$$\times \left(\frac{1}{\sigma^{2}}\right)^{\frac{\alpha_{0}+N}{2}+1} \exp\left\{-\frac{1}{2\sigma^{2}}\left(\sum_{i=1}^{N}(y_{i} - \bar{y})^{2} + \delta_{0} + \frac{\beta_{0}N}{\beta_{0} + N}(\bar{y} - \mu_{0})^{2}\right)\right\}.$$

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 $(\beta_0 \text{ and } N)$.

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

$$\pi(\mu|\mathbf{y}) \propto \int_0^\infty \left\{ \left(\frac{1}{\sigma^2}\right)^{\frac{\alpha_n+1}{2}+1} \exp\left\{-\frac{1}{2\sigma^2} (\beta_n(\mu-\mu_n)^2 + \delta_n)\right\} \right\} d\sigma^2$$

$$= \frac{\Gamma\left(\frac{\alpha_n+1}{2}\right)}{\left[\frac{\beta_n(\mu-\mu_n)^2+\delta_n}{2}\right]^{\frac{\alpha_n+1}{2}}}$$

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

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

where the second line due to having the kernel of an inverse gamma density with parameters $(\alpha_n + 1)/2$ and $-\frac{1}{2\sigma^2}(\beta_n(\mu - \mu_n)^2 + \delta_n)$.

This is the kernel of a Student's t distribution, $\mu|\mathbf{y} \sim t(\mu_n, \delta_n/\beta_n\alpha_n, \alpha_n)$, where $\mathbb{E}[\mu|\mathbf{y}] = \mu_n$ and $Var[\mu|\mathbf{y}] = \frac{\alpha_n}{\alpha_n-2} \left(\frac{\delta_n}{\beta_n\alpha_n}\right) = \frac{\delta_n}{(\alpha_n-2)\beta_n}$, $\alpha_n > 2$. Observe that the marginal posterior distribution for μ has heavier tails than the conditional posterior distribution due to incorporating uncertainty regarding σ^2 . The marginal likelihood is

$$\begin{split} 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\} \frac{(\delta_{0}/2)^{\alpha_{0}/2}}{\Gamma(\alpha_{0}/2)} \left(\frac{1}{\sigma^{2}} \right)^{\alpha_{0}/2+1} \\ &\times \exp\left\{ -\frac{\delta_{0}}{2\sigma^{2}} \right\} (2\pi\sigma^{2})^{-N/2} \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)^{-\left(\frac{N+1}{2}\right)} \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. \\ &\times \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)^{-\left(\frac{N+1}{2}\right)} \beta_{0}^{1/2} \Gamma\left(\frac{N+1+\alpha_{0}}{2} \right) \\ &\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)^{-\left(\frac{N+1}{2}\right)} \beta_{0}^{1/2} \Gamma\left(\frac{N+1+\alpha_{0}}{2} \right) \\ &\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}} \\ &= \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}} \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{split}$$

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

$$\pi(Y_{0}|\mathbf{y}) \propto \int_{-\infty}^{\infty} \int_{0}^{\infty} \left\{ \left(\frac{1}{\sigma^{2}}\right)^{1/2} \exp\left\{-\frac{1}{2\sigma^{2}}(y_{0}-\mu)^{2}\right\} \left(\frac{1}{\sigma^{2}}\right)^{1/2} \exp\left\{-\frac{\beta_{n}}{2\sigma^{2}}(\mu-\mu_{n})^{2}\right\} \right\} \\
\times \left(\frac{1}{\sigma^{2}}\right)^{\alpha_{n}/2+1} \exp\left\{-\frac{\delta_{n}}{2\sigma^{2}}\right\} d\sigma^{2} d\mu \\
= \int_{-\infty}^{\infty} \int_{0}^{\infty} \left\{ \left(\frac{1}{\sigma^{2}}\right)^{\frac{\alpha_{n}+2}{2}+1} \exp\left\{-\frac{1}{2\sigma^{2}}((y_{0}-\mu)^{2}+\beta_{n}(\mu-\mu_{n})^{2}+\delta_{n})\right\} \right\} d\sigma^{2} d\mu \\
\propto \int_{-\infty}^{\infty} \left[\beta_{n}(\mu-\mu_{n})^{2}+(y_{0}-\mu)^{2}+\delta_{n}\right]^{-\left(\frac{\alpha_{n}}{2}+1\right)} d\mu \\
= \int_{-\infty}^{\infty} \left[(\beta_{n}+1) \left(\mu-\left(\frac{\beta_{n}\mu_{n}+y_{0}}{\beta_{n}+1}\right)\right)^{2}+\frac{\beta_{n}(y_{0}-\mu_{n})^{2}}{\beta_{n}+1}+\delta_{n} \right]^{-\left(\frac{\alpha_{n}}{2}+1\right)} d\mu \\
= \int_{-\infty}^{\infty} \left[1+\frac{(\beta_{n}+1)^{2} \left(\mu-\left(\frac{\beta_{n}\mu_{n}+y_{0}}{\beta_{n}+1}\right)\right)^{2}}{\beta_{n}(y_{0}-\mu_{n})^{2}+(\beta_{n}+1)\delta_{n}} \right]^{-\left(\frac{\alpha_{n}}{2}+1\right)} d\mu \\
\times \left(\frac{\beta_{n}(y_{0}-\mu_{n})^{2}+(\beta_{n}+1)\delta_{n}}{\beta_{n}+1}\right)^{-\left(\frac{\alpha_{n}}{2}+1\right)} \\
\propto \left(\frac{\beta_{n}(y_{0}-\mu_{n})^{2}+(\beta_{n}+1)\delta_{n}}{(\beta_{n}+1)^{2}(\alpha_{n}+1)}\right)^{\frac{1}{2}} \left(\frac{\beta_{n}(y_{0}-\mu_{n})^{2}+(\beta_{n}+1)\delta_{n}}{\beta_{n}+1}\right)^{-\left(\frac{\alpha_{n}}{2}+1\right)} \\
\propto (\beta_{n}(y_{0}-\mu_{n})^{2}+(\beta_{n}+1)\delta_{n})^{\left(\frac{\alpha_{n}+1}{2}\right)} \\
\propto \left[1+\frac{\beta_{n}\alpha_{n}}{(\beta_{n}+1)\delta_{n}\alpha_{n}}(y_{0}-\mu_{n})^{2}\right]^{-\left(\frac{\alpha_{n}+1}{2}\right)},$$

where we have that $\left[1 + \frac{(\beta_n+1)^2\left(\mu - \left(\frac{\beta_n\mu_n + y_0}{\beta_n+1}\right)\right)^2}{\beta_n(y_0 - \mu_n)^2 + (\beta_n+1)\delta_n}\right]^{-\left(\frac{\alpha_n}{2} + 1\right)}$ is the ker-

nel of a Student's t density with degrees of freedom α_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 4.1 that the multivariate normal distribution is in the exponential family where

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

and

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

Then, its conjugate prior distribution should have the form

$$\pi(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp\left\{-\frac{b_0}{2} \left(tr\left(\boldsymbol{\mu}\boldsymbol{\mu}^{\top}\boldsymbol{\Sigma}^{-1}\right) + \log(|\boldsymbol{\Sigma}|)\right)\right\}$$

$$\times \exp\left\{\mathbf{a}_{01}^{\top}vec\left(\boldsymbol{\Sigma}^{-1}\right) + \mathbf{a}_{02}^{\top}vec\left(\boldsymbol{\mu}^{\top}\boldsymbol{\Sigma}^{-1}\right)\right\}$$

$$= |\boldsymbol{\Sigma}|^{-b_0/2} \exp\left\{-\frac{b_0}{2} \left(tr\left(\boldsymbol{\mu}^{\top}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}\right)\right) + tr\left(\mathbf{a}_{02}^{\top}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}\right)\right\}$$

$$\times \exp\left\{\mathbf{a}_{01}^{\top}vec\left(\boldsymbol{\Sigma}^{-1}\right) + \frac{\mathbf{a}_{02}^{\top}\boldsymbol{\Sigma}^{-1}\mathbf{a}_{02}}{2b_0} - \frac{\mathbf{a}_{02}^{\top}\boldsymbol{\Sigma}^{-1}\mathbf{a}_{02}}{2b_0}\right\}$$

$$= |\boldsymbol{\Sigma}|^{-b_0/2} \exp\left\{-\frac{b_0}{2} \left(\boldsymbol{\mu} - \frac{\mathbf{a}_{02}}{b_0}\right)^{\top} \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{\mu} - \frac{\mathbf{a}_{02}}{b_0}\right)\right\}$$

$$\times \exp\left\{-\frac{1}{2}tr\left(\left(\mathbf{A}_{01} - \frac{\mathbf{a}_{02}\mathbf{a}_{02}^{\top}}{b_0}\right) \boldsymbol{\Sigma}^{-1}\right)\right\}$$

$$= |\boldsymbol{\Sigma}|^{-1/2} \exp\left\{-\frac{b_0}{2} \left(\boldsymbol{\mu} - \frac{\mathbf{a}_{02}}{b_0}\right)^{\top} \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{\mu} - \frac{\mathbf{a}_{02}}{b_0}\right)\right\}$$

$$\times |\boldsymbol{\Sigma}|^{-(\alpha_0 + p + 1)/2} \exp\left\{-\frac{1}{2}tr\left(\left(\mathbf{A}_{01} - \frac{\mathbf{a}_{02}\mathbf{a}_{02}^{\top}}{b_0}\right) \boldsymbol{\Sigma}^{-1}\right)\right\},$$

where b_0 is the hypothetical sample size, and \mathbf{a}_{01} and \mathbf{a}_{02} are p^2 and p dimensional vectors of prior sufficient statistics, where $\mathbf{a}_{01} = -\frac{1}{2}vec(\mathbf{A}_{01})$ such that \mathbf{A}_{01} is a $p \times p$ positive semi-definite matrix. Setting $b_0 = 1 + \alpha_0 + p + 1$, we have that the first part in the last expression is the kernel of a multivariate normal density with mean $\boldsymbol{\mu}_0 = \mathbf{a}_{02}/b_0$ and covariance $\frac{\boldsymbol{\Sigma}}{b_0}$, that is, $\boldsymbol{\mu}|\boldsymbol{\Sigma} \sim N_p\left(\boldsymbol{\mu}_0, \frac{\boldsymbol{\Sigma}}{\beta_0}\right)$, $b_0 = \beta_0$. It makes sense these hyperparameters because \mathbf{a}_{02} is the hypothetical sum of prior observations and b_0 is the hypothetical prior sample size. In addition, the second expression in the last line is the kernel of a inverse Wishart distribution with scale matrix $\boldsymbol{\Psi}_0 = \left(\mathbf{A}_{01} - \frac{\mathbf{a}_{02}\mathbf{a}_{02}}{b_0}\right)$ and α_0 degrees of freedom, that is, $\boldsymbol{\Sigma} \sim IW_p(\boldsymbol{\Psi}_0, \alpha_0)$. Observe that $\boldsymbol{\Psi}_0$ has the same structure as the first part of the sufficient statistics in $T(\mathbf{y})$, just that it should be understood as coming from prior hypothetical observations.

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

$$\pi(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{Y}) \propto (2\pi)^{-pN/2} |\boldsymbol{\Sigma}|^{-N/2} \exp\left\{-\frac{1}{2} tr \left[\left(\mathbf{S} + N \left(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}\right) \left(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}\right)^{\top}\right) \boldsymbol{\Sigma}^{-1}\right]\right\}$$

$$\times |\boldsymbol{\Sigma}|^{-1/2} \exp\left\{-\frac{\beta_0}{2} tr \left[\left(\boldsymbol{\mu} - \boldsymbol{\mu}_0\right) (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^{\top} \boldsymbol{\Sigma}^{-1}\right]\right\} |\boldsymbol{\Sigma}|^{-(\alpha_0 + p + 1)/2}$$

$$\times \exp\left\{-\frac{1}{2} tr(\boldsymbol{\Psi}_0 \boldsymbol{\Sigma}^{-1})\right\}.$$

Taking into account that

$$N(\mu - \hat{\mu})(\mu - \hat{\mu})^{\top} + \beta_0(\mu - \mu_0)(\mu - \mu_0)^{\top} = (N + \beta_0)(\mu - \mu_n)(\mu - \mu_n)^{\top} + \frac{N\beta_0}{N + \beta_0}(\hat{\mu} - \mu_0)(\hat{\mu} - \mu_0)^{\top},$$

where $\mu_n=rac{N}{N+eta_0}\hat{m{\mu}}+rac{eta_0}{N+eta_0}m{\mu}_0$ is the posterior mean. We have

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

Then,
$$\boldsymbol{\mu}|\boldsymbol{\Sigma}, \mathbf{Y} \sim N_p\left(\boldsymbol{\mu}_n, \frac{1}{\beta_n}\boldsymbol{\Sigma}\right)$$
, and $\boldsymbol{\Sigma}|\mathbf{Y} \sim IW\left(\boldsymbol{\Psi}_n, \alpha_n\right)$ 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 μ is given by $\int_{\mathcal{S}} \pi(\mu, \Sigma) d\Sigma$ where \mathcal{S} is the space of positive semi-definite matrices. Then,

$$\pi(\boldsymbol{\mu}|\mathbf{Y}) \propto \int_{\mathcal{S}} \left\{ |\mathbf{\Sigma}|^{-(\alpha_n + p + 2)/2} \right\} d\mathbf{\Sigma}$$

$$\exp \left\{ -\frac{1}{2} tr \left[\left(\beta_n \left(\boldsymbol{\mu} - \boldsymbol{\mu}_n \right) \left(\boldsymbol{\mu} - \boldsymbol{\mu}_n \right)^\top + \boldsymbol{\Psi}_n \right) \boldsymbol{\Sigma}^{-1} \right] \right\} d\mathbf{\Sigma}$$

$$\propto \left| \left(\beta_n \left(\boldsymbol{\mu} - \boldsymbol{\mu}_n \right) \left(\boldsymbol{\mu} - \boldsymbol{\mu}_n \right)^\top + \boldsymbol{\Psi}_n \right) \right|^{-(\alpha_n + 1)/2}$$

$$= \left[\left| \boldsymbol{\Psi}_n \right| \times \left| 1 + \beta_n \left(\boldsymbol{\mu} - \boldsymbol{\mu}_n \right)^\top \boldsymbol{\Psi}_n^{-1} \left(\boldsymbol{\mu} - \boldsymbol{\mu}_n \right) \right| \right]^{-(\alpha_n + 1)/2}$$

$$\propto \left(1 + \frac{1}{\alpha_n + 1 - p} \left(\boldsymbol{\mu} - \boldsymbol{\mu}_n \right)^\top \left(\frac{\boldsymbol{\Psi}_n}{(\alpha_n + 1 - p)\beta_n} \right)^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_n) \right)^{-(\alpha_n + 1 - p + p)/2}$$

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.

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 $\boldsymbol{\Sigma}_n = \frac{\Psi_n}{(\alpha_n + 1 - p)\beta_n}$.

The marginal likelihood is given by

$$p(\mathbf{Y}) = \frac{\Gamma_p \left(\frac{v_n}{2}\right)}{\Gamma_n \left(\frac{\alpha_0}{2}\right)} \frac{|\mathbf{\Psi}_0|^{\alpha_0/2}}{|\mathbf{\Psi}_n|^{\alpha_n/2}} \left(\frac{\beta_0}{\beta_n}\right)^{p/2} (2\pi)^{-Np/2},$$

where Γ_p is the multivariate gamma function (see Exercise 5).

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

Example: Tangency portfolio of US tech stocks

The tangency portfolio is the portfolio that maximizes the Sharpe ratio, where this is the excess of return of a portfolio standardized by its risk.

We want to find the shares **w** of a portfolio that maximizes the Sharpe ratio, where $\mu_{i,T+\kappa} = \mathbb{E}(R_{i,T+\kappa} - R_{f,T+\kappa} \mid \mathcal{I}_T)$, $R_{i,T+\kappa}$ and $R_{f,T+\kappa}$ are the returns of stock i and a risk-free asset. Observe that we have the expected value at period $T + \kappa$ of the excess return conditional on information up to $T(\mathcal{I}_T)$, and $\Sigma_{T+\kappa}$ is the covariance of the excess returns, which is a measure of risk. In particular,

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

where the solution is

$$\mathbf{w}^* = rac{\mathbf{\Sigma}_{T+\kappa}^{-1} oldsymbol{\mu}_{T+\kappa}}{\mathbf{1}^{ op} \mathbf{\Sigma}_{T+\kappa}^{-1} oldsymbol{\mu}_{T+\kappa}}.$$

If we want to find the optimal portfolio for the next period under the assumption that the excess of returns follow a multivariate normal distribution, which is a common assumption in these applications, we can set $\kappa=1$, and use the predictive distribution of the excess of returns such that $\mu_{T+1}=\mu_n$ and $\Sigma_{T+1}=\frac{v_n}{v_n-2}(\beta_n+1)\Sigma_n$ given 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).

R code. Optimal tangency portfolio: Tech shares

```
1 library(quantmod)
2 library(xts)
3 library(ggplot2)
4 library(gridExtra)
5 # grid.arrange
6 graphics.off()
7 rm(list=ls())
8 # Data Range
9 sdate <- as.Date("2021-01-01")</pre>
10 edate <- as.Date("2022-09-30")</pre>
11 Date \leftarrow seq(sdate, edate, by = "day")
13 p <- length(tickers)</pre>
# AAPL: Apple, NFLX: Netflix, AMZN: Amazon,
# MSFT: Microsoft, GOOG: Google, META: Meta,
16 # TSLA: Tesla, NVDA: NVIDIA Corporation
17 # INTC: Intel, PYPL: PayPal
_{\rm 18} ss_stock <- getSymbols(tickers, from=sdate, to=edate, auto.
      assign = T)
19 ss_stock <- purrr::map(tickers,function(x) Ad(get(x)))</pre>
20 ss_stock <- as.data.frame(purrr::reduce(ss_stock, merge))</pre>
21 colnames(ss_stock) <- tickers</pre>
22 # This is to get stock prices
23 ss_rtn <- as.data.frame(apply(ss_stock, 2, function(x) {diff
      (log(x), 1)))
24 # Daily returns
25 t10yr <- getSymbols(Symbols = "DGS10", src = "FRED", from=
      sdate, to=edate, auto.assign = F)
_{\rm 26} # To get 10-Year US Treasury yield data from the
27 Federal Reserve Electronic Database (FRED)
28 t10yrd \leftarrow (1 + t10yr/100)^(1/365) -1
29 # Daily returns
30 t10yrd <- t10yrd[row.names(ss_rtn)]</pre>
31 Exc_rtn <- as.matrix(ss_rtn) - kronecker(t(rep(1, p)), as.</pre>
      matrix(t10yrd))
32 # Excesses of return
33 df <- as.data.frame(Exc_rtn)
34 df$Date <- as.Date(rownames(df))</pre>
35 # Get months
36 df $ Month <- months (df $ Date)
37 # Get years
38 df$Year <- format(df$Date, format="%y")
39 # Aggregate on months and year and get mean
40 Data <- sapply(1:p, function(i) {
   aggregate(df[, i] ~ Month + Year, df, mean)})
42 DataExcRtn <- matrix(0, length(Data[, 1] $Month), p)
43 for(i in 1:p){
    DataExcRtn[, i] <- as.numeric(Data[, i]$'df[, i]')</pre>
45 }
46 colnames (DataExcRtn) <- tickers
47 head(DataExcRtn)
```

R code. Optimal tangency portfolio: Tech shares

```
# Hyperparameters #
  N <- dim(DataExcRtn)[1]
  mu0 \leftarrow rep(0, p)
  beta0 <- 1
 Psi0 <- 100 * diag(p)
6 alpha0 <- p + 2
  # Posterior parameters #
  alphan <- N + alpha0
9 vn <- alphan + 1 - p
10 muhat <- colMeans(DataExcRtn)</pre>
nun \leftarrow N/(N + beta0) * muhat + beta0/(N + beta0) * mu0
12 S <- t(DataExcRtn - rep(1, N)%*%t(muhat))%*%(DataExcRtn -</pre>
      rep(1, N) %*%t(muhat))
13 Psin \leftarrow 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</pre>
      %*%mun)))
19 colnames(OptShare) <- tickers</pre>
20 OptShare
    AAPL NFLX
                 AMZN
                         GOOG INTC META MSFT
      PYPL
  -0.019 0.248 0.102 -0.034 0.173 0.23 -0.022 -0.016 0.035
```

We find that the optimal tangency portfolio is composed by 24.8%, 10.2%, 17.3%, 23%, 3.5% and 30.1% weights of Netflix, Amazon, Intel, Meta, NVIDIA and PayPal, and -1.9%, -3.4%, -2.2% and -1.6% weights of Apple, Google, Microsoft and Tesla. A negative weight means being short in financial jargon, that is, borrowing a stock to sell it.

${\bf 4.3 \quad Linear \; regression: The \; conjugate \; normal-normal/inverse \\ \; gamma \; model }$

In this setting we analyze the conjugate normal-normal/inverse gamma model which is the workhorse 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, $y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_K x_{iK} + \mu_i = \mathbf{x}_i^{\top} \beta + \mu_i$ where $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_K)^{\top}$ and $\mu_i \stackrel{iid}{\sim} N(0, \sigma^2)$ is an stochastic error such that $\mathbb{E}[\mu_i | \mathbf{x}_i] = 0$.

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

we can write the model in matrix form: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mu$, where $\mu \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$. This implies that $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$. Then, the likelihood function is

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

The conjugate priors for the parameters are

$$\beta | \sigma^2 \sim N(\beta_0, \sigma^2 \mathbf{B}_0),$$

 $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2).$

Then, the posterior distribution is

$$\pi(\boldsymbol{\beta}, \sigma^{2}|\mathbf{y}, \mathbf{X}) \propto (\sigma^{2})^{-\frac{N}{2}} \exp\left\{-\frac{1}{2\sigma^{2}}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^{\top}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right\}$$

$$\times (\sigma^{2})^{-\frac{K}{2}} \exp\left\{-\frac{1}{2\sigma^{2}}(\boldsymbol{\beta} - \boldsymbol{\beta}_{0})^{\top}\mathbf{B}_{0}^{-1}(\boldsymbol{\beta} - \boldsymbol{\beta}_{0})\right\}$$

$$\times \frac{(\delta_{0}/2)^{\alpha_{0}/2}}{\Gamma(\alpha_{0}/2)} \left(\frac{1}{\sigma^{2}}\right)^{\alpha_{0}/2+1} \exp\left\{-\frac{\delta_{0}}{2\sigma^{2}}\right\}$$

$$\propto (\sigma^{2})^{-\frac{K}{2}} \exp\left\{-\frac{1}{2\sigma^{2}}[\boldsymbol{\beta}^{\top}(\mathbf{B}_{0}^{-1} + \mathbf{X}^{\top}\mathbf{X})\boldsymbol{\beta} - 2\boldsymbol{\beta}^{\top}(\mathbf{B}_{0}^{-1}\boldsymbol{\beta}_{0} + \mathbf{X}^{\top}\mathbf{X}\hat{\boldsymbol{\beta}})]\right\}$$

$$\times \left(\frac{1}{\sigma^{2}}\right)^{(\alpha_{0}+N)/2+1} \exp\left\{-\frac{\delta_{0} + \mathbf{y}^{\top}\mathbf{y} + \boldsymbol{\beta}_{0}^{\top}\mathbf{B}_{0}^{-1}\boldsymbol{\beta}_{0}}{2\sigma^{2}}\right\},$$

where $\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}\mathbf{y}$ is the maximum likelihood estimator. Adding and subtracting $\boldsymbol{\beta}_n^{\top}\mathbf{B}_n^{-1}\boldsymbol{\beta}_n$ to complete the square, where $\mathbf{B}_n = (\mathbf{B}_0^{-1} + \mathbf{X}^{\top}\mathbf{X})^{-1}$ and $\boldsymbol{\beta}_n = \mathbf{B}_n(\mathbf{B}_0^{-1}\boldsymbol{\beta}_0 + \mathbf{X}^{\top}\mathbf{X}\hat{\boldsymbol{\beta}})$,

$$\pi(\boldsymbol{\beta}, \sigma^{2} | \mathbf{y}, \mathbf{X}) \propto \underbrace{(\sigma^{2})^{-\frac{K}{2}} \exp\left\{-\frac{1}{2\sigma^{2}} (\boldsymbol{\beta} - \boldsymbol{\beta}_{n})^{\top} \mathbf{B}_{n}^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_{n})\right\}}_{1} \times \underbrace{(\sigma^{2})^{-\left(\frac{\alpha_{n}}{2} + 1\right)} \exp\left\{-\frac{\delta_{n}}{2\sigma^{2}}\right\}}_{2}.$$

The first expression is the kernel of a normal density function, $\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \sigma^2 \mathbf{B}_n)$. 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 $\delta_n = \delta_0 + \mathbf{y}^{\mathsf{T}}\mathbf{y} + \boldsymbol{\beta}_0^{\mathsf{T}} \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 - \boldsymbol{\beta}_n^{\mathsf{T}} \mathbf{B}_n^{-1} \boldsymbol{\beta}_n$.

Taking into account that

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

where $(\mathbf{B}_0^{-1} + \mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{B}_0^{-1} = \mathbf{I}_{\mathbf{K}} - (\mathbf{B}_0^{-1} + \mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}\mathbf{X}$ [97]. Setting $\mathbf{W} = (\mathbf{B}_0^{-1} + \mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}\mathbf{X}$ we have $\boldsymbol{\beta}_n = (\mathbf{I}_{\mathbf{K}} - \mathbf{W})\boldsymbol{\beta}_0 + \mathbf{W}\boldsymbol{\hat{\beta}}$, that is, the posterior mean of $\boldsymbol{\beta}$ 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 $\mathbf{B}_0^{-1} \to \mathbf{0}_{\mathbf{K}}$, we obtain $\mathbf{W} \to I_K$, such that $\boldsymbol{\beta}_n \to \hat{\boldsymbol{\beta}}$, 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.⁴ 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.

Now, we calculate the posterior marginal distribution of β ,

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

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

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

⁴A particular case of the Woodbury matrix identity.

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 σ^2 , the integral is equal to 1, and therefore

$$\pi(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) = \frac{\Gamma((\alpha_n + K)/2)}{(s/2)^{(\alpha_n + K)/2}}$$

$$\propto s^{-(\alpha_n + K)/2}$$

$$= [\delta_n + (\boldsymbol{\beta} - \boldsymbol{\beta}_n)^{\top} \mathbf{B}_n^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_n)]^{-(\alpha_n + K)/2}$$

$$= \left[1 + \frac{(\boldsymbol{\beta} - \boldsymbol{\beta}_n)^{\top} \left(\frac{\delta_n}{\alpha_n} \mathbf{B}_n\right)^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_n)}{\alpha_n}\right]^{-(\alpha_n + K)/2}$$

$$\propto \left[1 + \frac{(\boldsymbol{\beta} - \boldsymbol{\beta}_n)^{\top} \mathbf{H}_n^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_n)}{\alpha_n}\right]^{-(\alpha_n + K)/2},$$

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

Observe that as we have incorporated the uncertainty of the variance, the posterior for β 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

$$\begin{split} p(\mathbf{y}) &= \int_{0}^{\infty} \int_{R^{K}} \pi(\beta | \sigma^{2}) \pi(\sigma^{2}) p(\mathbf{y} | \beta, \sigma^{2}, \mathbf{X}) d\sigma^{2} d\beta \\ &= \int_{0}^{\infty} \pi(\sigma^{2}) \frac{1}{(2\pi\sigma^{2})^{N/2}} \exp\left\{-\frac{1}{2\sigma^{2}}m\right\} \frac{1}{(2\pi\sigma^{2})^{K/2} |\mathbf{B}_{0}|^{1/2}} \\ &\times \int_{R^{K}} \exp\left\{-\frac{1}{2\sigma^{2}} (\beta - \beta_{n})^{\top} \mathbf{B}_{n}^{-1} (\beta - \beta_{n})\right\} d\sigma^{2} d\beta \\ &= \int_{0}^{\infty} \pi(\sigma^{2}) \frac{1}{(2\pi\sigma^{2})^{N/2}} \exp\left\{-\frac{1}{2\sigma^{2}}m\right\} \frac{|\mathbf{B}_{n}|^{1/2}}{|\mathbf{B}_{0}|^{1/2}} d\sigma^{2} \\ &= \int_{0}^{\infty} \frac{(\delta_{0}/2)^{\alpha_{0}/2}}{\Gamma(\alpha_{0}/2)} \left(\frac{1}{\sigma^{2}}\right)^{\alpha_{0}/2+1} \exp\left\{\left(-\frac{\delta_{0}}{2\sigma^{2}}\right)\right\} \frac{1}{(2\pi\sigma^{2})^{N/2}} \exp\left\{-\frac{1}{2\sigma^{2}}m\right\} \frac{|\mathbf{B}_{n}|^{1/2}}{|\mathbf{B}_{0}|^{1/2}} d\sigma^{2} \\ &= \frac{1}{(2\pi)^{N/2}} \frac{(\delta_{0}/2)^{\alpha_{0}/2}}{\Gamma(\alpha_{0}/2)} \frac{|\mathbf{B}_{n}|^{1/2}}{|\mathbf{B}_{0}|^{1/2}} \int_{0}^{\infty} \left(\frac{1}{\sigma^{2}}\right)^{\frac{\alpha_{0}+N}{2}+1} \exp\left\{\left(-\frac{\delta_{0}+m}{2\sigma^{2}}\right)\right\} d\sigma^{2} \\ &= \frac{1}{\pi^{N/2}} \frac{\delta_{0}^{\alpha_{0}/2}}{\delta_{n}^{\alpha_{0}/2}} \frac{|\mathbf{B}_{n}|^{1/2}}{|\mathbf{B}_{0}|^{1/2}} \frac{\Gamma(\alpha_{n}/2)}{\Gamma(\alpha_{0}/2)}. \end{split}$$

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

$$BF_{12} = \frac{p(\mathbf{y}|\mathcal{M}_{1})}{p(\mathbf{y}|\mathcal{M}_{2})}$$

$$= \frac{\frac{\delta_{10}^{\alpha_{10}/2}}{\delta_{1n}^{\alpha_{1n}/2}} \frac{|\mathbf{B}_{1n}|^{1/2}}{|\mathbf{B}_{10}|^{1/2}} \frac{\Gamma(\alpha_{1n}/2)}{\Gamma(\alpha_{10}/2)}}{\frac{\delta_{20}^{\alpha_{20}/2}}{\delta_{2n}^{\alpha_{2n}/2}} \frac{|\mathbf{B}_{2n}|^{1/2}}{|\mathbf{B}_{20}|^{1/2}} \frac{\Gamma(\alpha_{2n}/2)}{\Gamma(\alpha_{20}/2)}},$$

Observe that *ceteris paribus*, the model having better fit, coherence be-

where subscripts 1 and 2 refer to each model, respectively.

tween sample and prior information regarding location parameters, higher prior to posterior precision and less parameters is favored by the Bayes factor. Observe that the Bayes factor rewards model fit as the sum of squared errors is in δ_n , the better fit (lower sum of squared errors), the better the Bayes factor. In addition, a weighted distance between sample and prior location parameters also appears in δ_n , the greater this distance, the worse is model support. The ratio of determinants between posterior and prior covariance matrices is also present, the higher this ratio, the better for the Bayes factor supporting a model due to information gains. To see the effect of model's parsimony, let's take the common situation in applications where $\mathbf{B}_{j0} = c\mathbf{I}_{K_j}$ then $|\mathbf{B}_{j0}| = c^{K_j}$ such that $\left(\frac{|\mathbf{B}_{20}|}{|\mathbf{B}_{10}|}\right)^{1/2} = \left(\frac{c^{K_2/2}}{c^{K_1/2}}\right)$, if $K_2/K_1 > 1$ and $c \to \infty$, the latter implying a non-informative prior, then $BF_{12} \to \infty$, this means infinite evidence supporting the parsimonious model no matter what sample information says. Comparing models having the same number of regressors $(K_1 = K_2)$ is not a safe ground as $|\mathbf{B}_0|$ depending on measure units of the regressors such that conclusions regarding model selection depending on this, which is not a nice property. This prevents against using non-informative priors when performing model selection in the Bayesian framework. Observe that this is not the case when $\alpha_0 \to 0$ and $\delta_0 \to 0$, which implies a non-informative prior for the variance parameter.⁵ We observe here that $\Gamma(\alpha_{j0})$ cancels out, $\alpha_{jn} \to N$ and $\delta_{jn} \to (\mathbf{y} - \mathbf{X}_j \hat{\boldsymbol{\beta}}_j)^\top (\mathbf{y} - \mathbf{X}_j \hat{\boldsymbol{\beta}}_j) + (\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_{j0})^\top ((\mathbf{X}_j^\top \mathbf{X}_j)^{-1} + \mathbf{B}_{j0})^{-1} (\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_{j0}),$ therefore there is not effect. This is due to σ^2 being a common parameter in both models. In general, we can use non-informative priors for common parameters to all models, but we cannot use non-informative priors for non-common

parameters when performing model selection using the Bayes factor.

The posterior predictive is equal to

⁵[38] prevents against this common practice.

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

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,

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

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

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

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

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

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 $\mathbf{C} = \mathbf{I}_{N_0} - \mathbf{X}_0 (\mathbf{B}_n^{-1} + \mathbf{X}_0 \mathbf{B}_n^{-1})$

$$\begin{split} \mathbf{X}_{0}^{\top}\mathbf{X}_{0})^{-1}\mathbf{X}_{0}^{\top} &= \mathbf{I}_{N_{0}} - \mathbf{X}_{0}\mathbf{M}^{-1}\mathbf{X}_{0}^{\top},^{6} \text{ and } \boldsymbol{\beta}_{**} = \mathbf{C}^{-1}\mathbf{X}_{0}\mathbf{M}^{-1}\mathbf{B}_{n}^{-1}\boldsymbol{\beta}_{n}, \text{ then} \\ \boldsymbol{\beta}_{n}^{\top}\mathbf{B}_{n}^{-1}\boldsymbol{\beta}_{n} + \mathbf{Y}_{0}^{\top}\mathbf{Y}_{0} - \boldsymbol{\beta}_{*}^{\top}\mathbf{M}\boldsymbol{\beta}_{*} &= \boldsymbol{\beta}_{n}^{\top}\mathbf{B}_{n}^{-1}\boldsymbol{\beta}_{n} + \mathbf{Y}_{0}^{\top}\mathbf{Y}_{0} - (\boldsymbol{\beta}_{n}^{\top}\mathbf{B}_{n}^{-1} + \mathbf{Y}_{0}^{\top}\mathbf{X}_{0})\mathbf{M}^{-1}(\mathbf{B}_{n}^{-1}\boldsymbol{\beta}_{n} + \mathbf{X}_{0}^{\top}\mathbf{Y}_{0}) \\ &= \boldsymbol{\beta}_{n}^{\top}(\mathbf{B}_{n}^{-1} - \mathbf{B}_{n}^{-1}\mathbf{M}^{-1}\mathbf{B}_{n}^{-1})\boldsymbol{\beta}_{n} + \mathbf{Y}_{0}^{\top}\mathbf{C}\mathbf{Y}_{0} \\ &- 2\mathbf{Y}_{0}^{\top}\mathbf{C}\mathbf{C}^{-1}\mathbf{X}_{0}\mathbf{M}^{-1}\mathbf{B}_{n}^{-1}\boldsymbol{\beta}_{n} + \boldsymbol{\beta}_{**}^{\top}\mathbf{C}\boldsymbol{\beta}_{**} - \boldsymbol{\beta}_{**}^{\top}\mathbf{C}\boldsymbol{\beta}_{**} \\ &= \boldsymbol{\beta}_{n}^{\top}(\mathbf{B}_{n}^{-1} - \mathbf{B}_{n}^{-1}\mathbf{M}^{-1}\mathbf{B}_{n}^{-1})\boldsymbol{\beta}_{n} + (\mathbf{Y}_{0} - \boldsymbol{\beta}_{**})^{\top}\mathbf{C}(\mathbf{Y}_{0} - \boldsymbol{\beta}_{**}) \\ &- \boldsymbol{\beta}_{**}^{\top}\mathbf{C}\boldsymbol{\beta}_{**}, \end{split}$$

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,

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

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

Example: Demand of electricity

We study in this example the determinants of monthly demand of electricity by Colombian households. There is information of 2103 households, particularly, average price (USD/kWh), indicators of socioeconomic conditions of the neighborhood where the household is located (IndSocio1 is the lowest and IndSocio3 is the highest), an indicator if the household is located in a municipality that is above 1000 meters above the sea level, the number of rooms in the house, the number of members of the households, presence of children in the household (1 is yes), and monthly income (USD). The specification is

log(Electricity_i) =
$$\beta_1$$
 log(price_i) + β_2 IndSocio1_i + β_3 IndSocio2_i + β_4 Altitude_i
+ β_5 Nrooms_i + β_6 HouseholdMem_i + β_7 Children_i
+ β_8 log(Income_i) + β_9 + μ .

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) is that the posterior mean of the own-price of electricity demand is -1.09, and the 95% symmetric credible interval is (-1.47, -0.71). Households in neighborhoods of low socioeconomic conditions and located in municipalities 1000 meters above the sea level consume less electricity, 32.7% and 19.7% on average, respectively. An additional

$$6$$
Using $(A + BDC)^{-1} = A^{-1} - A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1}$

room implies a 8.7% increase in electricity consumption, and another household member increases consumption in 5.9% on average. The income elasticity mean estimate is 0.074, which means that 10% increase of income increases 0.74% 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,

```
log(Electricity<sub>i</sub>) = \beta_1IndSocio1<sub>i</sub> + \beta_2IndSocio2<sub>i</sub> + \beta_3Altitude<sub>i</sub> + \beta_4Nrooms<sub>i</sub> + \beta_5HouseholdMem<sub>i</sub> + \beta_6Children<sub>i</sub> + \beta_7 log(Income<sub>i</sub>) + \beta_8 + \mu.
```

In particular, we check what happens as c_0 increases from 10^0 to 10^{20} . We see that when $c_0 = 1$, $BF_{12} = 8.68 \times 10^{+16}$, which means very strong evidence in favor of the model including the price of electricity. However, as c_0 increases, the Bayes factor decreases, which means evidence supporting the model 2, for instance, $BF_{12} = 3.11 \times 10^{-4}$ when $c_0 = 10^{20}$. This is an example of the problem of using non-informative priors to calculate the Bayes factor; there is very strong evidence to support the parsimonious model when $c_0 \to \infty$.

We can get the posterior predictive distribution of the monthly electricity demand 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, and an electricity price equal to USD/kWh 0.15. Figure 4.1 shows the histogram of the predictive posterior distribution, the highest posterior density credible interval at 95% is between kWh 44.4 and kWh 373.9, and the posterior mean is kWh 169.4.

R code. Demand of electricity, posterior predictive distribution

```
1 rm(list = ls())
2 set.seed(010101)
3 # Electricity demand
4 DataUt <- read.csv("DataApplications/Utilities.csv", sep = "
       ,", header = TRUE, fileEncoding = "latin1")
5 DataUtEst <- DataUt %>%
    filter(Electricity != 0)
7 attach(DataUtEst)
8 # Dependent variable: Monthly consumption (kWh) in log
9 Y <- log(Electricity)
10 # Regressors quantity including intercept
11 X <- cbind(LnPriceElect, IndSocio1, IndSocio2, Altitude,
Nrooms, HouseholdMem, Children, Lnincome, 1) _{\rm 12} # LnPriceElect: Price per kWh (USD) in log
13 # IndSocio1, IndSocio2, IndSocio3: Indicators socio-economic
        condition (1) is the lowest and (3) the highest
_{14} # Altitude: Indicator of household location (1 is more than
      1000 meters above sea level)
15 # Nrooms: Number of rooms in house
16 # HouseholdMem: Number of household members
17 # Children: Indicator por presence of children in household
       (1)
18 # Lnincome: Monthly income (USD) in log
19 k <- dim(X)[2]
20 N <- dim(X)[1]
21 # Hyperparameters
22 d0 <- 0.001/2
23 a0 <- 0.001/2
24 b0 <- rep(0, k)
25 B0 <- 1000*diag(k)
26 # Posterior parameters
27 bhat <- solve(t(X)%*%X)%*%t(X)%*%Y
28 Bn <- as.matrix(Matrix::forceSymmetric(solve(solve(B0) + t(X
      )%*\mbox{\ensuremath{\%}X))) # Force this matrix to be symmetric
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)%
      *%solve(Bn)%*%bn)
31 an <- a0 + N
32 Hn <- Bn*dn/an
33 # Posterior draws
_{\rm 34} S <- 10000 # Number of draws from posterior distributions
35 sig2 <- MCMCpack::rinvgamma(S,an/2,dn/2)</pre>
36 summary (coda::mcmc(sig2))
```

R code. Demand of electricity, posterior distribution

```
1 Iterations = 1:10000
2 Thinning interval = 1
3 Number of chains = 1
4 Sample size per chain = 10000
6 1. Empirical mean and standard deviation for each
7 variable, plus standard error of the mean:
                   SD
                            Naive SE Time-series SE
10 2.361e-01
                 7.617e-03 7.617e-05 7.617e-05
12 2. Quantiles for each variable:
13
14 2.5%
          25%
                50%
                       75% 97.5%
15 0.2217 0.2309 0.2360 0.2412 0.2513
17 Betas <- LaplacesDemon::rmvt(S, bn, Hn, an)
18 summary(coda::mcmc(Betas))
19 Iterations = 1:10000
20 Thinning interval = 1
21 Number of chains = 1
22 Sample size per chain = 10000
24 1. Empirical mean and standard deviation for each
25 variable, plus standard error of the mean:
                          SD
                                Naive SE Time-series SE
                Mean
28 LnPriceElect -1.09043 0.19459 0.0019459 0.0019459
               -0.32783 0.05294 0.0005294 0.0005294
29 IndSocio1
30 IndSocio2
               -0.05737 0.04557 0.0004557 0.0004557
31 Altitude
               -0.19780 0.02386 0.0002386 0.0002429
                0.08731 0.01094 0.0001094 0.0001119
32 Nrooms
33 HouseholdMem 0.05987 0.01334 0.0001334 0.0001334
                0.05696 0.03043 0.0003043 0.0003043
34 Children
                0.07447 0.01223 0.0001223 0.0001223
35 Lnincome
                2.52296 0.35077 0.0035077 0.0035077
38 2. Quantiles for each variable:
                    2.5%
                              25%
                                       50%
                                               75%
                                                       97.5%
41 LnPriceElect -1.472069 -1.22432 -1.08961 -0.95703 -0.71429
               42 IndSocio1
43 IndSocio2
               -0.244759 -0.21372 -0.19783 -0.18164 -0.15094
44 Altitude
45 Nrooms
                0.066432 0.07985 0.08709 0.09480 0.10864
46 HouseholdMem 0.033623 0.05089 0.05975
                                           0.06889 0.08596
47 Children
               -0.002259
                         0.03637
                                  0.05698
                                           0.07736
                                                    0.11681
                0.050536 0.06614
48 Lnincome
                                  0.07449
                                           0.08283
                                                    0.09852
                1.835507 2.28703 2.52165 2.76364 3.21199
49
50
```

R code. Demand of electricity, Bayes factor

```
1 # Log marginal function (multiply by -1 due to minimization)
2 LogMarLikLM <- function(X, c0){</pre>
    k \leftarrow dim(X)[2]
    N <- dim(X)[1]
    # Hyperparameters
    B0 \leftarrow c0*diag(k)
    b0 \leftarrow rep(0, k)
     # Posterior parameters
    bhat <- solve(t(X)%*%X)%*%t(X)%*%Y
    # Force this matrix to be symmetric
    Bn <- as.matrix(Matrix::forceSymmetric(solve(solve(B0) + t</pre>
      (X)%*%X)))
     bn <- Bn\%*\%(solve(B0)\%*\%b0 + t(X)\%*\%X\%*\%bhat)
    dn \leftarrow as.numeric(d0 + t(Y))%*%Y+t(b0)%*%solve(B0)%*%b0-t(bn)
      )%*%solve(Bn)%*%bn)
    an \leftarrow a0 + N
14
     # Log marginal likelihood
15
    logpy \leftarrow (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)
    return(-logpy)
18 }
19 cs <- c(10<sup>0</sup>, 10<sup>3</sup>, 10<sup>6</sup>, 10<sup>10</sup>, 10<sup>12</sup>, 10<sup>15</sup>, 10<sup>20</sup>)
20 # Observe -1 to recover the right sign
21 LogML <- sapply(cs, function(c) {-LogMarLikLM(c0=c, X = X)})</pre>
22 # Regressor without price
_{\rm 23} Xnew <- cbind(IndSocio1, IndSocio2, Altitude, Nrooms,
       HouseholdMem, Children, Lnincome, 1)
24 # Observe -1 to recover the right sign
25 LogMLnew <- sapply(cs, function(c) {-LogMarLikLM(c0=c,X =</pre>
       Xnew)})
26 # Bayes factor
27 BF <- exp(LogML - LogMLnew)
28 BF
29 8.687567e+16 1.006679e+05 3.108415e+03 3.108340e+01 3.108343
       e+00 9.829443e-02 3.108343e-04
30 # Empirical Bayes: Obtain c0 maximizing the log
31 marginal likelihood
32 c0 <- c0
33 EB <- optim(c0, fn = LogMarLikLM, method = "Brent", lower =
       0.0001, upper = 10^6, X = X)
34 EB$par
35 3.254822
36 EB$value
37 1404.108
38 EBnew <- optim(c0, fn = LogMarLikLM, method = "Brent", lower
        = 0.0001, upper = 10^6, X = Xnew)
39 EBnew $par
40 10.00597
41 EBnew$value
42 1422.199
43 # Change of order to take into account the -1 in the
       LogMarLikLM function
44 BFEM <- exp(EBnew$value - EB$value)
45 BFEM
46 71897938
```

$R\ code.\ Demand\ of\ electricity,\ predictive\ distribution$

```
Predictive distribution
  Xpred \leftarrow c(log(0.15), 1, 0, 0, 2, 3, 1, log(500), 1)
  Mean <- Xpred%*%bn
4 Hn <- dn*(1+t(Xpred)%*%Bn%*%Xpred)/an
5 ExpKwH <- exp(LaplacesDemon::rmvt(S, Mean, Hn, an))</pre>
6 summary (ExpKwH)
            24.06
  1st Qu.: 121.70
  Median : 169.37
          : 189.60
10 Mean
11 3rd Qu.: 234.19
12 Max.
          :1243.68
13 HDI <- HDInterval::hdi(ExpKwH, credMass = 0.95) # Highest
      posterior density credible interval
14 HDI
15 lower 44.40203
16 upper 373.86494
17 hist(ExpKwH, main = "Histogram: Monthly demand of
      electricity", xlab = "Monthly kWh", col = "blue", breaks
```

Histogram: Monthly demand of electricity



FIGURE 4.1
Histogram using the posterior predictive distribution of electricity demand

4.4 Multivariate linear regression: The conjugate normalnormal/inverse Wishart model

Let's study the multivariate regression setting where there are N-dimensional vectors \mathbf{y}_m , $m=1,2,\ldots,M$ such that $\mathbf{y}_m=\mathbf{X}\boldsymbol{\beta}_m+\mu_m$, \mathbf{X} is the set of common regressors, and μ_m is the N-dimensional vector of stochastic errors for each equation such that $\mathbf{U}=[\mu_1 \ \mu_2 \ \ldots \ \mu_M] \sim MN_{N,M}(\mathbf{0},\mathbf{I}_N,\boldsymbol{\Sigma})$, that is, a matrix variate normal distribution where $\boldsymbol{\Sigma}$ is the covariance matrix of each i-th row of \mathbf{U} , $i=1,2,\ldots,N$, and we are assuming independence between the rows. Then, $vec(\mathbf{U}) \sim N_{N\times M}(\mathbf{0},\boldsymbol{\Sigma}\otimes\mathbf{I_N})$.

This framework can be written in matrix form

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

Therefore, $\mathbf{Y} \sim N_{N \times M}(\mathbf{XB}, \mathbf{\Sigma} \otimes \mathbf{I_N})^8$

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

where $\mathbf{S} = (\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}})^{\top}(\mathbf{Y} - \mathbf{X}\widehat{\mathbf{B}})$, $\widehat{\mathbf{B}} = (\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}\mathbf{Y}$ (see Exercise 9). The conjugate prior for this models is $\pi(\mathbf{B}, \mathbf{\Sigma}) = \pi(\mathbf{B}|\mathbf{\Sigma})\pi(\mathbf{\Sigma})$ where $\pi(\mathbf{B}|\mathbf{\Sigma}) \sim N_{K \times M}(\mathbf{B}_0, \mathbf{V}_0, \mathbf{\Sigma})$ and $\pi(\mathbf{\Sigma}) \sim IW(\mathbf{\Psi}_0, \alpha_0)$, that is,

 $^{^7}vec$ denotes the vectorization operation, and \otimes denotes the kronecker product.

⁸We can write down the former expression in a more familiar way using vectorization properties, $\underbrace{vec(Y)}_{\mathbf{y}} = \underbrace{(\mathbf{I}_M \otimes \mathbf{X})}_{\mathbf{Z}} \underbrace{vec(\mathbf{B})}_{\boldsymbol{\beta}} + \underbrace{vec(\mathbf{U})}_{\boldsymbol{\mu}}, \text{ where } \mathbf{y} \sim N_{N \times M}(\mathbf{Z}\boldsymbol{\beta}, \boldsymbol{\Sigma} \otimes \mathbf{I_N}).$

$$\pi(\mathbf{B}, \mathbf{\Sigma}) \propto |\mathbf{\Sigma}|^{-K/2} \exp\left\{-\frac{1}{2} tr \left[(\mathbf{B} - \mathbf{B}_0)^{\top} \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) \mathbf{\Sigma}^{-1} \right] \right\}$$
$$\times |\mathbf{\Sigma}|^{-(\alpha_0 + M + 1)/2} \exp\left\{-\frac{1}{2} tr \left[\mathbf{\Psi}_0 \mathbf{\Sigma}^{-1} \right] \right\}.$$

The posterior distribution is given by

$$\begin{split} \pi(\mathbf{B}, \mathbf{\Sigma} | \mathbf{Y}, \mathbf{X}) &\propto p(\mathbf{Y} | \mathbf{B}, \mathbf{\Sigma}, \mathbf{X}) \pi(\mathbf{B} | \mathbf{\Sigma}) \pi(\mathbf{\Sigma}) \\ &\propto |\mathbf{\Sigma}|^{-\frac{N+K+\alpha_0+M+1}{2}} \\ &\times \exp\left\{-\frac{1}{2} tr\left[(\mathbf{\Psi_0} + \mathbf{S} + (\mathbf{B} - \mathbf{B_0})^{\top} \mathbf{V_0^{-1}} (\mathbf{B} - \mathbf{B_0}) \right.\right. \\ &\left. + (\mathbf{B} - \widehat{\mathbf{B}})^{\top} \mathbf{X}^{\top} \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}})) \mathbf{\Sigma}^{-1}\right]\right\}. \end{split}$$

Completing the squares on ${\bf B}$ and collecting the remaining terms in the bracket yields

$$\mathbf{\Psi}_0 + \mathbf{S} + (\mathbf{B} - \mathbf{B}_0)^{\top} \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) + (\mathbf{B} - \widehat{\mathbf{B}})^{\top} \mathbf{X}^{\top} \mathbf{X} (\mathbf{B} - \widehat{\mathbf{B}}) = (\mathbf{B} - \mathbf{B}_n)^{\top} \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) + \mathbf{\Psi}_n,$$
where

$$\mathbf{B}_{n} = (\mathbf{V}_{0}^{-1} + \mathbf{X}^{\top} \mathbf{X})^{-1} (\mathbf{V}_{0}^{-1} \mathbf{B}_{0} + \mathbf{X}^{\top} \mathbf{Y}) = (\mathbf{V}_{0}^{-1} + \mathbf{X}^{\top} \mathbf{X})^{-1} (\mathbf{V}_{0}^{-1} \mathbf{B}_{0} + \mathbf{X}^{\top} \mathbf{X} \widehat{\mathbf{B}}),$$

$$\mathbf{V}_{n} = (\mathbf{V}_{0}^{-1} + \mathbf{X}^{\top} \mathbf{X})^{-1},$$

Thus, the posterior distribution can be written as

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

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

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

$$\pi(\mathbf{B}|\mathbf{Y},\mathbf{X}) \propto \int_{\mathcal{S}} |\mathbf{\Sigma}|^{-(\alpha_n + K + M + 1)/2}$$

$$\times \exp\left\{-\frac{1}{2}tr\left\{\left[(\mathbf{B} - \mathbf{B}_n)^{\top}\mathbf{V}_n^{-1}(\mathbf{B} - \mathbf{B}_n) + \mathbf{\Psi}_n\right]\mathbf{\Sigma}^{-1}\right\}\right\} d\mathbf{\Sigma}$$

$$\propto |(\mathbf{B} - \mathbf{B}_n)^{\top}\mathbf{V}_n^{-1}(\mathbf{B} - \mathbf{B}_n) + \mathbf{\Psi}_n|^{-(K + \alpha_n)/2}$$

$$= \left[|\mathbf{\Psi}_n| \times |\mathbf{I}_K + \mathbf{V}_n^{-1}(\mathbf{B} - \mathbf{B}_n)\mathbf{\Psi}_n^{-1}(\mathbf{B} - \mathbf{B}_n)^{\top}|\right]^{-(\alpha_n + 1 - M + K + M - 1)/2}$$

$$\propto |\mathbf{I}_K + \mathbf{V}_n^{-1}(\mathbf{B} - \mathbf{B}_n)\mathbf{\Psi}_n^{-1}(\mathbf{B} - \mathbf{B}_n)^{\top}|^{-(\alpha_n + 1 - M + K + M - 1)/2}.$$

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

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

The marginal likelihood is the following,

$$\begin{split} p(\mathbf{Y}) &= \int_{\mathcal{B}} \int_{\mathcal{S}} \left\{ (2\pi)^{-NM/2} |\mathbf{\Sigma}|^{-N/2} \exp\left\{ -\frac{1}{2} tr \left[\mathbf{S} + (\mathbf{B} - \hat{\mathbf{B}})^{\top} \mathbf{X}^{\top} \mathbf{X} (\mathbf{B} - \hat{\mathbf{B}}) \right] \mathbf{\Sigma}^{-1} \right\} \\ &\times (2\pi)^{-KM/2} |\mathbf{V}_0|^{-M/2} |\mathbf{\Sigma}|^{-K/2} \exp\left\{ -\frac{1}{2} tr \left[(\mathbf{B} - \mathbf{B}_0)^{\top} \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) \mathbf{\Sigma}^{-1} \right] \right\} \\ &\times \frac{|\Psi_0|^{\alpha_0/2}}{2^{\alpha_0M/2} \Gamma_M(\alpha_0/2)} |\mathbf{\Sigma}|^{-(\alpha_0+M+1)/2} \exp\left\{ -\frac{1}{2} tr \left[\Psi_0 \mathbf{\Sigma}^{-1} \right] \right\} \right\} d\mathbf{\Sigma} d\mathbf{B} \\ &= (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\Psi_0|^{\alpha_0/2}}{2^{\alpha_0M/2} \Gamma_M(\alpha_0/2)} \\ &\times \int_{\mathcal{B}} \int_{\mathcal{S}} \left\{ |\mathbf{\Sigma}|^{-(\alpha_0+N+K+M+1)/2} \exp\left\{ -\frac{1}{2} tr \left[\mathbf{S} + (\mathbf{B} - \hat{\mathbf{B}})^{\top} \mathbf{X}^{\top} \mathbf{X} (\mathbf{B} - \hat{\mathbf{B}}) + (\mathbf{B} - \mathbf{B}_0)^{\top} \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) + \Psi_0 \right] \mathbf{\Sigma}^{-1} \right\} \right\} d\mathbf{\Sigma} d\mathbf{B} \\ &= (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\Psi_0|^{\alpha_0/2}}{2^{\alpha_0M/2} \Gamma_M(\alpha_0/2)} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2) \\ &\times \int_{\mathcal{B}} \left| \mathbf{S} + (\mathbf{B} - \hat{\mathbf{B}})^{\top} \mathbf{X}^{\top} \mathbf{X} (\mathbf{B} - \hat{\mathbf{B}}) + (\mathbf{B} - \mathbf{B}_0)^{\top} \mathbf{V}_0^{-1} (\mathbf{B} - \mathbf{B}_0) + \Psi_0 \right]^{-(\alpha_n+K)/2} d\mathbf{B} \\ &= (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\Psi_0|^{\alpha_0/2}}{2^{\alpha_0M/2} \Gamma_M(\alpha_0/2)} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2) \\ &\times \int_{\mathcal{B}} \left| (\mathbf{B} - \hat{\mathbf{B}}_n)^{\top} \mathbf{V}_n^{-1} (\mathbf{B} - \hat{\mathbf{B}}_n) + \Psi_n \right|^{-(\alpha_n+K)/2} d\mathbf{B} \\ &= (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\Psi_0|^{\alpha_0/2}}{2^{\alpha_0M/2} \Gamma_M(\alpha_0/2)} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2) \\ &\times \int_{\mathcal{B}} \left[|\Psi_n| \times |\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+K)/2} d\mathbf{B} \\ &= |\Psi_n|^{-(\alpha_n+K)/2} (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\Psi_0|^{\alpha_0/2} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_n+K)/2)}{2^{\alpha_0M/2} \Gamma_M(\alpha_0/2)} \\ &\times \int_{\mathbf{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+K)/2} d\mathbf{B} \\ &= |\Psi_n|^{-(\alpha_n+K)/2} (2\pi)^{-M(N+K)/2} |\mathbf{V}_0|^{-M/2} \frac{|\Psi_0|^{\alpha_0/2} 2^{M(\alpha_n+K)/2} \Gamma_M((\alpha_0/2))}{2^{\alpha_0M/2} \Gamma_M(\alpha_0/2)} \\ &\times \pi^{MK/2} \frac{\Gamma_M((\alpha_n+1-M+K+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_0/2)}{\Gamma_M(\alpha_0/2)} \pi^{-MN/2}. \end{aligned}$$

The third equality follows from having the kernel of a inverse Wishart distribution, the fifth from the Silvester's theorem, and the seventh from having the kernel of a matrix t distribution.

Summary 91

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

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

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 $T_{N_0,M}(\alpha_n - M + 1, \mathbf{X}_0\mathbf{B}_n, \mathbf{I}_{N_0} + \mathbf{X}_0\mathbf{V}_n\mathbf{X}_0^\top, \mathbf{\Psi}_n)$ (see Exercise 6). Observe that the prediction is centered at $\mathbf{X}_0\mathbf{B}_n$, and the covariance matrix of $vec(\mathbf{Y}_0)$ is $\frac{(\mathbf{I}_{N_0} + \mathbf{X}_0\mathbf{V}_n\mathbf{X}_0^\top)\otimes \mathbf{\Psi}_n}{\alpha_n - M - 1}$.

4.5 Summary

We introduce the conjugate family models for discrete and continuous data. These models are the basic Bayesian framework to due its mathematical tractability as we get closed-form expressions for the posterior distributions, the marginal likelihood, and the predictive distribution. We also present the Bayesian linear univariate and multivariate regression frameworks under conjugate families. This is the cornerstone to perform regression analysis in the Bayesian setting.

4.6 Exercises

- 1. Write in the canonical form the distribution of the Bernoulli example, 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 having known size n_i and same unknown probability θ . Show that $p(\mathbf{y}|\theta)$ is in the exponential family, and 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 a *exponential distribution*. Show that $p(\mathbf{y}|\lambda)$ is in the exponential family, and find

the posterior distribution, marginal likelihood and predictive distribution of the exponential-gamma model.

- 4. Given $\mathbf{y} \sim N_N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 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 $\mathbf{Y}_0|\mathbf{Y} \sim T_{N_0,M}(\alpha_n M + 1, \mathbf{X}_0 \mathbf{B}_n, \mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{V}_n \mathbf{X}_0^\top, \mathbf{\Psi}_n)$.
- 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 $(\mathbf{Y} \mathbf{X}\mathbf{B})^{\top}(\mathbf{Y} \mathbf{X}\mathbf{B}) = \mathbf{S} + (\mathbf{B} \widehat{\mathbf{B}})^{\top}\mathbf{X}^{\top}\mathbf{X}(\mathbf{B} \widehat{\mathbf{B}})$ where $\mathbf{S} = (\mathbf{Y} \mathbf{X}\widehat{\mathbf{B}})^{\top}(\mathbf{Y} \mathbf{X}\widehat{\mathbf{B}}), \hat{\mathbf{B}} = (\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}\mathbf{Y}$ in the multivariate regression model.

10. What is the probability that the Sun will rise tomorrow?

This is the most famous Richard Price's example developed in the Appendix of the Bayes' theorem paper [7]. Here, we implicitly use Laplace's Rule of Succession to solve this question. In particular, if we were a priori uncertain about the probability the Sun will rise on a specified day, we can assume a prior uniform 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, $Y_i \sim \mathcal{N}(\mu, \sigma)$. The sample mean and variance are

Exercises 93

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 μ .
- •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 $\log(\text{electricity price}_i)$, $\log(\text{water price}_i)$, $\log(\text{gas price}_i)$, $\operatorname{IndSocio1}_i$, $\operatorname{IndSocio2}_i$, $\operatorname{Altitude}_i$, Nrooms_i , $\operatorname{HouseholdMem}_i$, $\operatorname{Children}_i$, and $\log(\operatorname{Income}_i)$, where electricity, water and gas are monthly consumption of electricity (kWh), water (m³) and gas (m³), and other definitions are given in the Example of Section 4.3. Omit households that do not consume any of the utilities in this exercise.

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

- (a) Find the posterior mean estimates and the highest posterior density intervals at 95% of $\bf B$ and $\bf \Sigma$. Use the marginal distribution and the conditional distribution to obtain the posterior estimates of $\bf B$, and compare the results.
- (b) 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?
- (c) 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/M³ 0.70, and a gas price equal to USD/M³ 0.75.

5

Simulation methods

5.1	The inverse transform method
5.2	Method of composition
5.3	Accept and reject algorithm
5.4	Importance sampling
5.5.1	Markov chain Monte Carlo methods Some theory Gibbs sampler Metropolis-Hastings
5.6	Sequential Monte Carlo
5.7	Hamiltonian Monte Carlo

Convergence diagnostics

Part II Regression models: A GUIded tour

Graphical user interface

This chapter presents our graphical user interface (GUI) to carry out Bayesian regression analysis in a very friendly environment without any programming skills (drag and drop). Our GUI is based on an interactive web application using *shiny* [18], and packages like *MCMCpack* [67] and *bayesm* [88] from **R** software [79], and is designed for teaching and applied purposes at an introductory level. In the next chapters of the second part of this book we carry out some applications to highlight the potential of our GUI for applied researchers and practitioners.

6.1 Introduction

Our GUI allows performing inference using Bayesian regression analysis without requiring programming skills. The latter seems to be a significant impediment to increasing the use of the Bayesian framework [105, 53].

There are other available graphical user interfaces for carrying out Bayesian regression analysis. ShinyStan [98] is a very flexible open source program, but users are required to have some programming skills. BugsXLA [105] is open source, but less flexible. However, users do not need to have programming skills. Bayesian regression: Nonparametric and parametric models [53] is a very flexible and friendly GUI that is based on MATLAB Compiler for a 64-bit Windows computer. Its focus is on Bayesian nonparametric regressions, and it can be thought of for users who have mastered basic parametric models, such as the ones that we show in our GUI. There are also MATLAB toolkit, Stata and BayES, but these are not open sources.

We developed our GUI based on an interactive web application using shiny [18], and some libraries in **R** [78]. The specific libraries and commands that are used in our GUI can be seen in Table 15.1. It has ten univariate models, four multivariate, time series models, three hierarchical longitudinal, and seven Bayesian model averaging frameworks. In addition, it gives basic summaries and diagnostics of the posterior chains, as well as the posterior chains themselves, and different plots, such as trace, autocorrelation and densities.

In terms of its flexibility and possibilities, our GUI lies between ShinyStan and BugsXLA: users are not required to have any programming skills, but it

is not as advanced as [53]'s software. However, our GUI can be run in any operating system. Our GUI, which we call BEsmarter,¹ is freely available at https://github.com/besmarter/BSTApp; so users have access to all our code and datasets.

Simulated and applied datasets are in the folders **DataSim** (see Table 15.2 for details), and **DataApp** (see Table 15.3 for details) of our *GitHub* repository. The former folder also includes the files that were used to simulate different processes, so, the population parameters are available, and as a consequence these files can be used as a pedagogical tool to show some statistical properties of the inferential frameworks available in our GUI. The latter folder contains the datasets used in the applications of this second part of the book. Users should use these datasets as templates to structure their own datasets.

To display our GUI type

```
R code. How to display our graphical user interface

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

in the \mathbf{R} package console or any \mathbf{R} code editor to run our $\mathrm{GUI.^2}$

After this, users can see a new window where a presentation of our research team is displayed. In addition, the top panel in Figure 6.1 shows the class of models that can be estimated in our GUI.

6.2 Univariate models

After our GUI is deployed (see Figure 6.1), the user should select *Univariate Models* in the top panel. Then, the Figure 6.2 is displayed, and the user can see the radio button on the left hand side that shows the specific models inside this generic class. In particular, users can see that the normal model is selected from inside the class of univariate models.

Then, the right hand side panel displays a widget to upload the input dataset, which should be a csv file with headers in the first row. Users also

 $^{^1\}mathrm{Bayesian}$ econometrics: Simulations, models and applications to research, teaching and encoding with responsibility.

²We strongly recommend to type this directly, rather than copy and paste. This is due to an issue with the quotation mark.

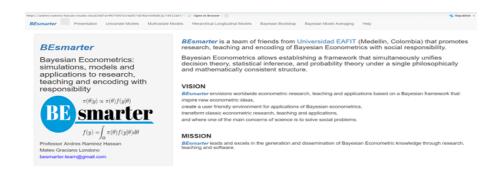


FIGURE 6.1

Display of our graphical user interface.



FIGURE 6.2

Univariate models: Specification.



FIGURE 6.3

Univariate models: Results.

should select the kind of separator used in the input file: comma, semicolon, or tab (use the folders **DataSim** and **DataApp** for the input file templates). Once users upload the dataset, they can see a data preview. Range sliders help to set the number of iterations of the Markov chain Monte Carlo algorithm, the amount of burn-in, and the thinning parameter can be selected as well (see next chapters of this second part of the book for technical details). After this, users should specify the equation. This can be done with the formula builder, where users can select the dependent variable, and the independent variables, and then click on the Build formula tab. Users can see in the Main Equation space the formula expressed in the format used by R software (see Main equation box in Figure 6.2, $y \sim x1 + x2 + x3$). Users can modify this if necessary, for instance, including higher order or interaction terms, other transformations are also allowed. This is done directly in the Main Equation space taking into account that this extra terms should follow formula command structure.³ Note that the class of univariate models includes the intercept by default, except ordered probit, where the specification has to do this explicitly, that is, ordered probit models do not admit an intercept, for identification issues (see details below).⁴ Hence, users should write down specifically this fact $(y \sim x1 + x2 + x3 - 1)$. Finally, users should define the hyperparameters of the prior; for instance, in the normal-inverse gamma model, these are the mean vector, covariance matrix, shape, and scale parameters (see Figure 6.3). However, users should take into account that our GUI has non-informative hyperparameters by default in all our modelling frameworks, so the last part is not a requirement.

After this specification process, users should click the *Go!* button to initiate the estimation. Our GUI displays the summary statistics and convergence

 $^{^3}$ See https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/3.6.2/topics/formula/stats/versions/stats/stats/stats/versions/stats

⁴An *identification* issue means that multiple values for the model parameters give rise to the same value for the likelihood function.

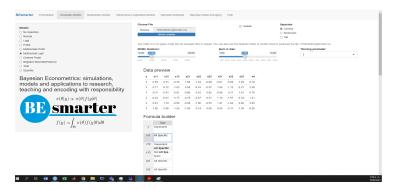


FIGURE 6.4

Univariate models: Multinomial.

diagnostics after this process is finished (see Figure 6.3). There are also widgets to download posterior chains (csv file) and graphs (pdf and eps files). Note that the order of the coefficients in the results (summary, posterior chains, and graphs) is first for the location parameters, and then for the scale parameters.

Multinomial models (probit and logit) require a dataset file to have in the first column the dependent variable, then alternative specific regressors (for instance alternatives' prices), and finally, non-alternative regressors (for instance, income). The formula builder specifies the dependent variable, and independent variables that are alternative specific and non-alternative specific (see technical details in next chapter). Specification also requires defining the base category, number of alternatives (this is also required in ordered probit), number of alternative specific regressors, and number of non-alternative regressors (see Figure 6.4). Multinomial logit also allows defining a tuning parameter, the number of degrees of freedom in this case, for the Metropolis-Hastings algorithm (see technical details in next chapter). This is a feature in our GUI when the estimation of the models is based on the Metropolis-Hastings algorithm. The order of the coefficients in the results of these models is first the intercepts (cte_l appearing in the summary display, l-th alternative), and then the non-alternative specific regressors (NAS_{il} appearing in the summary display, l-th alternative and j-th non-alternative regressor), and lastly, the coefficients for the alternative specific regressors (AS_i appearing in the summary display, j-th alternative specific regressor). Note that the non-alternative specific regressors associated with the base category are equal to zero (they do not appear in the results). In addition, some coefficients of the main diagonal of the covariance matrix are constant due to identification issues in multinomial and multivariate probit models.

In the case of the negative binomial model, users should set a dispersion parameter (see the negative binomial model in the next chapter). User should also set the censorship points and quantiles in the Tobit and quantile models, respectively.

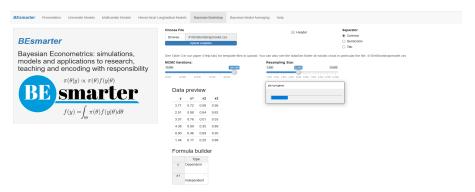


FIGURE 6.5

Univariate models: Bootstrap.

Bayesian bootstrap only requires uploading a dataset, specifying the number of iterations of the MCMC, the resampling size, and the equation (see Figure 6.5). The input file has the same structure as the file used in the univariate normal model.

6.3 Multivariate models

After our GUI is deployed (see Figure 6.1), the user should select *Multivariate Models* in the top panel. Then, the Figure 6.6 is displayed, and the user can see the radio button on the left hand side that shows the specific models inside this generic class.

Figure 6.6 displays the multivariate regression setting. In this case, the input file should have first the dependent variables, and then the regressors. If there are intercepts in each equation, there should be a column of 1's after the dependent variables in the input file. The user also has to set the number of dependent variables, the number of regressors, if necessary include the intercept, and the values of the hyperparameters (see Figure 6.6).

The input file in seemingly unrelated regressions should have first the dependent variables, and then the regressors by equation, including the intercept in each equation if necessary (column of 1's). Users should define the number of dependent variables (equations), the number of total regressors, that is, the sum of all regressors associated with the equation (if necessary include intercepts, each intercept is an additional regressor), and the number of regressors by equation (if necessary include the intercept). Users can also set the values of the hyperparameters if there is prior information.

The results of the simple multivariate and seemingly unrelated regressions



FIGURE 6.6

Multivariate models: Simple multivariate.

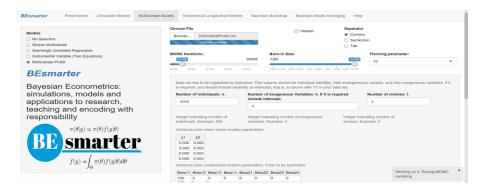


FIGURE 6.7

Multivariate models: Multivariate probit.

show first the posterior location parameters by equation, and then the posterior covariance matrix.

In the instrumental variable setting, users should specify the main equation and the instrumental equation. This setting includes intercepts by default. The first variable on the right hand side in the main equation has to be the variable with endogeneity issues. In the instrumental equation box, the dependent variable is the variable with endogeneity issues as a function of the instruments. Users can also specify the values of the hyperparameters if they have prior information. The input file should have the dependent variable, the endogenous regressor, the instruments, and the exogenous regressors. The results first list the posterior estimates of the endogenous regressor, then the location parameters of the auxiliary regression (instrumental equation), and the location parameters of the exogenous regressors. Last is the posterior covariance matrix.

The multivariate probit model requires an input dataset ordered by unit,

for instance three choices implies repeat each unit three times. The first column has to be the identification of each unit; users should use ordered integers, then the dependent variable, just one vector, composed of 0's and 1's, then the regressors, which should include a column of 1's for the intercepts. Users should set the number of units, number of regressors, and number of choices (see Figure 6.7). The results first display the posterior location parameters by equation, and then the posterior covariance matrix.

6.4 Time series model

6.5 Longitudinal (panel) models

After our GUI is deployed (see Figure 6.1), the user should select *Hierarchical Longitudinal Models* in the top panel. Then, the Figure 6.8 is 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 have first the dependent variable, then the regressors and a cross sectional identifier (i = 1, 2, ..., m). It is not a requirement to have a balanced dataset: n_i can be different for each i (see Chapter 10 for technical details). Users should specify the fixed part equation and the random part equation, both in \mathbf{R} format. In case of only requiring random intercepts, do not introduce anything in the latter part (see Figure 6.8). Users should also type the name of the cross sectional identifier variable. 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 covariance matrix.

6.6 Bayesian model average

After our GUI is deployed (see Figure 6.1), the user should select *Bayesian Model Averaging* in the top panel. Then, the Figure 6.9 is 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 based on a Gaussian distribution can be carried



FIGURE 6.8

Hierarchical longitudinal models: Specification.

out using the Bayesian information criterion (BIC) approximation, Markov chain Monte Carlo model composition (MC3), or instrumental variables (see Figure 6.9). The former two approaches require an input dataset where the first column is the dependent variable, and then, the potentially important regressors. Users should set the band width model selection parameter (O_R) and number of iterations for BIC and MC3, respectively (see Chapter 11 for technical details). The results include the posterior inclusion probability (p! = 0), expected value (EV), and standard deviation (SD) of the coefficients associated with each regressor. The BIC framework also displays the most relevant models, including the number of regressors, the coefficient of determination (R^2) , the BIC, and the posterior model probability. 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, 0 otherwise. 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 one has the dependent variable in the first column, the next columns are the regressors with endogeneity issues, and then the exogeneous regressors. The user should include a column of 1's if an intercept is required. The second input file has all the instruments. Users should also introduce the number of regressors with endogeneity issues (see Figure 6.10).

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

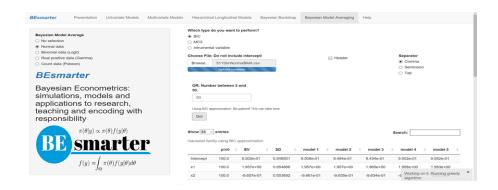


FIGURE 6.9

Bayesian model averaging: Specification and results.



FIGURE 6.10

Bayesian model averaging: Instrumental variable specification.

Warning 109

parameters of the main equation, β_l , $l = 1, 2, ..., dim \{\beta\}$, the location parameters of the auxiliary equations, $\gamma_{j,i}$, $j = 1, 2, ..., dim \{\beta_s\}$ where $dim \{\beta_s\}$ is the number of regressors with endogeneity issues, $i = 1, 2, ..., dim \{\gamma\}$, 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 11 for technical details).

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 band width model selection parameters, which are 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 see 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 of the Gaussian model.

6.7 Warning

Users should also note that sometimes our GUI shuts down. In our experience, this is due to computational issues using the implicit commands that we call when estimating some models, for instance, computationally singular systems, missing values where TRUE/FALSE needed, L-BFGS-B needs finite values of "fn", NA/NaN/Inf values, or Error in backsolve. Sometimes these issues can be solved by adjusting the dataset, for instance, avoiding high levels of multicollinearity. It should also be taken into account that when warning messages are displayed in our GUI, there is a high chance that there are convergence issues of the posterior chains. So, the results are not trustworthy. Users can identify these problems by checking the console of their RStudio sections, where the specific folder/file where the issue happened is specified. In any case, we would appreciate your feedback to improve and enhance our GUI.

We also should say there are many ways to improve the codes that we present in the following five chapters. 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 with the codes that we present in the next chapters. We could improve the computational times of our codes using parallel computing and the *rcpp* package, but this requires more advanced skills that we do not cover in this book.

We describe how to perform Bayesian inference in some of the most common univariate models: normal-inverse gamma, logit, probit, multinomial probit and logit, ordered probit, negative binomial, tobit, quantile regression, and Bayesian bootstrap in linear models. We show the posterior distributions of the parameters and some applications. In addition, we show how to perform inference in various models using three levels of programming skills: our graphical user interface (GUI), packages from **R**, and programming the posterior distributions. The first requires no programming skills, the second requires an intermediate level, and the third demands more advanced skills. We also include mathematical and computational exercises.

We can run our GUI typing

```
R code. How to display our graphical user interface

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

in the \mathbf{R} package console or any \mathbf{R} code editor.

7.1 The Gaussian linear model

The Gaussian linear model specifies $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mu$ such that $\mu \sim N(\mathbf{0}, \sigma^2 \mathbf{I_N})$ is an stochastic error, \mathbf{X} is a $N \times K$ matrix of regressors, $\boldsymbol{\beta}$ is a K-dimensional vector of location coefficients, σ^2 is the variance of the model (scale parameter), \mathbf{y} is a N-dimensional vector of a dependent variable, and N is the sample size. We describe this model using the conjugate family in Section 4.3, that is, $\pi(\boldsymbol{\beta}, \sigma^2) = \pi(\boldsymbol{\beta}|\sigma^2) \times \pi(\sigma^2)$, and this allowed to get the posterior marginal distribution for $\boldsymbol{\beta}$ and σ^2 .

We assume independent prior in this section, that is, $\pi(\beta, \sigma^2) = \pi(\beta) \times$

 $\pi(\sigma^2)$, where $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$ and $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$, $\alpha_0/2$ and $\delta_0/2$ are the shape and rate parameters. This setting allows getting the posterior conditional distributions, that is, $\pi(\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X})$ and $\pi(\sigma^2|\boldsymbol{\beta}, \mathbf{y}, \mathbf{X})$, which in turn allows to use the Gibbs sampler algorithm to perform posterior inference of $\boldsymbol{\beta}$ and σ^2 .

The likelihood function in this model is

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

Then, the conditional posterior distributions are

$$\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \sigma^2 \mathbf{B}_n),$$

and

$$\sigma^2 | \boldsymbol{\beta}, \mathbf{y}, \mathbf{X} \sim IG(\alpha_n/2, \delta_n/2),$$

where
$$\mathbf{B}_n = (\mathbf{B}_0^{-1} + \sigma^{-2}\mathbf{X}^{\top}\mathbf{X})^{-1}$$
, $\boldsymbol{\beta}_n = \mathbf{B}_n(\mathbf{B}_0^{-1}\boldsymbol{\beta}_0 + \sigma^{-2}\mathbf{X}^{\top}\mathbf{y})$, $\boldsymbol{\alpha}_n = \alpha_0 + N$ and $\boldsymbol{\delta}_n = \boldsymbol{\delta}_0 + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^{\top}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ (see Exercise 1 in this chapter).¹

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 which is in folder **DataApp** in our github repository https://github.com/besmarter/BSTApp. This dataset was used by [96] to finding the determinants of high performance soccer players in the five most important national leagues in Europe.

The specification of the model is

$$\log(\text{Value}_i) = \beta_1 + \beta_2 \text{Perf}_i + \beta_3 \text{Age}_i + \beta_4 \text{Age}_i^2 + \beta_5 \text{NatTeam}_i + \beta_6 \text{Goals}_i + \beta_7 \text{Exp}_i + \beta_8 \text{Exp}_i^2 + \mu_i.$$

where *Value* is the market value in Euros (2017), *Perf* is a measure of performance, *Age* is the players' age in years, *NatTem* is an indicator variable that takes the value of 1 if the player has been on the national team, *Goals* is the number of goals scored by the player during his career, and *Exp* is his experience in years.

We assume that the dependent variable distributes normal, then we use a normal-inverse gamma model using vague conjugate priors where $\mathbf{B}_0 = 1000\mathbf{I}_8$, $\boldsymbol{\beta}_0 = \mathbf{0}_8$, $\alpha_0 = 0.001$ and $\delta_0 = 0.001$. We perform a Gibbs sampler with 5,000 MCMC iterations plus a burn-in equal to 5,000, and a thinning parameter equal to 1.

Once our GUI is displayed (see beginning of this chapter), we should follow Algorithm A1 to run linear Gaussian models in our GUI (see Chapter 6 for details):

¹This model can be extended to consider 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.

Algorithm A1 Linear Gaussian 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

We can see in the next ${\bf R}$ codes how to perform the linear Gaussian model using the command MCMCregress of the MCMCpack package, and programming the Gibbs sampler ourselves. We should get similar results using the three approaches: GUI, package and our function. In fact, our GUI relies on the MCMCregress command. For instance, the value of a top soccer player in Europe increases 134% (exp(0.85) -1)) on average when he has played in the national team, the credible interval at 95% is (86%, 197%).

R code. The value of soccer players, using R packages

```
1 rm(list = ls())
2 set.seed(010101)
3 ###################### Linear regression: Value of
      4 Data <- read.csv("DataApplications/1ValueFootballPlayers.csv
      ", sep = ",", header = TRUE, fileEncoding = "latin1")
5 attach(Data)
6 y <- log(Value)
7 # Value: Market value in Euros (2017) of soccer players
8 # Regressors quantity including intercept
9 X <- cbind(1, Perf, Age, Age2, NatTeam, Goals, Exp, Exp2)
10 # Perf: Performance. Perf2: Performance squared. Age: Age;
      Age: Age squared.
11 # NatTeam: Indicator of national team. Goals: Scored goals.
      Goals2: Scored goals squared
12 # Exp: Years of experience. Exp2: Years of experience
      squared. Assists: Number of assists
13 k <- dim(X)[2]
14 N <- dim(X)[1]
15 # Hyperparameters
16 d0 <- 0.001/2
17 a0 <- 0.001/2
18 b0 <- rep(0, k)
19 c0 <- 1000
20 B0 <- c0*diag(k)
21 B0i <- solve(B0)
22 # MCMC parameters
23 mcmc <- 5000
24 burnin <- 5000
25 tot <- mcmc + burnin
26 thin <- 1
27 # Posterior distributions using packages: MCMCpack sets the
      model in terms of the precision matrix
28 posterior <- MCMCpack::MCMCregress(y~X-1, b0=b0, B0 = B0i,</pre>
      c0 = a0, d0 = d0, burnin = burnin, mcmc = mcmc, thin =
      thin)
29 summary(coda::mcmc(posterior))
30 Iterations = 1:5000
31 Thinning interval = 1
32 Number of chains = 1
33 Sample size per chain = 5000
34 1. Empirical mean and standard deviation for each variable,
35 plus standard error of the mean:
                       SD Naive SE Time-series SE
            Mean
37 X
            3.695499 2.228060 3.151e-02
                                              3.151e-02
                                              6.079e-05
38 XPerf
            0.035445 0.004299 6.079e-05
39 XAge
            0.778410 0.181362 2.565e-03
                                              2.565e-03
            -0.016617 0.003380 4.781e-05
40 XAge2
                                              4.781e-05
41 XNatTeam 0.850362 0.116861 1.653e-03
                                              1.689e-03
42 XGoals
            0.009097 0.001603 2.266e-05
                                              2.266e-05
43 XExp
            0.206208 0.062713 8.869e-04
                                              8.428e-04
44 XExp2
            -0.006992 0.002718 3.844e-05
                                              3.719e-05
45 sigma2
            0.969590 0.076091 1.076e-03
                                              1.076e-03
```

R. code. The value of soccer players, programming our Gibbs sampler

```
1 # Posterior distributions programming the Gibbs sampling
2 # Auxiliary parameters
3 XtX <- t(X) %*%X
4 bhat <- solve(XtX)%*%t(X)%*%y
5 an <- a0 + N
6 # Gibbs sampling functions
7 PostSig2 <- function(Beta){</pre>
8 dn <- d0 + t(y - X%*%Beta)%*%(y - X%*%Beta)
    sig2 \leftarrow invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
    return(sig2)
11 }
12 PostBeta <- function(sig2){</pre>
Bn <- solve(B0i + sig2^(-1)*XtX)</pre>
    bn <- Bn%*%(B0i%*%b0 + sig2^(-1)*XtX%*%bhat)
    Beta <- MASS::mvrnorm(1, bn, Bn)
16
    return (Beta)
17 }
18 PostBetas <- matrix(0, mcmc+burnin, k)
19 PostSigma2 <- rep(0, mcmc+burnin)</pre>
20 Beta <- rep(0, k)
21 for(s in 1:tot){
   sig2 <- PostSig2(Beta = Beta)
    PostSigma2[s] <- sig2
23
    Beta <- PostBeta(sig2 = sig2)
    PostBetas[s,] <- Beta
25
26 }
27 keep <- seq((burnin+1), tot, thin)</pre>
28 PosteriorBetas <- PostBetas[keep,]</pre>
29 colnames(PosteriorBetas) <- c("Intercept", "Perf", "Age", "
      Age2", "NatTeam", "Goals", "Exp", "Exp2")
30 summary(coda::mcmc(PosteriorBetas))
31 Iterations = 1:5000
32 Thinning interval = 1
33 Number of chains = 1
34 Sample size per chain = 5000
35 1. Empirical mean and standard deviation for each variable,
36 plus standard error of the mean:
            Mean
                        SD Naive SE Time-series SE
38 Intercept 3.663230 2.194363 3.103e-02 3.103e-02
              0.035361 0.004315 6.102e-05
                                                6.102e-05
39 Perf
40 Age
             0.780374 0.178530 2.525e-03
                                                2.525e-03
             -0.016641 0.003332 4.713e-05
                                                 4.713e-05
41 Age2
42 NatTeam
             0.850094 0.119093 1.684e-03
                                                 1.684e-03
             0.009164 0.001605 2.270e-05
43 Goals
                                                 2.270e-05
44 Exp
             0.205965 0.062985 8.907e-04
                                                 8.596e-04
            -0.007006 0.002731 3.862e-05
                                                 3.701e-05
45 Exp2
46 PosteriorSigma2 <- PostSigma2[keep]
47 summary(coda::mcmc(PosteriorSigma2))
48 Iterations = 1:5000
49 Thinning interval = 1
50 Number of chains = 1
51 Sample size per chain = 5000
52 1. Empirical mean and standard deviation for each variable,
53 plus standard error of the mean:
                              Naive SE Time-series SE
54 Mean
                    SD
                  0.077316
55 0.973309
                                  0.001093
                                                  0.001116
```

7.2 The logit model

In the logit model the dependent variable is binary, $Y_i = \{1, 0\}$, then it follows a Bernoulli distribution, $Y_i \stackrel{ind}{\sim} B(\pi_i)$, that is, $p(Y_i = 1) = \pi_i$, such that $\pi_i = \frac{\exp\{\mathbf{x}_i^{\mathsf{T}}\boldsymbol{\beta}\}}{1+\exp\{\mathbf{x}_i^{\mathsf{T}}\boldsymbol{\beta}\}}$.

The likelihood function of the logit model is

$$p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{X}) = \prod_{i=1}^{N} \pi_i^{y_i} (1 - \pi_i)^{1 - y_i}$$
$$= \prod_{i=1}^{N} \left(\frac{\exp\left\{\mathbf{x}_i^{\top} \boldsymbol{\beta}\right\}}{1 + \exp\left\{\mathbf{x}_i^{\top} \boldsymbol{\beta}\right\}} \right)^{y_i} \left(\frac{1}{1 + \exp\left\{\mathbf{x}_i^{\top} \boldsymbol{\beta}\right\}} \right)^{1 - y_i}.$$

We can specify a Normal distribution as prior, $\beta \sim N(\beta_0, \mathbf{B}_0)$. Then, the posterior distribution is

$$\pi(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) \propto \prod_{i=1}^{N} \left(\frac{\exp\left\{\mathbf{x}_{i}^{\top} \boldsymbol{\beta}\right\}}{1 + \exp\left\{\mathbf{x}_{i}^{\top} \boldsymbol{\beta}\right\}} \right)^{y_{i}} \left(\frac{1}{1 + \exp\left\{\mathbf{x}_{i}^{\top} \boldsymbol{\beta}\right\}} \right)^{1 - y_{i}} \times \exp\left\{ -\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\beta}_{0})^{\top} \mathbf{B}_{\mathbf{0}}^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_{\mathbf{0}}) \right\}.$$

The logit model does not have a standard posterior distribution. Then, a random walk Metropolis–Hastings algorithm can be used to obtain draws from the posterior distribution. A potential proposal is a multivariate Normal centered at the current value, 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 from the maximum likelihood estimation [66].²

Observe that $\log(p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{X})) = \sum_{i=1}^{N} y_i \mathbf{x}_i^{\top} \boldsymbol{\beta} - \log(1 + \exp(\mathbf{x}_i^{\top} \boldsymbol{\beta}))$. We can use this expression when calculating the acceptance parameter in the computational implementation of the Metropolis-Hastings algorithm. In particular, the acceptance parameter is

$$\alpha = \min \left\{ 1, \exp(\log(p(\mathbf{y}|\boldsymbol{\beta}^c, \mathbf{X})) + \log(\pi(\boldsymbol{\beta}^c)) - (\log(p(\mathbf{y}|\boldsymbol{\beta}^{(s-1)}, \mathbf{X})) + \log(\pi(\boldsymbol{\beta}^{(s-1)}))) \right\},$$

where β^c and $\beta^{(s-1)}$ are the draws from the proposal distribution and the previous iteration of the Markov chain, respectively.³

Example: Simulation exercise

 $^{^2}$ Tuning parameters should be set in a way such that one obtains reasonable diagnostic criteria and acceptation rates.

³Formulating the acceptance rate using log helps to mitigate computational problems.

The logit model 117

Let's do a simulation exercise to check the performance of the algorithm. Set $\beta = \begin{bmatrix} 0.5 & 0.8 & -1.2 \end{bmatrix}^{\top}$, $x_{ik} \sim N(0,1)$, k = 2, 3 and i = 1, 2, ..., 10000. We set as hyperparameters $\beta_0 = \begin{bmatrix} 0 & 0 \end{bmatrix}^{\top}$ and $\mathbf{B}_0 = 1000\mathbf{I}_3$. The tune

parameter for the Metropolis-Hastings algorithm is equal to 1.

Once our GUI is displayed (see beginning of this chapter), we should follow Algorithm A2 to run logit models in our GUI (see Chapter 6 for details):

Algorithm A2 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
- 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
- 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 can see in the next \mathbf{R} codes how to perform the logit model using the command MCMClogit of the MCMCpack package, and programming the Metropolis-Hastings algorithm ourselves.

We should get similar results using the three approaches: GUI, package and our function. Our GUI relies on the MCMClogit command. In particular, we obtain an acceptance rate of 0.46, and the diagnostics suggest that the posterior chains behave well. In general, the 95% credible intervals encompass the population values, and the mean and median are very close to these values.

$R.\ code.\ Simulation\ of\ the\ logit\ model\ estimation\ using\ R\ packages$

```
################## Logit: Simulation
      ############################
2 # Simulate data
3 rm(list = ls())
4 set.seed(010101)
5 N <- 10000 # Sample size
_{6} B <- c(0.5, 0.8, -1.2) # Population location parameters
7 x2 <- rnorm(N) # Regressor
8 x3 <- rnorm(N) # Regressor</pre>
9 X <- cbind(1, x2, x3) # Regressors
10 XB <- X%*%B
11 PY \leftarrow exp(XB)/(1 + exp(XB)) # Probability of Y = 1
12 Y <- rbinom(N, 1, PY) # Draw Y's
13 table(Y) # Frequency
# write.csv(cbind(Y, x2, x3), file = "DataSimulations/")
      LogitSim.csv") # Export data
_{15} # MCMC parameters
16 iter <- 5000; burnin <- 1000; thin <- 5; tune <- 1
17 # Hyperparameters
18 K <- dim(X)[2]
19 b0 <- rep(0, K)
20 c0 <- 1000
21 BO <- c0*diag(K)
22 B0i <- solve(B0)
23 # Posterior distributions using packages: MCMCpack sets the
      model in terms of the precision matrix
24 RegLog <- MCMCpack::MCMClogit(Y~X-1, mcmc = iter, burnin =
      burnin, thin = thin, b0 = b0, B0 = B0i, tune = tune)
25 summary(RegLog)
26 Iterations = 1001:5996
27 Thinning interval = 5
_{28} Number of chains = 1
29 Sample size per chain = 1000
30 1. Empirical mean and standard deviation for each variable,
31 plus standard error of the mean:
        Mean
                   SD Naive SE Time-series SE
       0.4896 0.02550 0.0008064
                                       0.001246
33 X
34 Xx2 0.8330 0.02730 0.0008632
                                        0.001406
35 Xx3 -1.2104 0.03049 0.0009643
                                       0.001536
36 2. Quantiles for each variable:
        2.5%
                 25%
                         50%
                                  75%
                                        97.5%
       0.4424
               0.4728
                       0.4894
                                0.5072
                                        0.5405
38 X
39 Xx2 0.7787
               0.8159
                       0.8327
                                0.8505
                                        0.8852
40 Xx3 -1.2758 -1.2296 -1.2088 -1.1902 -1.1513
```

R. code. Simulation of the logit model estimation programming our M-H algorithm

```
# Posterior distributions programming the Metropolis-
      Hastings algorithm
2 MHfunc \leftarrow 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){
    Xm <- cbind(1, X) # Regressors</pre>
    K <- dim(Xm)[2] # Number of location parameters</pre>
    BETAS <- matrix(0, iter + burnin, K) # Space for posterior
       chains
    Reg <- glm(y ~ Xm - 1, family = binomial(link = "logit"))</pre>
      # Maximum likelihood estimation
    BETA <- Reg$coefficients # Maximum likelihood parameter
      estimates
    tot <- iter + burnin # Total iterations M-H algorithm
    COV <- vcov(Reg) # Maximum likelihood covariance matrix
9
    COVt <- tau^2*solve(solve(BO) + solve(COV)) # Covariance
      matrix for the proposal distribution
    Accep <- rep(0, tot) # Space for calculating the
      acceptance rate
    # Create progress bar in case that you want to see
      iterations progress
    pb <- winProgressBar(title = "progress bar", min = 0,</pre>
13
    max = tot, width = 300)
14
    for(it in 1:tot){
      BETAc \leftarrow BETA + MASS::mvrnorm(n = 1, mu = rep(0, K),
16
      Sigma = COVt) # Candidate location parameter
      likecand <- sum((Xm%*%BETAc) * Y - apply(Xm%*%BETAc, 1,
17
      function(x) log(1 + exp(x))) # Log likelihood for the
      candidate
      likepast <- sum((Xm%*%BETA) * Y - apply((Xm%*%BETA), 1,
      function(x) log(1 + exp(x)))) # Log likelihood for the
      actual draw
      priorcand <- (-1/2)*crossprod((BETAc - b0), solve(B0))%*</pre>
      %(BETAc - b0) # Log prior for candidate
      priorpast <- (-1/2)*crossprod((BETA - b0), solve(B0))%*%</pre>
20
      (BETA - b0) # Log prior for actual draw
      alpha <- min(1, exp((likecand - priorcand) - (likepast -
       priorpast))) #Probability of selecting candidate
      u \leftarrow runif(1) # Decision rule for selecting candidate
      if(u < alpha){</pre>
23
24
        BETA <- BETAc # Changing reference for candidate if
      selected
         Accep[it] <- 1 # Indicator if the candidate is
      accepted
26
      BETAS[it, ] <- BETA # Saving draws
27
      setWinProgressBar(pb, it, title=paste( round(it/tot*100,
28
       0),
      "% done"))
29
    }
30
31
    close(pb)
    keep <- seq(burnin, tot, thin)</pre>
32
    return(list(Bs = BETAS[keep[-1], ], AceptRate = mean(Accep
       [keep[-1]])))
34 }
```

R. code. Simulation of the logit model programming our M-H algorithm, results

```
Posterior <- MHfunc(y = Y, X = cbind(x2, x3), iter = iter,</pre>
      burnin = burnin, thin = thin) # Running our M-H function
       changing some default parameters.
paste("Acceptance rate equal to", round(Posterior$AceptRate,
       2), sep = " ")
3 "Acceptance rate equal to 0.46"
4 PostPar <- coda::mcmc(Posterior$Bs)</pre>
5 # Names
6 colnames(PostPar) <- c("Cte", "x1", "x2")</pre>
7 # Summary posterior draws
8 summary(PostPar)
9 Iterations = 1:1000
10 Thinning interval = 1
11 Number of chains = 1
12 Sample size per chain = 1000
13 1. Empirical mean and standard deviation for each variable,
14 plus standard error of the mean:
15 Mean
          SD Naive SE Time-series SE
16 Cte 0.4893 0.02427 0.0007674
                                       0.001223
17 x1 0.8309 0.02699 0.0008536
                                       0.001440
                                       0.001423
18 x2 -1.2107 0.02943 0.0009308
19 2. Quantiles for each variable:
      2.5%
               25%
                       50%
                               75%
                                     97.5%
20
21 Cte 0.4431 0.4721 0.4899 0.5059 0.5344
22 x1 0.7817 0.8123 0.8305 0.8505 0.8833
23 x2 -1.2665 -1.2309 -1.2107 -1.1911 -1.1538
_{24} # Trace and density plots
25 plot(PostPar)
26 # Autocorrelation plots
27 coda::autocorr.plot(PostPar)
28 # Convergence diagnostics
29 coda::geweke.diag(PostPar)
30 Fraction in 1st window = 0.1
31 Fraction in 2nd window = 0.5
32 Cte
         x 1
                x2
33 -0.975 -3.112 1.326
34 coda::raftery.diag(PostPar,q=0.5,r=0.05,s = 0.95)
35 Quantile (q) = 0.5
36 Accuracy (r) = +/- 0.05
_{37} Probability (s) = 0.95
38 Burn-in Total Lower bound Dependence
39 (M)
                               factor (I)
           (N) (Nmin)
40 Cte 6
               731 385
                                   1.90
41 x1 6
               703
                                   1.83
                    385
42 x2 6
               725
                    385
                                   1.88
43 coda::heidel.diag(PostPar)
44 Stationarity start p-value
45 test
               iteration
_{
m 46} Cte passed
                   1
                              0.4436
47 x1 passed
                    101
                              0.3470
                              0.0872
48 x2 passed
                    1
49 Halfwidth Mean
                  Halfwidth
50 test
51 Cte passed
                 0.489 0.00240
52 x1 passed
                 0.832 0.00268
                -1.211 0.00279
53 x2 passed
```

7.3 The probit model

The probit model also has as dependent variable a binary outcome. In this case, there is a latent variable $(y_i^*, \text{ unobserved})$ that defines the structure of the estimation problem. In particular,

$$Y_i = \begin{cases} 0, \ Y_i^* \le 0 \\ 1, \ Y_i^* > 0 \end{cases},$$

such that $Y_i^* = \mathbf{x}_i^{\top} \boldsymbol{\beta} + \mu_i$, $\mu_i \stackrel{i.i.d.}{\sim} N(0,1)$. This implies $P(Y_i = 1) = \pi_i = 0$ $\Phi(\mathbf{x}_{i}^{\top}\boldsymbol{\beta}).$

[2] implemented data augmentation [100] to apply a Gibbs sampling algorithm in this model. Augmenting this model with Y_i^* , we can have the likelihood contribution from observation i, $p(y_i|y_i^*) = \mathbb{1}_{y_i=0} \mathbb{1}_{y_i^* \leq 0} + \mathbb{1}_{y_i=1} \mathbb{1}_{y_i^* > 0}$, where $\mathbb{1}_A$ is an indicator function that takes the value of 1 when condition A is satisfied.

The posterior distribution is $\pi(\boldsymbol{\beta}, \boldsymbol{y^*}|\boldsymbol{y}, \boldsymbol{X}) \propto \prod_{i=1}^{N} \left[\mathbbm{1}_{y_i=0} \mathbbm{1}_{y_i^* \leq 0} + \mathbbm{1}_{y_i=1} \mathbbm{1}_{y_i^* > 0}\right] \times$ $N_N(y^*|X\beta, I_n) \times N_K(\beta|\beta_0, B_0)$ when taking a Gaussian distribution as prior $\boldsymbol{\beta} \sim N_k(\boldsymbol{\beta}_0, \boldsymbol{B}_0)$. This implies

$$y_i^* | \boldsymbol{\beta}, \boldsymbol{y}, \boldsymbol{X} \sim \left\{ \begin{aligned} TN_{(-\infty,0]}(\mathbf{x}_i^{\top} \boldsymbol{\beta}, \mathbf{1}) &, \ \mathbf{y_i} = \mathbf{0} \\ TN_{(0,\infty)}(\mathbf{x}_i^{\top} \boldsymbol{\beta}, \mathbf{1}) &, \ \mathbf{y_i} = \mathbf{1} \end{aligned} \right\},$$

$$\boldsymbol{\beta}|\boldsymbol{y}^*, \boldsymbol{X} \sim N(\boldsymbol{\beta}_n, \boldsymbol{B}_n),$$

where
$$B_n = (B_0^{-1} + X^{\top}X)^{-1}$$
, and $\beta_n = B_n(B_0^{-1}\beta_0 + X^{\top}y^*)$.

Example: Determinants of hospitalization

We use the dataset named **2HealthMed.csv**, which is in folder **DataApp** in our github repository (https://github.com/besmarter/BSTApp and was used by [83]. Our dependent variable is a binary indicator with a value equal to 1 if an individual was hospitalized in 2007, and 0 otherwise.

The specification of the model is

$$Hosp_{i} = \beta_{1} + \beta_{2}SHI_{i} + \beta_{3}Female_{i} + \beta_{4}Age_{i} + \beta_{5}Age_{i}^{2} + \beta_{6}Est2_{i} + \beta_{7}Est3_{i} + \beta_{8}Fair_{i} + \beta_{9}Good_{i} + \beta_{10}Excellent_{i},$$

where SHI is a binary variable equal to 1 if the individual is in a subsidized

⁴The variance in this model is set to 1 due to identification restrictions. Observe that $P(Y_i = 1 | \mathbf{x_i}) = \mathbf{P}(\mathbf{Y_i^*} > \mathbf{0} | \mathbf{x_i}) = \mathbf{P}(\mathbf{x_i^\top} \boldsymbol{\beta} + \mu_i > \mathbf{0} | \mathbf{x_i}) = \mathbf{P}(\mu_i > -\mathbf{x_i^\top} \boldsymbol{\beta} | \mathbf{x_i}) = \mathbf{P}(\mathbf{c} \times \mu_i > -\mathbf{c} \times \mathbf{x_i^\top} \boldsymbol{\beta} | \mathbf{x_i}) \ \forall c > 0$. Multiplying for a positive constant does not affect the probability of $Y_i = 1$. 5TN denotes a truncated normal density.

health care program and 0 otherwise, Female is an indicator of gender, Age in years, Est2 and Est3 are indicators of socioeconomic status, the reference is Est1, which is the lowest, and self perception of health status where bad is the reference.

Let's set $\beta_0 = \mathbf{0}_{10}$, $\mathbf{B}_0 = \mathbf{I}_{10}$, iterations, burn-in and thinning parameters equal to 10000, 1000 and 1, respectively. We can use the Algorithm A1 to run the probit model in our GUI. We should select *Probit* model in stage 2. Our GUI relies in the command *rbprobitGibbs* from the package *bayesm* to perform inference in the Probit model. The following \mathbf{R} code shows how to run this example using the command *rbprobitGibbs*. We asked to program a Gibbs sampler algorithm to perform inference in the probit model in the exercises.

We find evidence that gender and self-perceived health status affect the probability of hospitalization. Women have a higher probability of being hospitalized than men, and a better perception of health status decreases this probability.

R. code. Determinants of hospitalization

```
nydata <- read.csv("DataApplications/2HealthMed.csv", header</pre>
        = T, sep = ",")
2 attach (mydata)
3 str(mydata)
4 K <- 10 # Number of regressors
5 b0 <- rep(0, K) # Prio mean
6 BOi <- diag(K) # Prior precision (inverse of covariance)
7 Prior <- list(betabar = b0, A = B0i) # Prior list</pre>
8 y <- Hosp # Dependent variables</pre>
9 X <- cbind(1, SHI, Female, Age, Age2, Est2, Est3, Fair, Good
       , Excellent) # Regressors
10 Data \leftarrow list(y = y, X = X) # Data list
11 Mcmc <- list(R = 10000, keep = 1, nprint = 0) # MCMC
       parameters
12 RegProb <- bayesm::rbprobitGibbs(Data = Data, Prior = Prior,</pre>
       Mcmc = Mcmc) # Inference using bayesm package
13 PostPar <- coda::mcmc(RegProb$betadraw) # Posterior draws
14 colnames(PostPar) <- c("Cte", "SHI", "Female", "Age", "Age2", "Est2", "Est3", "Fair", "Good", "Excellent") # Names
15 summary(PostPar) # Posterior summary
16 Iterations = 1:10000
17 Thinning interval = 1
18 Number of chains = 1
19 Sample size per chain = 10000
20 2. Quantiles for each variable:
                                     50%
           2.5%
                        25%
                                                 75%
                                                           97.5%
22 Cte
              -1.22e+00 -1.03e+00 -9.43e-01 -8.50e-01 -0.671744
              -1.24e-01 -4.63e-02 -6.30e-03
23 SHI
                                               3.26e-02 0.104703
24 Female
              2.80e-02 9.65e-02
                                   1.28e-01
                                               1.60e-01
                                                           0.223123
25 Age
             -7.55e-03 -2.50e-03
                                    1.25e-04
                                                2.80e-03
                                                          0.007646
             -4.98e-05 9.05e-06 4.02e-05 7.07e-05
-1.89e-01 -1.23e-01 -8.84e-02 -5.32e-02
                                                           0.000128
26 Age2
27 Est2
                                                           0.012714
             -2.13e-01 -1.03e-01 -4.73e-02 1.01e-02 0.109527
28 Est3
29 Fair
             -7.09e-01 -5.69e-01 -4.93e-01 -4.16e-01 -0.269494
              -1.42e+00 -1.28e+00 -1.20e+00 -1.12e+00 -0.982533
30 Good
31 Excellent -1.33e+00 -1.15e+00 -1.06e+00 -9.74e-01 -0.795881
```

7.4 The multinomial probit model

The multinomial probit model is used to model the choice of the l-th alternative over a set L mutually exclusive options. We observe

$$y_{il} = \begin{cases} 1, & y_{il}^* \ge \max\{\boldsymbol{y}_i^*\} \\ 0, & \text{otherwise} \end{cases},$$

such that $\boldsymbol{y}_i^* = \boldsymbol{X}_i \boldsymbol{\delta} + \boldsymbol{\mu}_i$, $\boldsymbol{\mu}_i$ $\overset{i.i.d.}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$, \boldsymbol{y}_i^* is an unobserved latent L dimensional vector, $\boldsymbol{X}_i = \begin{bmatrix} (1 \ \boldsymbol{c}_i^\top) \otimes \boldsymbol{I}_L \ \boldsymbol{A}_i \end{bmatrix}$ is an $L \times j$ matrix of regressors for each alternative, $l = 1, 2, \ldots, L, j = L \times (1 + dim \{\boldsymbol{c}_i\}) + a, \boldsymbol{c}_i$ is a vector of the individuals' specific characteristics, \boldsymbol{A}_i is an $L \times a$ matrix of alternative-varying regressors, a is the number of alternative-varying regressors, and $\boldsymbol{\delta}$ is a j dimensional vector of parameters.

We take into account simultaneously the alternative-varying regressors (alternative attributes) and alternative-invariant regressors (individual characteristics). 6 \boldsymbol{y}_i^* can be stacked up into a multiple regression with correlated stochastic errors, $\boldsymbol{y}^* = \boldsymbol{X}\boldsymbol{\delta} + \boldsymbol{\mu}$, where $\boldsymbol{y}^* = [\boldsymbol{y}_1^{*\top}, \boldsymbol{y}_2^{*\top}, \dots, \boldsymbol{y}_N^{*\top}], \boldsymbol{X} = [\boldsymbol{X}_1^{\top}, \boldsymbol{X}_2^{\top}, \dots, \boldsymbol{X}_N^{\top}]^{\top}$, and $\boldsymbol{\mu} = [\boldsymbol{\mu}_1^{\top}, \boldsymbol{\mu}_2^{\top}, \dots, \boldsymbol{\mu}_N^{\top}]^{\top}$. Following the practice of expressing \boldsymbol{y}_{il}^* relative to \boldsymbol{y}_{iL}^* by letting $\boldsymbol{w}_i = \boldsymbol{v}_{il}^{\top}$.

Following the practice of expressing y_{il}^* relative to y_{iL}^* by letting $\boldsymbol{w}_i = [w_{i1}, w_{i2}, \ldots, w_{iL-1}]^\top$, $w_{il} = y_{il}^* - y_{iL}^*$, we can write $\boldsymbol{w}_i = \boldsymbol{R}_i \boldsymbol{\beta} + \boldsymbol{\epsilon}_i$, $\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \boldsymbol{\Omega})$, where $\boldsymbol{R}_i = [(1 \ \boldsymbol{c}_i^\top) \otimes \boldsymbol{I}_{L-1} \ \boldsymbol{\Delta} \boldsymbol{A}_i]$ is an $L-1 \times k$ matrix where $\Delta \boldsymbol{A}_i = \boldsymbol{A}_{li} - \boldsymbol{A}_{Li}, l = 1, 2, \ldots, L-1$, that is, the last row of \boldsymbol{A}_i is subtracted from each row of \boldsymbol{A}_i , and $\boldsymbol{\beta}$ is a k dimensional vector, $k = (L-1) \times (1 + dim\{\boldsymbol{c}_i\}) + a$.

Observe that $\boldsymbol{\beta}$ contains the same last a elements as $\boldsymbol{\delta}$, that is, alternative specific attributes coefficients, but the first $(L-1)\times(1+\dim\{\boldsymbol{c}_i\})$ -th elements are $\delta_{jl}-\delta_{jL},\,j=1+\dim\{\boldsymbol{c}_i\},\,l=1,2,\ldots,L-1$, that is, the difference between the coefficients of each qualitative response and the L-th alternative for the individuals' characteristics. This makes it difficult to interpret the multinomial probit coefficients.

Note that in multinomial models, for each alternative specific attribute, it is only required to estimate one coefficient for all alternatives, whereas for individuals' characteristics (non-alternative specific regressors), it is necessary to estimate L-1 coefficients (the coefficient of the base alternative is set equal to 0).

The likelihood function in this model is $p(\boldsymbol{\beta}, \boldsymbol{\Omega} | \boldsymbol{y}, \boldsymbol{R}) = \prod_{i=1}^{N} \prod_{l=1}^{L} p_{il}^{y_{il}}$ where $p_{il} = p(y_{il}^* \ge \max(\boldsymbol{y}_i^*))$.

We assume independent priors, $\beta \sim N(\beta_0, B_0)$ and $\Omega^{-1} \sim W(\alpha_0, \Sigma_0)$. We can employ Gibbs sampling in this model because this is a standard

⁶Note that this model is not identified if Σ is unrestricted. The likelihood function is the same if a scalar random variable is added to each of the L latent regressions.

⁷W denotes the Wishart density.

Bayesian linear regression model when data augmentation in \boldsymbol{w} is used. The posterior conditional distributions are

$$eta|\Omega, oldsymbol{w} \sim N(oldsymbol{eta}_n, oldsymbol{B}_n),$$
 $oldsymbol{\Omega}^{-1}|oldsymbol{eta}, oldsymbol{w} \sim W(lpha_n, oldsymbol{\Sigma}_n),$ where $oldsymbol{B}_n = (oldsymbol{B}_0^{-1} + oldsymbol{X}^{*\top} oldsymbol{X}^*)^{-1}, \ oldsymbol{eta}_n = oldsymbol{B}_n(oldsymbol{B}_0^{-1} oldsymbol{eta}_0 + oldsymbol{X}^{*\top} oldsymbol{w}^*), \ oldsymbol{\Omega}^{-1} = oldsymbol{C}^{\top} oldsymbol{C}, \ oldsymbol{X}_i^{*\top} = oldsymbol{C}^{\top} oldsymbol{w}_i, \ oldsymbol{X}^* = oldsymbol{X}^* oldsymbol{W}_i, \ oldsymbol{X}^* = oldsymbol{C}^{\top} oldsymbol{W}_i, \ oldsymbol{W}^* = oldsymbol{W}_i, \ oldsymbol{X}^* = oldsymbol{C}^{\top} oldsymbol{W}_i, \ oldsymbol{X}^* = oldsymbol{X}^* oldsymbol{X}^* oldsymbol{W}_i, \ oldsymbol{W}^* = oldsymbol{W}_i, \ oldsymbol{W}_i, \ oldsymbol{W}^* = oldsymbol{W}_i, \ oldsymbol{W}_i, \ oldsymbol{$

 $(\Sigma_0 + \sum_{i=1}^N (\boldsymbol{w}_i - \boldsymbol{R}_i \boldsymbol{\beta})^{\top} (\boldsymbol{w}_i - \boldsymbol{R}_i \boldsymbol{\beta}))^{-1}$. We can collapse the multinomial vector \boldsymbol{y}_i into the indicator variable $d_i = \sum_{l=1}^{L-1} l \times \mathbb{1}_{\max(\boldsymbol{w}_l) = w_{il}}$. Then the distribution of $\boldsymbol{w}_i | \boldsymbol{\beta}, \boldsymbol{\Omega}^{-1}, d_i$ is an L-1 dimensional Gaussian distribution truncated over the appropriate cone in \mathcal{R}^{L-1} . [68] propose drawing from the univariate conditional distributions $w_{il} | \boldsymbol{w}_{i,-l}, \boldsymbol{\beta}, \boldsymbol{\Omega}^{-1}, d_i \sim TN_{I_{il}}(m_{il}, \tau_{ll}^2)$, where

$$I_{il} = \begin{cases} w_{il} > \max(\boldsymbol{w}_{i,-l}, 0), & d_i = l \\ w_{il} < \max(\boldsymbol{w}_{i,-l}, 0), & d_i \neq l \end{cases},$$

and permuting the columns and rows of Ω^{-1} so that the *l*-th column and row is the last,

$$\boldsymbol{\Omega}^{-1} = \begin{bmatrix} \boldsymbol{\Omega}_{-l,-l} & \boldsymbol{\omega}_{-l,l} \\ \boldsymbol{\omega}_{l,-1} & \boldsymbol{\omega}_{l,l} \end{bmatrix}^{-1} = \begin{bmatrix} \boldsymbol{\Omega}_{-l,-l}^{-1} + \tau_{ll}^{-2} \boldsymbol{f} \boldsymbol{f}^{\top} & -\boldsymbol{f} \tau_{ll}^{-2} \\ -\tau_{ll}^{-2} \boldsymbol{f}^{\top} & \tau_{ll}^{-2} \end{bmatrix}$$

where $\boldsymbol{f} = \boldsymbol{\Omega}_{-l,-l}^{-1} \boldsymbol{\omega}_{-l,l}$, $\tau_{ll}^2 = \omega_{ll} - \boldsymbol{\omega}_{l,-l} \boldsymbol{\Omega}_{-l,-1}^{-1} \boldsymbol{\omega}_{-l,l}$, $m_{il} = \boldsymbol{r}_{il}^{\top} \boldsymbol{\beta} + \boldsymbol{f}^{\top} (\boldsymbol{w}_{i,-l} - \boldsymbol{R}_{i,-l} \boldsymbol{\beta})$, $\boldsymbol{w}_{i,-l}$ is an L-2 dimensional vector of all components of \boldsymbol{w}_i excluding w_{il} , \boldsymbol{r}_{il} is the l-th row of \boldsymbol{R}_i , $l = 1, 2, \dots, L-1$.

The identified parameters are obtained by normalizing with respect to one of the diagonal elements $\frac{1}{\omega_{1,1}^{0.5}} \boldsymbol{\beta}$ and $\frac{1}{\omega_{1,1}} \boldsymbol{\Omega}.^9$

A warning is required here! This model is an example where we have to make decisions about setting the model in an identified parameter space or unidentified space. The mixing properties of the posterior draws can be better in the latter case [69], this means less computational burden. However, we should to recover the identified space in a final stage. In addition, we should take into account that defining priors in the unidentified space may have unintended consequences on the posterior distributions of the identified space [73]. The multinomial probit model that is presented in this section is set in the unidentified space [68]. A version of the multinomial probit in the identified

⁸Observe that the identification issue in this model is due to scaling w_{il} by a positive constant does not change the value of d_i .

⁹Our GUI is based on the *bayesm* package that takes into account this identification restriction to display the outcomes of the posterior chains.

space is presented by [69].

Example: Choice of fishing mode

We used in this application the dataset 3Fishing.csv from [16, p. 491]. The dependent variable is mutually exclusive alternatives regarding fishing modes (mode), where beach is equal to 1, pier is equal to 2, private boat is equal to 3, and chartered boat (baseline alternative) is equal to 4. In this model, we have

$$\boldsymbol{X}_i = \begin{cases} 1 & 0 & 0 & 0 & \text{Income}_i & 0 & 0 & 0 & \text{Price}_{i,1} & \text{Catch rate}_{i,1} \\ 0 & 1 & 0 & 0 & 0 & \text{Income}_i & 0 & 0 & \text{Price}_{i,2} & \text{Catch rate}_{i,2} \\ 0 & 0 & 1 & 0 & 0 & 0 & \text{Income}_i & 0 & \text{Price}_{i,3} & \text{Catch rate}_{i,3} \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & \text{Income}_i & \text{Price}_{i,4} & \text{Catch rate}_{i,4} \end{cases}.$$

In this example chartered boat is the base category, the number of choice categories is four, there are two alternative-specific regressors (price and catch rate), and one non alternative-specific regressor (income). This setting involves the estimation of eight location parameters (β): three intercepts, three for income, one for price, and one for catch rate. This is the order of the posterior chains in our GUI. Note that the location coefficients are set equal to 0 for the baseline category. For multinomial models, we strongly recommend using the last category as the baseline.

We also get posterior estimates for a 3×3 covariance matrix (four alternatives minus one), where the element (1,1) is equal to 1 due to identification restrictions, and elements 2 and 4 are the same, as well as 3 and 7, and 6 and 8, due to symmetry. Observe that this identification restriction implies NaN values in [41] and [48] tests for element (1,1) of the covariance matrix, and just eight dependence factors associated with the remaining elements of the covariance matrix.

Once our GUI is displayed (see beginning of this chapter), we should follow Algorithm A3 to run multinomial probit models in our GUI (see Chapter 6 for details), which in turn uses the command rmnpGibbs from the bayesm package.

We ran 100,000 MCMC iterations plus 10,000 as burn-in with a thinning parameter equal to 5, where all priors use default values for the hyperparameters in our GUI. We found that the 95% credible intervals of the coefficient associated with income for beach and private boat alternatives are equal to (8.58e-06, 8.88e-05) and (3.36e-05, 1.45e-04). This suggests that the probability of choosing these alternatives increases compared to a chartered boat when income increases. In addition, an increase in the price or a decrease in the catch rate for specific fishing alternatives imply lower probabilities of choosing them as the 95% credible intervals are (-9.91e-03, -3.83e-03) and (1.40e-01, 4.62e-01), respectively. However, the posterior chain diagnostics suggest there are convergence issues with the posterior draws (see exercise 5).

 $^{^{10}}$ This is the order in the pdf, eps and csv files that can be downloaded from our GUI.

Algorithm A3 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 $Build\ formula\ button\ to\ generate\ the\ formula\ in\ {\bf R}\ syntax.$
- 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

R. code. Choice of fishing mode, results

```
Iterations = 10005:110000
  Thinning interval = 5
  Number of chains = 1
  Sample size per chain = 20000
  Quantiles for each variable:
           2.5%
                       25%
                                   50%
                                                        97.5%
                                               75%
           -5.83e-01 -4.08e-01 -3.22e-01 -2.37e-01
                                                     -7.93e-02
                                2.16e-02
           -1.93e-01 -4.14e-02
                                           7.93e-02
           -8.15e-01 -5.43e-01 -4.29e-01 -3.33e-01
           8.58e-06
                     3.61e-05
                                4.95e-05
                                           6.27e-05
  NAS_1_2
           -3.24e-05
                     -7.04e-06
                                 5.52e - 06
                                           1.93e - 05
                                                      5.17e - 05
           3.36e - 05
                      6.38e - 05
                                 8.08e - 05
           -9.91e-03
                     -7.90e-03
                                -6.86e-03
                                          -5.93e-03
13 AS 1
                                                      -3.83e-03
14 AS_2
            1.40e-01
                      2.25e-01
                                2.72e-01
```

7.5 The multinomial logit model

The multinomial logit model is used to model mutually exclusive discrete outcomes or qualitative response variables. However, this model assumes the independence of irrelevant alternatives (IIA), meaning that the choice between two alternatives does not depend on a third alternative. We consider the multinomial mixed logit model (not to be confused with the random parameters logit model), which accounts for both alternative-varying regressors (conditional) and alternative-invariant regressors (multinomial) simultaneously. 11

In this setting there are L mutually exclusive alternatives, and the dependent variable y_{il} is equal to 1 if the lth alternative is chosen by individual i, and 0 otherwise, $l = \{1, 2, ..., L\}$. The likelihood function is $p(\beta|\boldsymbol{y},\boldsymbol{X}) = \prod_{i=1}^{N} \prod_{l=1}^{L} p_{il}^{y_{il}}$, where the probability that individual i chooses the alternative l is given by $p_{il} := p(y_i = l|\boldsymbol{\beta},\boldsymbol{X}) = \frac{\exp\{\boldsymbol{x}_{il}^{\top}\boldsymbol{\beta}_l\}}{\sum_{j=1}^{L} \exp\{\boldsymbol{x}_{ij}^{\top}\boldsymbol{\beta}_j\}}$, \boldsymbol{y} and \boldsymbol{X} are the vector and matrix of the dependent variable and regressors, and $\boldsymbol{\beta}$ is the vector containing all the coefficients. Remember that coefficients associated with alternative-invariant regressors are set to 0 for the baseline category, and the coefficients associated with the alternative-varying regressors are the same for all the categories. In addition, we assume $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \boldsymbol{B}_0)$ as prior distribution. Thus, the posterior distribution is $\pi(\boldsymbol{\beta}|\boldsymbol{y},\boldsymbol{X}) \propto p(\boldsymbol{\beta}|\boldsymbol{y},\boldsymbol{X}) \times \pi(\boldsymbol{\beta})$.

As the multinomial logit model does not have a standard posterior distribution, [89] propose a "tailored" independent Metropolis–Hastings algorithm where the proposal distribution is a multivariate Student's t distribution with v degrees of freedom (tuning parameter), mean equal to the maximum likelihood estimator, and scale equal to the inverse of the Hessian matrix.

Example: Simulation exercise

Let's do a simulation exercise to check the performance of the Metropolis-Hastings algorithm to perform inference in the multinomial logit model. Assume a situation where there are three alternatives, one alternative-invariant regressor plus the intercept, and three alternative-varying regressors. The population parameters are $\beta_1 = [1 - 2.5 \ 0.5 \ 0.8 \ -3]^{\top}$, $\beta_2 = [1 - 3.5 \ 0.5 \ 0.8 \ -3]^{\top}$ and $\beta_3 = [0 \ 0 \ 0.5 \ 0.8 \ -3]^{\top}$, the first two elements of the vectors are associated with the intercept and the alternative-invariant regressor, and the last three elements with the alternative-varying regressors. The sample size is 1000, and all regressors are simulated from standard normal distributions.

We can deploy our GUI using the command line at the beginning of this chapter. We should follow Algorithm A4 to run multinomial logit models in our GUI (see Chapter 6 for details):

The following code in \mathbf{R} shows how to implement the M-H algorithm from scratch. The first part simulates the dataset, the second part builds the

¹¹The multinomial mixed logit model can be implemented as a conditional logit model.

Algorithm A4 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 $Build\ formula\ button\ to\ generate\ the\ formula\ in\ {f R}\ syntax.$
- 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

loglikelihood function, and the third part implements the M-H algorithm. We use vague priors centered on zero, and covariance matrix $1000\,I_7$. We observe that the posterior estimates closely match the population parameters, and all 95% credible intervals contain the population parameters.

R. code. Simulation of the multinomial logit model

```
remove(list = ls())
  2 set.seed(12345)
  3 # Simulation of data
  4 N<-1000 # Sample Size
  5 B < -c(0.5, 0.8, -3); B1 < -c(-2.5, -3.5, 0); B2 < -c(1, 1, 0)
  6 # Alternative specific attributes of choice 1, for instance,
                         price, quality and duration of choice {\bf 1}
       X1<-matrix(cbind(rnorm(N,0,1),rnorm(N,0,1),rnorm(N,0,1)),N,</pre>
                      length(B))
             Alternative specific attributes of choice 2, for instance,
                        price, quality and duration of choice 2
  9 X2<-matrix(cbind(rnorm(N,0,1),rnorm(N,0,1),rnorm(N,0,1)),N,</pre>
                      length(B))
10 # Alternative specific attributes of choice 3, for instance,
                         price, quality and duration of choice 3
11 X3<-matrix(cbind(rnorm(N,0,1),rnorm(N,0,1),rnorm(N,0,1)),N,
                      length(B))
12 X4<-matrix(rnorm(N,1,1),N,1)
 \begin{tabular}{ll} & \tt V1 &-B2 [1] + \tt X1 \% * \% B + B1 [1] * \tt X4 ; & \tt V2 &-B2 [2] + \tt X2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt X2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt X2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt X2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V2 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V3 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V3 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V3 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V3 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V3 \% * \% B + B1 [2] * \tt X4 ; & \tt V3 &-B2 (2) + \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \tt V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [2] * \textmd V3 \% * \% B + B1 [
                       [3] + X3\% * %B + B1 [3] * X4
14 suma \leftarrow exp(V1) + exp(V2) + exp(V3)
15 p1 \leftarrow exp(V1)/suma; p2 \leftarrow exp(V2)/suma; p3 \leftarrow exp(V3)/suma
16 p <-cbind (p1, p2, p3)</pre>
17 y <- apply(p,1, function(x)sample(1:3, 1, prob = x, replace =
                         TRUE))
 y1 < -y = 1; y2 < -y = 2; y3 < -y = 3
```

R. code. Simulation of the multinomial logit model

```
1 # Log likelihood
2 log.L<- function(Beta){</pre>
     V1 \leftarrow Beta[1] + Beta[3] * X4 + X1 \% * Beta[5:7]
    V2 <-Beta[2] +Beta[4] *X4+X2%*%Beta[5:7]
    V3<- X3%*%Beta[5:7]
     suma \leftarrow exp(V1) + exp(V2) + exp(V3)
     p11 \leftarrow exp(V1)/suma; p22 \leftarrow exp(V2)/suma; p33 \leftarrow exp(V3)/suma
     suma2 <- NULL
     for(i in 1:N){
       suma1 <-y1[i] *log(p11[i])+y2[i] *log(p22[i])+y3[i] *log(p33
       [i])
       suma2 <-c(suma2, suma1)}</pre>
     logL <- sum (suma2)
     return(-logL)
13
14 }
15 # Parameters: Proposal
16 k <- 7
17 res.optim<-optim(rep(0, k), log.L, method="BFGS", hessian=</pre>
       TRUE)
18 MeanT <- res.optim$par</pre>
19 ScaleT <- as.matrix(Matrix::forceSymmetric(solve(res.optim$</pre>
       hessian))) # Force this matrix to be symmetric
20 # Hyperparameters: Priors
21 B0 \leftarrow 1000*diag(k); b0 \leftarrow rep(0, k)
22 MHfunction <- function(iter, tuning){</pre>
    Beta <- rep(0, k); Acept <- NULL
     BetasPost <- matrix(NA, iter, k)</pre>
24
     pb <- winProgressBar(title = "progress bar", min = 0, max</pre>
       = iter, width = 300)
     for(s in 1:iter){
       LogPostBeta <- -log.L(Beta) + mvtnorm::dmvnorm(Beta,</pre>
       mean = b0, sigma = B0, log = TRUE)
       BetaC <- c(LaplacesDemon::rmvt(n=1, mu = MeanT, S =</pre>
       ScaleT, df = tuning))
       LogPostBetaC <- -log.L(BetaC) + mvtnorm::dmvnorm(BetaC,</pre>
       mean = b0, sigma = B0, log = TRUE)
       alpha <- min(exp((LogPostBetaC-mvtnorm::dmvt(BetaC,</pre>
       delta = MeanT, sigma = ScaleT, df = tuning, log = TRUE))
       -(LogPostBeta-mvtnorm::dmvt(Beta, delta = MeanT, sigma =
        ScaleT, df = tuning, log = TRUE))),1)
       u <- runif(1)
       if(u <= alpha){</pre>
         Acepti <- 1; Beta <- BetaC
34
35
         Acepti <- 0; Beta <- Beta
36
37
       BetasPost[s, ] <- Beta; Acept <- c(Acept, Acepti)</pre>
       setWinProgressBar(pb, s, title=paste( round(s/iter*100,
38
       0),"% done"))
39
40
     close(pb); AcepRate <- mean(Acept)</pre>
     Results <- list(AcepRate = AcepRate, BetasPost = BetasPost
42
     return(Results)
43 }
```

R. code. Simulation of the multinomial logit model

```
# MCMC parameters
2 mcmc <- 10000; burnin <- 1000; thin <- 5; iter <- mcmc +
      burnin; keep <- seq(burnin, iter, thin); tuning <- 6 #
      Degrees of freedom
3 ResultsPost <- MHfunction(iter = iter, tuning = tuning)</pre>
4 summary(coda::mcmc(ResultsPost$BetasPost[keep[-1], ]))
5 Iterations = 1:2000
  Thinning interval = 1
7 Number of chains = 1
8 Sample size per chain = 2000
9 1. Empirical mean and standard deviation for each variable,
10 plus standard error of the mean:
                  SD Naive SE Time-series SE
12 [1,]
        0.9711 0.20162 0.004508
                                       0.004508
        0.9742 0.20934 0.004681
14 [3,] -2.4350 0.18950 0.004237
                                       0.004137
       -3.4195 0.24656 0.005513
15 [4,]
16 [5,]
        0.5253 0.07396 0.001654
                                       0.001654
        0.8061 0.08007 0.001790
17 [6,]
                                       0.001790
18 [7,] -3.0853 0.17689 0.003955
                                       0.003955
19 2. Quantiles for each variable:
                 25%
                          50%
        0.5862 0.8367
                        0.9650 1.1017 1.3683
21 var1
22 var2
        0.5679 0.8310 0.9681 1.1151 1.3761
23 var3 -2.8239 -2.5607 -2.4291 -2.3050 -2.0812
24 var4 -3.9176 -3.5806 -3.4074 -3.2496 -2.9423
        0.3840
                0.4761
                        0.5250
                                0.5759
        0.6555 0.7494
                        0.8064 0.8616 0.9604
26 var6
  var7 -3.4476 -3.1991 -3.0777 -2.9641 -2.7500
```

7.6 Ordered probit model

The ordered probit model is used when there is a natural order in the categorical response variable. In this case, there is a latent variable $y_i^* = \boldsymbol{x}_i^{\top} \boldsymbol{\beta} + \mu_i, \ \mu_i \overset{i.i.d.}{\sim} N(0,1)$ such that $y_i = l$ if and only if $\alpha_{l-1} < y_i^* \le \alpha_l, l = \{1, 2, \dots, L\}$, where $\alpha_0 = -\infty, \ \alpha_1 = 0$ and $\alpha_L = \infty$. Then,

 $^{^{12}}$ Identification issues necessitate setting the variance in this model equal to 1 and $\alpha_1 = 0$. Observe that multiplying y_i^* by a positive constant or adding a constant to all of the cut-offs and subtracting the same constant from the intercept does not affect y_i .

 $p(y_i = l) = \Phi(\alpha_l - \boldsymbol{x}_i^{\top} \boldsymbol{\beta}) - \Phi(\alpha_{l-1} - \boldsymbol{x}_i^{\top} \boldsymbol{\beta})$, and the likelihood function is $p(\boldsymbol{\beta}, \boldsymbol{\alpha} | \boldsymbol{y}, \boldsymbol{X}) = \prod_{i=1}^{N} p(y_i = l | \boldsymbol{\beta}, \boldsymbol{\alpha}, \boldsymbol{X})$. There are independent priors of this model, $\pi(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \pi(\boldsymbol{\beta}) \times \pi(\boldsymbol{\gamma})$, where

There are independent priors of this model, $\pi(\beta, \gamma) = \pi(\beta) \times \pi(\gamma)$, where $\beta \sim N(\beta_0, \mathbf{B}_0)$ and $\gamma \sim N(\gamma_0, \Gamma_0)$, $\gamma = [\gamma_2, \gamma_3, \dots, \gamma_{L-1}]^{\mathsf{T}}$, such that $\alpha = \left[\exp\{\gamma_2\}, \sum_{l=2}^3 \exp\{\gamma_l\}, \dots, \sum_{l=2}^{L-1} \exp\{\gamma_l\}\right]^{\mathsf{T}}$. The latter structure imposes the ordinal condition in the cut-offs.

This model does not have a standard conditional posterior distribution for γ (α), but it does have a standard conditional distribution for β once data augmentation is used. Then, we can use a Metropolis-within-Gibbs sampling algorithm. In particular, we use Gibbs sampling algorithms to draw β and y^* ,

$$\boldsymbol{\beta}|\boldsymbol{y}^*, \boldsymbol{\alpha}, \boldsymbol{X} \sim N(\boldsymbol{\beta}_n, \boldsymbol{B}_n),$$

where
$$\boldsymbol{B}_n = (\boldsymbol{B}_0^{-1} + \boldsymbol{X}^{\top} \boldsymbol{X})^{-1}, \ \boldsymbol{\beta}_n = \boldsymbol{B}_n (\boldsymbol{B}_0^{-1} \boldsymbol{\beta}_0 + \boldsymbol{X}^{\top} \boldsymbol{y}^*), \text{ and } y_i^* | \boldsymbol{\beta}, \boldsymbol{\alpha}, \boldsymbol{y}, \boldsymbol{X} \sim TN_{(\alpha_{y_i-1}, \alpha_{y_i})}(\boldsymbol{x}_i^{\top} \boldsymbol{\beta}, \boldsymbol{1}).$$

We use a random-walk Metropolis–Hastings algorithm for γ that has as proposal a Gaussian distribution with mean equal to the current value, and covariance matrix $s^2(\Gamma_0^{-1} + \hat{\Sigma}_{\gamma}^{-1})^{-1}$, where s > 0 is a tuning parameter, and $\hat{\Sigma}_{\gamma}$ is the sample covariance matrix associated with γ from the maximum likelihood estimation.

Example: Determinants of preventive health care visits

We used the file named 2HealthMed.csv in this applications. In particular, the dependent variable is MedVisPrevOr, which is an ordered variable equal to 1 if the individual did not visit a physician for preventive reasons, 2 if the individual visited once in that year, and so on, until it is equal to 6 for visiting five or more times. The latter category is 1.6% of the sample. Observe that the dependent variable has six categories.

In this example, the set of regressors is given by SHI, which an indicator of being in the subsidized health care system (1 means being in the system), sex (Female), age (linear and squared), socioeconomic conditions indicator (Est2 and Est3), the lowest is the baseline category, self perception of health status (Fair, Good and Excellent), where Bad is the baseline, and education level, primary (PriEd), high school (HighEd), vocational (VocEd), and university (UnivEd), no education is the baseline category.

We ran this application with 50,000 MCMC iterations plus 10,000 as burnin, and thinning parameter equal to 5. This setting means 10,000 effective posterior draws. We set $\beta_0 = \mathbf{0}_{11}$, $\mathbf{B}_0 = 1000\mathbf{I}_{11}$, $\gamma_0 = \mathbf{0}_4$, $\Gamma_0 = \mathbf{I}_4$, and the tuning parameter is 1.

We can run the ordered probit models in our GUI following the steps in the Algorithm A5.

The following \mathbf{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.

Algorithm A5 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 vectors and covariance matrices. This step is not necessary as by default our GUI uses non-informative priors
- 8: Select the tuning parameter for the Metropolis-Hastings algorithm
- 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. Determinants of preventive health care visits

```
rm(list = ls())
  set.seed (010101)
  Data <- read.csv("DataApplications/2HealthMed.csv", sep
       ", header = TRUE, fileEncoding = "latin1")
  attach (Data)
  y <- MedVisPrevOr
    MedVisPrevOr: Oredered variable for preventive visits to
      doctors in one year: 1 (none), 2 (once), ... 6 (five or
     <- cbind(SHI, Female, Age, Age2, Est2, Est3, Fair, Good,</pre>
       Excellent, PriEd, HighEd, VocEd, UnivEd)
    \leftarrow dim(X)[2]
  L <- length(table(y))
10 # Hyperparameters
11 b0 <- rep(0, k); c0 <- 1000; B0 <- c0*diag(k)
12 \text{ gamma0} \leftarrow \text{rep}(0, L-2); \text{ Gamma0} \leftarrow \text{diag}(L-2)
13 # MCMC parameters
14 mcmc <- 60000+1; thin <- 5; tuningPar <- 1/(L-2)^0.5
15 DataApp \leftarrow list(y = y, X = X, k = L)
16 Prior <- list(betabar = b0, A = solve(B0), dstarbar = gamma0
       , Ad = Gamma0)
17 mcmcpar <- list(R = mcmc, keep = 5, s = tuningPar)
18 PostBeta <- bayesm::rordprobitGibbs(Data = DataApp, Prior =</pre>
      Prior, Mcmc = mcmcpar)
```

R. code. Determinants of preventive health care visits, results

```
1 BetasPost <- coda::mcmc(PostBeta[["betadraw"]])</pre>
colnames(BetasPost) <- c("SHI", "Female", "Age", "Age2", "</pre>
      Est2", "Est3", "Fair", "Good", "Excellent", "PriEd", "
      HighEd", "VocEd", "UnivEd")
3 summary(BetasPost)
4 Iterations = 1:12000
5 Thinning interval = 1
6 Number of chains = 1
7 Sample size per chain = 12000
8 1. Empirical mean and standard deviation for each variable,
9 plus standard error of the mean:
10 Mean
              SD Naive SE Time-series SE
11 SHI
             0.0654824 2.281e-02 2.082e-04
                                                 3.357e - 04
12 Female
             -0.0374788 1.908e-02 1.742e-04
                                                 1.742e-04
             0.0190336 1.869e-03 1.706e-05
                                                 4.576e-05
13 Age
14 Age2
             -0.0002328 2.438e-05 2.225e-07
                                                 6.690e-07
15 Est2
             0.0949445 2.226e-02 2.032e-04
                                                 4.659e-04
16 Est3
             -0.1383965 3.411e-02 3.114e-04
                                                 3.459e-04
17 Fair
             0.6451828 5.375e-02 4.907e-04
                                                 3.924e - 03
             0.7343932 4.955e-02 4.523e-04
18 Good
                                                 4.491e-03
             0.9826531 6.393e-02 5.836e-04
19 Excellent
                                                 5.261e-03
             0.0309418 2.376e-02 2.169e-04
                                                 2.221e-04
20 PriEd
21 HighEd
             -0.1805753 2.910e-02 2.656e-04
                                                 3.456e-04
22 VocEd
             0.1395760 9.640e-02 8.800e-04
                                                 9.291e-04
             -0.2218120 1.189e-01 1.086e-03
                                                 1.086e-03
23 UnivEd
24 2. Quantiles for each variable:
                                  50%
                                             75%
                                                      97.5%
          2.5%
                      25%
                     0.04995 0.06540 0.08085
26 SHI
             0.02090
                                                  0.11021
             27 Female
                                                  0.00023
             0.01550 0.01781 0.01902 0.02023
                                                  0.02268
28 Age
29 Age2
            -0.00028 -0.00024 -0.00023 -0.00021 -0.00018
                                                 0.13933
30 Est2
             0.05149 0.08004 0.09482 0.10968
             -0.20559 -0.16144 -0.13815 -0.11563
31 Est3
                                                 -0.07179
             0.55799 0.61295
32 Fair
                               0.64148
                                        0.67268
                                                  0.74395
33 Good
             0.66690 0.70808
                               0.73032 0.75406
                                                 0.81064
34 Excellent 0.88919 0.94770
                               0.97836 1.01026
                                                  1.08460
             -0.01584 0.01493
                                        0.04718
35 PriEd
                               0.03101
                                                  0.07732
             -0.23782 -0.20035 -0.18021 -0.16073
36 HighEd
                                                 -0.12435
            -0.04911 0.07474
37 VocEd
                               0.13811
                                        0.20414
                                                  0.33331
            -0.45381 -0.30239 -0.22193 -0.14148
38 UnivEd
39 # Convergence diagnostics
40 coda::geweke.diag(BetasPost)
41 coda::raftery.diag(BetasPost,q=0.5,r=0.05,s = 0.95)
42 coda::heidel.diag(BetasPost)
43 # Cut offs
44 Cutoffs <- PostBeta[["cutdraw"]]
45 summary (Cutoffs)
46 coda::geweke.diag(Cutoffs)
47 coda::heidel.diag(Cutoffs)
48 coda::raftery.diag(Cutoffs[,-1],q=0.5,r=0.05,s = 0.95)
```

The results suggest that older individuals (at decreasing rate) in the subsidized health program, characterized in the second socioeconomic status with increasing good self perception of health condition, and not having high school as their highest education degree, have a higher probability of visiting a physician for preventive health aims. Convergence diagnostics look well, except for the self health perception draws.

We also got the posterior estimates of the cutoffs in the ordered probit model. These estimates are necessary to calculate the probability that an individual is in a specific category of visiting physicians. Due to identification restrictions, the first cutoff is set equal to 0. That is why we have NaN values in [41] and [48] tests, and we observe only four values in the [82] test, which correspond to the remaining free cutoffs. It seems that these cutoff estimates have some convergence issues when taking as diagnostic tool the [82] test. Their dependence factors are also very high.

7.7 Negative binomial model

The dependent variable in the negative binomial model is a nonnegative integer or count. In contrast to the Poisson model, the negative binomial model takes into account over-dispersion. The Poisson model has equal mean and variance (equi-dispersion).

We assume that $y_i \stackrel{i.n.d.}{\sim} NB(\gamma, \theta_i)$, that is, the density function for individual i is $\frac{\Gamma(y_i + \gamma)}{\Gamma(\gamma)y_i!}(1 - \theta_i)^{y_i}\theta_i^{\gamma}$, where the success probability is $\theta_i = \frac{\gamma}{\lambda_i + \gamma}$, $\lambda_i = \exp\left\{x_i^{\top}\beta\right\}$ is the mean, and $\gamma = \exp\left\{\alpha\right\}$ is the target for number of successful trials, or dispersion parameter.

We assume independent priors for this model are $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \boldsymbol{B}_0)$ and $\alpha \sim G(\alpha_0, \delta_0)$.¹³

This model does not have standard conditional posterior distributions, so [89] use a random-walk Metropolis–Hastings algorithm where the proposal distribution for $\boldsymbol{\beta}$ is Gaussian centered at the current stage with covariance matrix $s_{\boldsymbol{\beta}}^2 \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\beta}}$ where $s_{\boldsymbol{\beta}}$ is a tuning parameter and $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\beta}}$ is the maximum likelihood covariance estimator. In addition, the proposal for α is normal centered at the current value, with variance $s_{\alpha}^2 \hat{\sigma}_{\alpha}^2$ where s_{α} is a tuning parameter and $\hat{\sigma}_{\alpha}^2$ is the maximum likelihood variance estimator.

Example: Simulation exercise

Let's do a simulation exercise to check the performance of the M-H algorithms in the negative binomial model. There are two regressors, $x_{i1} \sim U(0,1)$ and $x_{i1} \sim N(0,1)$, and the intercept. The dispersion parameter is $\gamma = \exp\{1.2\}$), and $\boldsymbol{\beta} = \begin{bmatrix}1 & 1\end{bmatrix}^{\mathsf{T}}$. The sample size is 1,000.

 $^{^{13}}G$ denotes a gamma density.

We run this simulation using 10,000 MCMC iterations, a burn-in equal to 1,000, and a thinning parameter equal to 5. We set vague priors for the location parameters, particularly, $\beta_0 = \mathbf{0}_3$ and $B_0 = 1000 I_3$, and $\alpha_0 = 0.5$ and $\delta_0 = 0.1$, which are the default values in the rnegbinRw command from bayesm package in \mathbf{R} . In addition, the tuning parameters of the Metropolis–Hastings algorithms are $s_\beta = 2.93/k^{1/2}$ and $s_\alpha = 2.93$, which are also the default parameters in rnegbinRw, k is the number of location parameters.

We can run the negative binomial models in our GUI following the steps in the Algorithm A6.

Algorithm A6 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 ${\bf R}$ syntax. You can modify the formula in the ${\bf Main\ equation}$ box using valid arguments of the formula command structure in ${\bf R}$
- 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 \mathbf{R} code shows how to perform inference in the negative binomial model programming the M-H algorithms from scratch. We ask to estimate this example using the rnegbinRw command in exercise 8.

We observe from the results that all 95% credible intervals encompass the population parameters, and the posterior means are very close to the population parameters.

$R.\ code.\ Simulation\ of\ the\ negative\ binomial\ model$

```
1 rm(list = ls())
2 set.seed(010101)
_{\rm 3} N <- 2000 # Sample size
4 x1 <- runif(N); x2 <- rnorm(N)
5 X <- cbind(1, x1, x2); k <- dim(X)[2]; B <- rep(1, k)
6 alpha <- 1.2; gamma <- exp(alpha); lambda <- exp(X%*%B)
7 y <- rnbinom(N, mu = lambda, size = gamma)</pre>
8 # log likelihood
9 logLik <- function(par){
10 alpha <- par[1]; beta <- par[2:(k+1)]</pre>
    gamma <- exp(alpha)</pre>
    lambda <- exp(X%*%beta)
    return(-logLikNB)
15 }
16 # Parameters: Proposal
17 par0 <- rep(0.5, k+1)
res.optim <- optim(par0, logLik, method="BFGS", hessian=TRUE
19 res.optim$par
20 res.optim$convergence
21 Covar <- solve(res.optim$hessian)</pre>
22 CovarBetas <- Covar[2:(k+1),2:(k+1)]</pre>
23 VarAlpha <- Covar[1:1]
24 # Hyperparameters: Priors
25 B0 \leftarrow 1000*diag(k); b0 \leftarrow rep(0, k)
26 alpha0 <- 0.5; delta0 <- 0.1
```

R. code. Simulation of the negative binomial model, M-H algorithm

```
1 # Metropolis-Hastings function
2 MHfunction <- function(iter, sbeta, salpha){</pre>
    Beta <- rep(0, k); Acept1 <- NULL; Acept2 <- NULL
    BetasPost <- matrix(NA, iter, k); alpha <- 1</pre>
    alphaPost <- rep(NA, iter); par <- c(alpha, Beta)</pre>
    pb <- winProgressBar(title = "progress bar", min = 0, max</pre>
      = iter, width = 300)
    for(s in 1:iter){
      LogPostBeta <- -logLik(par) + dgamma(alpha, shape =
      alpha0, scale = delta0, log = TRUE) + mvtnorm::dmvnorm(
      Beta, mean = b0, sigma = B0, log = TRUE)
      BetaC <- c(MASS::mvrnorm(1, mu = Beta, Sigma = sbeta^2*</pre>
      CovarBetas))
       parC <- c(alpha, BetaC)</pre>
      LogPostBetaC <- -logLik(parC) + dgamma(alpha, shape =
11
      alpha0, scale = delta0, log = TRUE) + mvtnorm::dmvnorm(
      BetaC, mean = b0, sigma = B0, log = TRUE)
      alpha1 <- min(exp((LogPostBetaC - mvtnorm::dmvnorm(BetaC</pre>
       , mean = Beta, sigma = sbeta^2*CovarBetas, log = TRUE))
       -(LogPostBeta - mvtnorm::dmvnorm(Beta, mean = Beta,
       sigma = sbeta^2*CovarBetas, log = TRUE))),1)
      u1 <- runif(1)
       if(u1 <= alpha1){Acept1i <- 1; Beta <- BetaC}else{</pre>
14
         Acept1i <- 0; Beta <- Beta
16
       par <- c(alpha, Beta)</pre>
17
      LogPostBeta <- -logLik(par) + dgamma(alpha, shape =</pre>
18
       alpha0, scale = delta0, log = TRUE) + mvtnorm::dmvnorm(
      Beta, mean = b0, sigma = B0, log = TRUE)
       alphaC <- rnorm(1, mean = alpha, sd = salpha*VarAlpha
19
       ^0.5)
       parC <- c(alphaC, Beta)</pre>
20
       LogPostBetaC <- -logLik(parC) + dgamma(alphaC, shape =
21
      alpha0, scale = delta0, log = TRUE) + mvtnorm::dmvnorm(
Beta, mean = b0, sigma = B0, log = TRUE)
       alpha2 <- min(exp((LogPostBetaC - dnorm(alphaC, mean =</pre>
       alpha, sd = salpha*VarAlpha^0.5, log = TRUE))-(
      LogPostBeta - dnorm(alpha, mean = alpha, sd = salpha*
VarAlpha^0.5, log = TRUE))),1)
       u2 <- runif(1)
       if(u2 <= alpha2){Acept2i <- 1; alpha <- alphaC}else{</pre>
         Acept2i <- 0; alpha <- alpha
25
      BetasPost[s, ] <- Beta; alphaPost[s] <- alpha</pre>
28
       Acept1 <- c(Acept1, Acept1i); Acept2 <- c(Acept2,</pre>
29
       Acept2i)
       setWinProgressBar(pb, s, title=paste( round(s/iter*100,
      0),"% done"))
31
    close(pb)
32
33
    AcepRateBeta <- mean(Acept1); AcepRateAlpha <- mean(Acept2
    Results <- list(AcepRateBeta = AcepRateBeta, AcepRateAlpha
        = AcepRateAlpha, BetasPost = BetasPost, alphaPost =
       alphaPost)
    return(Results)
35
36 }
```

R. code. Simulation of the negative binomial model, results

```
1 # MCMC parameters
2 mcmc <- 10000
3 burnin <- 1000
4 thin <- 5
5 iter <- mcmc + burnin
6 keep <- seq(burnin, iter, thin)
7 sbeta <- 2.93/sqrt(k); salpha <- 2.93
8 # Run M-H
9 ResultsPost <- MHfunction(iter = iter, sbeta = sbeta, salpha</pre>
       = salpha)
10 ResultsPost $ AcepRateBeta
11 ResultsPost$AcepRateAlpha
summary(coda::mcmc(ResultsPost$BetasPost[keep[-1], ]))
13 Iterations = 1:2000
14 Thinning interval = 1
15 Number of chains = 1
_{16} Sample size per chain = 2000
17 1. Empirical mean and standard deviation for each variable,
18 plus standard error of the mean:
        Mean
                  SD Naive SE Time-series SE
20 [1,] 1.0270 0.04799 0.0010730
                                     0.0014727
21 [2,] 0.9981 0.07752 0.0017333
                                      0.0024262
22 [3,] 0.9677 0.02343 0.0005239
                                      0.0007182
23 2. Quantiles for each variable:
                25%
                     50%
                              75% 97.5%
        2.5%
25 var1 0.9343 0.9943 1.0255 1.0592 1.122
26 var2 0.8445 0.9448 0.9980 1.0520 1.144
27 var3 0.9242 0.9512 0.9678 0.9839 1.013
28 summary(coda::mcmc(ResultsPost$alphaPost[keep[-1]]))
29 Iterations = 1:2000
30 Thinning interval = 1
31 Number of chains = 1
32 Sample size per chain = 2000
33 1. Empirical mean and standard deviation for each variable,
34 plus standard error of the mean:
                                   Naive SE Time-series SE
        Mean
                         SD
                 0.058769
                                 0.001314
                                             0.001427
36 1.282664
37 2. Quantiles for each variable:
38 2.5% 25% 50% 75% 97.5%
39 1.173 1.242 1.282 1.320 1.407
```

 $Tobit \ model$ 141

7.8 Tobit model

The dependent variable is partially observed in Tobit models due to sampling schemes, whereas the regressors are completely observed. In particular,

$$y_i = \left\{ \begin{array}{l} L \ , \quad y_i^* < L \\ y_i^* \ , \ L \le y_i^* < U \\ U \ , \quad y_i^* \ge U \end{array} \right\},$$

where $y_i^* \overset{i.n.d.}{\sim} N(\boldsymbol{x}_i^{\top} \boldsymbol{\beta}, \sigma^2).^{14}$

We use conjugate independent priors $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \boldsymbol{B}_0)$ and $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$, and data augmentation using \boldsymbol{y}_C^* such that $\boldsymbol{y}_{C_i}^* \stackrel{i.n.d.}{\sim} N(\boldsymbol{x}_i^{\top} \boldsymbol{\beta}, \sigma^2)$, $y_{C_i} = \left\{ y_{C_i^L}^* \cup y_{C_i^U}^* \right\}$ are lower and upper censored data. This allows implementing the Gibbs sampling algorithm [20]. Then,

$$\pi(\boldsymbol{\beta}, \sigma^2, \boldsymbol{y^*} | \boldsymbol{y}, \boldsymbol{X}) \propto \prod_{i=1}^{N} \left[\mathbb{1}_{y_i = L} \mathbb{1}_{y_{C_i^L}^* < L} + \mathbb{1}_{L \le y_i < U} + \mathbb{1}_{y_i = U} \mathbb{1}_{y_{C_i^U}^* \ge U} \right] \times N(y_i^* | \boldsymbol{x}_i^\top \boldsymbol{\beta}, \sigma^2) \times N(\boldsymbol{\beta} | \boldsymbol{\beta}_0, \boldsymbol{B}_0) \times IG(\sigma^2 | \alpha_0 / 2, \delta_0 / 2)$$

The posterior distributions are

$$\begin{aligned} y_{C_i}^*|\boldsymbol{\beta}, \sigma^2, \boldsymbol{y}, \boldsymbol{X} \sim & \begin{cases} TN_{(-\infty,L)}(\boldsymbol{x}_i^\top \boldsymbol{\beta}, \sigma^2) \;, \; y_i = L \\ TN_{[U,\infty)}(\boldsymbol{x}_i^\top \boldsymbol{\beta}, \sigma^2) \;\;, \; y_i = U \end{cases} , \\ \boldsymbol{\beta}|\sigma^2, \boldsymbol{y}, \boldsymbol{X} \sim \mathcal{N}(\boldsymbol{\beta}_n, \sigma^2 \boldsymbol{B}_n), \\ \sigma^2|\boldsymbol{\beta}, \boldsymbol{y}, \boldsymbol{X} \sim IG(\alpha_n/2, \delta_n/2), \\ \text{where } \boldsymbol{B}_n \; = \; (\boldsymbol{B}_0^{-1} + \sigma^{-2} \boldsymbol{X}^\top \boldsymbol{X})^{-1}, \; \boldsymbol{\beta}_n \; = \; \boldsymbol{B}_n(\boldsymbol{B}_0^{-1} \boldsymbol{\beta}_0 + \sigma^{-2} \boldsymbol{X}^\top \boldsymbol{y}^*), \\ \boldsymbol{\alpha}_n = \boldsymbol{\alpha}_0 + N \; \text{and} \; \boldsymbol{\delta}_n = \boldsymbol{\delta}_0 + (\boldsymbol{y}^* - \boldsymbol{X} \boldsymbol{\beta})^\top (\boldsymbol{y}^* - \boldsymbol{X} \boldsymbol{\beta}). \end{aligned}$$

Example: The market value of soccer players in Europe continues

We continue the example of the market value of soccer players from Section 7.1. We specify the same equation, but assume the sample is censored from below, and have just information of soccer players whose market value is higher than one million euros. The dependent variable is log(ValueCens), and the left censoring point is 13.82.

The Algorithm A7 shows how to estimate Tobit models in our GUI. Our GUI uses the command MCMCtobit from the package MCMCpack.

We run this application using the same hyperparameters that we set in

¹⁴We can set L or U equal to $-\infty$ or ∞ to model data censored in just one side.

Algorithm A7 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

the example of Section 7.1. All results seem similar to those in the example of linear models. In addition, the posterior chains seem to achieve good diagnostics.

Tobit model 143

R. code. The value of soccer player with left censoring

```
1 rm(list = ls()); set.seed(010101)
2 Data <- read.csv("DataApplications/1ValueFootballPlayers.csv</pre>
      ", sep = ",", header = TRUE, fileEncoding = "latin1")
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 \leftarrow dim(X)[1]
8 # Hyperparameters
9 d0 <- 0.001/2; a0 <- 0.001/2
10 b0 <- rep(0, k); c0 <- 1000; B0 <- c0*diag(k)
11 B0i <- solve(B0)
_{12} # MCMC parameters
13 mcmc <- 50000
14 burnin <- 10000
15 tot <- mcmc + burnin
16 thin <- 1
17 # Posterior distributions using packages: MCMCpack sets the
      model in terms of the precision matrix
18 posterior <- MCMCpack::MCMCtobit(y~X-1, b0=b0, B0 = B0i, c0</pre>
       = a0, d0 = d0, burnin = burnin, mcmc = mcmc, thin =
      thin, below = 13.82, above = Inf)
19 summary(coda::mcmc(posterior))
20 Iterations = 1:50000
21 Thinning interval = 1
22 Number of chains = 1
23 Sample size per chain = 50000
24 1. Empirical mean and standard deviation for each variable,
25 plus standard error of the mean:
                     SD Naive SE Time-series SE
          Mean
27 X
            1.045830 2.641038 1.181e-02
                                              1.673e-02
                                               2.247e-05
            0.033990 0.004515 2.019e-05
28 XPerf
            1.025099 0.213368 9.542e-04
29 XAge
                                               1.335e-03
            -0.021871 0.004004 1.791e-05
30 XAge2
                                               2.542e-05
31 XNatTeam 0.847495 0.125924 5.631e-04
                                               6.393e-04
32 XGoals
            0.010088 0.001649 7.377e-06
                                               7.691e-06
            0.174383 0.069948 3.128e-04
зз ХЕхр
                                               3.846e-04
            -0.005652 0.002970 1.328e-05
                                               1.561e-05
34 XExp2
            0.982906 0.095965 4.292e-04
                                               6.727e-04
35 sigma2
36 2. Quantiles for each variable:
          2.5%
                      25%
                                50%
                                           75%
                                                    97.5%
            -4.174794 -0.725691
                                 1.076420
                                            2.840533
38 X
                                                     6.1935618
39 XPerf
            0.025110 0.030949
                                 0.033980
                                            0.037003
                                                      0.0428650
            0.608620 0.880648 1.023043
40 XAge
                                           1.168486
                                                      1.4480001
41 XAge2
            -0.029801 -0.024556 -0.021822 -0.019164 -0.0140990
42 XNatTeam 0.603953 0.762394 0.846461
                                            0.932056
                                                     1.0960274
43 XGoals
            0.006875
                       0.008977
                                 0.010091
                                            0.011197
                                                      0.0133323
44 XExp
            0.038752
                       0.127167
                                 0.173880
                                            0.221355
                                                      0.3122043
            -0.011483 -0.007623 -0.005654 -0.003662
45 XExp2
                                                      0.0001615
            0.811953 0.915246 0.977257
                                           1.043158
46 sigma2
```

R. code. The value of soccer player with left censoring, Gibbs sampler

```
1 # Gibbs sampling functions
2 XtX <- t(X)%*%X
3 PostBeta <- function(Y1, sig2){</pre>
    Bn <- solve(B0i + sig2^(-1)*XtX)</pre>
    bn <- Bn%*%(B0i%*%b0 + sig2^(-1)*t(X)%*%Y1)
    Beta <- MASS::mvrnorm(1, bn, Bn)</pre>
    return (Beta)
8 }
9 PostYl <- function(Beta, L, U, i){</pre>
    Ylmean <- X[i,]%*%Beta
    if(y[i] == L){
      Yli <- truncnorm::rtruncnorm(1, a = -Inf, b = L, mean =
      Ylmean, sd = sig2^0.5)
      if(y[i] == U){
        Yli <- truncnorm::rtruncnorm(1, a = U, b = Inf, mean =
       Ylmean, sd = sig2^0.5)
      }else{
16
         Yli <- y[i]
18
    }
19
    return(Yli)
20
21 }
22 PostSig2 <- function(Beta, Y1){</pre>
   dn \leftarrow d0 + t(Y1 - X%*\%Beta)%*%(Y1 - X%*\%Beta)
23
    sig2 <- invgamma::rinvgamma(1, shape = an/2, rate = dn/2)
    return(sig2)
25
26 }
27 PostBetas <- matrix(0, mcmc+burnin, k); Beta <- rep(0, k)
28 PostSigma2 <- rep(0, mcmc+burnin); sig2 <- 1</pre>
29 L <- log(1000000); U <- Inf
30 # create progress bar in case that you want to see
      iterations progress
31 pb <- winProgressBar(title = "progress bar", min = 0, max =
      tot, width = 300)
32 for(s in 1:tot){
   Yl <- sapply(1:N, function(i){PostYl(Beta = Beta, L = L, U
       = U, i)})
    Beta <- PostBeta(Y1 = Y1, sig2)</pre>
    sig2 <- PostSig2(Beta = Beta, Y1 = Y1)</pre>
    PostBetas[s,] <- Beta; PostSigma2[s] <- sig2</pre>
    setWinProgressBar(pb, s, title=paste( round(s/tot*100, 0),
        "% done"))
38 }
39 close(pb)
40 keep <- seq((burnin+1), tot, thin)
41 PosteriorBetas <- PostBetas[keep,]
42 colnames(PosteriorBetas) <- c("Intercept", "Perf", "Age", "
      Age2", "NatTeam", "Goals", "Exp", "Exp2")
43 summary (coda::mcmc(PosteriorBetas))
44 summary(coda::mcmc(PostSigma2[keep]))
```

7.9 Quantile regression

In quantile regression the location parameters vary according to the quantile of the dependent variable. Let $q_{\tau}(\boldsymbol{x}_i) = \boldsymbol{x}_i^{\top} \boldsymbol{\beta}_{\tau}$ denote the τ -th $(0 < \tau < 1)$ quantile regression function of y_i given \boldsymbol{x}_i such that $y_i = \boldsymbol{x}_i^{\top} \boldsymbol{\beta}_{\tau} + \mu_i$ where $\int_{-\infty}^{0} f_{\tau}(\mu_i) d\mu_i = \tau$, that is, the τ -th quantile of μ_i is 0.

In particular, [57] propose $f_{\tau}(\mu_i) = \tau(1-\tau) \exp\left\{-\mu_i(\tau - \mathbb{1}_{\mu_i < 0})\right\}$ (asymmetric Laplace distribution), where $\mu_i(\tau - \mathbb{1}_{\mu_i < 0})$ is the check (loss) function. These authors propose the location-scale mixture of normals with a representation given by $\mu_i = \theta e_i + \psi \sqrt{e_i} z_i$ where $\theta = \frac{1-2\tau}{\tau(1-\tau)}$, $\psi^2 = \frac{2}{\tau(1-\tau)}$, $e_i \sim E(1)$ and $z_i \sim N(0,1)$, $e_i \perp z_i$. As a consequence of this representation and the fact that the sample is i.i.d., $p(\boldsymbol{y}|\boldsymbol{\beta}_{\tau},\boldsymbol{e},\boldsymbol{X}) \propto \left(\prod_{i=1}^n e_i^{-1/2}\right) \exp\left\{-\sum_{i=1}^N \frac{(y_i - \boldsymbol{x}_i^{\top} \boldsymbol{\beta}_{\tau} - \theta e_i)^2}{2\psi^2 e_i}\right\}$.

Taking as prior a normal distribution for β_{τ} , that is, $\beta_{\tau} \sim N(\beta_{\tau 0}, \mathbf{B}_{\tau 0})$, and using data augmentation for \mathbf{e} , we can implement a Gibbs sampling algorithm in this model. The posterior distributions are

$$\boldsymbol{\beta}_{\tau}|\boldsymbol{e},\boldsymbol{y},\boldsymbol{X} \sim N(\boldsymbol{\beta}_{n\tau},\boldsymbol{B}_{n\tau}),$$

$$e_{i}|\boldsymbol{\beta}_{\tau},\boldsymbol{y},\boldsymbol{X} \sim GIG(1/2,\alpha_{ni},\delta_{ni}),^{16}$$
where $\boldsymbol{B}_{n\tau} = \left(\boldsymbol{B}_{\tau0}^{-1} + \sum_{i=1}^{N} \frac{\boldsymbol{x}_{i}\boldsymbol{x}_{i}^{\top}}{\psi^{2}e_{i}}\right)^{-1}, \boldsymbol{\beta}_{n\tau} = \boldsymbol{B}_{n\tau} \left(\boldsymbol{B}_{\tau0}^{-1}\boldsymbol{\beta}_{\tau0} + \sum_{i=1}^{N} \frac{\boldsymbol{x}_{i}(y_{i}-\theta e_{i})}{\psi^{2}e_{i}}\right),$

$$\alpha_{ni} = (y_{i} - \boldsymbol{x}_{i}^{\top}\boldsymbol{\beta}_{\tau})^{2}/\psi^{2} \text{ and } \delta_{ni} = 2 + \theta^{2}/\psi^{2}.$$

Example: The market value of soccer players in Europe continues

We continue the example of the market value of soccer players from Section 7.1. Now, we wan to know if the marginal effect of having been in the national team changes with the quantile of the market value of top soccer players in Europe. Thus, we have same regressors as in the example in the previous section, but analyze the effects in the 0.5-th an d0.9-th quantiles of the *NatTeam*.

The Algorithm A8 shows how to estimate Tobit models in our GUI. Our GUI uses the command *MCMCquantreg* from the package *MCMCpack*. The next are code shows to perform this using this package.

The results show that at the median market value, the 95% credible interval for the coefficient associated with *national team* 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. It appears that being on the national team increases the market value of more expensive players more significantly on average, although there is some overlap in the credible intervals.

 $^{^{15}}E$ denotes an exponential density.

¹⁶GIG denotes a generalized inverse Gaussian density.

Algorithm A8 Quantile regression

- 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 ${\bf R}$ syntax. You can modify the formula in the **Main equation** box using valid arguments of the *formula* command structure in ${\bf R}$
- 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

$R.\ code.\ The\ value\ of\ soccer\ player,\ quantile\ regression$

$R.\ code.\ The\ value\ of\ soccer\ player,\ quantile\ regression$

```
1 # Quantile
2 q <- 0.5
3 posterior05 <- MCMCpack::MCMCquantreg(y^X-1, tau = q, b0=b0
      , BO = BOi, burnin = burnin, mcmc = mcmc, thin = thin,
      below = 13.82, above = Inf)
4 summary(coda::mcmc(posterior05))
5 1. Empirical mean and standard deviation for each variable,
6 plus standard error of the mean:
            Mean
                       SD Naive SE Time-series SE
            7.374348 2.916758 1.304e-02
                                              2.291e-02
9 XPerf
            0.029325 0.005938 2.655e-05
                                              5.241e-05
10 XAge
            0.550633 0.241596 1.080e-03
                                              1.903e-03
11 XAge2
            -0.012027 0.004597 2.056e-05
                                              3.643e-05
12 XNatTeam 0.685275 0.170768 7.637e-04
                                              1.587e-03
            0.010608 0.002425 1.085e-05
13 XGoals
                                              1.951e-05
14 XExp
            0.092561 0.085499 3.824e-04
                                              6.799e-04
15 XExp2
           -0.002979 0.003877 1.734e-05
                                              2.941e-05
16 2. Quantiles for each variable:
            2.5%
                       25%
                                  50%
                                             75%
            1.74594
                     5.405772 7.351090
18 X
                                         9.2994982 13.216024
            0.01753 0.025340 0.029354
19 XPerf
                                          0.0333155
            0.06845 0.390780 0.553187
                                         0.7139430
                                                     1.016664
20 XAge
21 XAge2
            -0.02087 -0.015141 -0.012095 -0.0089849
                                                     -0.002813
22 XNatTeam 0.34645
                     0.572081
                               0.686735
                                         0.7996086
                                                     1.016189
            0.00578
                     0.009055
23 XGoals
                               0.010562
                                          0.0121751
                                                     0.015403
24 XExp
           -0.06761 0.034149 0.089632 0.1482128
25 XExp2
            -0.01094 -0.005456 -0.002891 -0.0004099
                                                     0.004466
26 q <- 0.9
27 posterior09 <- MCMCpack::MCMCquantreg(y^X1-1, tau = q, b0=b0
      , BO = BOi, burnin = burnin, mcmc = mcmc, thin = thin,
      below = 13.82, above = Inf)
28 summary(coda::mcmc(posterior09))
29 1. Empirical mean and standard deviation for each variable,
30 plus standard error of the mean:
                       SD Naive SE Time-series SE
31
            Mean
32 X
            8.860043 5.933902 2.654e-02
                                              6.686e-02
            0.019663 0.010241 4.580e-05
                                              1.140e-04
33 XPerf
34 XAge
            0.532823 0.483213 2.161e-03
                                              5.397e-03
35 XAge2
            -0.012328 0.008864 3.964e-05
                                              9.620e-05
36 XNatTeam
           1.033384 0.294271 1.316e-03
                                              3.389e - 03
37 XGoals
            0.008781 0.004340 1.941e-05
                                              4.991e-05
зв ХЕхр
            0.132760 0.177677 7.946e-04
                                              2.125e-03
39 XExp2
           -0.002713 0.007639 3.416e-05
                                              8.531e-05
40 2. Quantiles for each variable:
            2.5%
                       25%
                                  50%
                                            75%
                                                   97.5%
           -2.7084122 4.829341 8.821031 12.850002 20.66191
42 X
43 XPerf
            -0.0001863
                       0.012782
                                 0.019605
                                            0.026495
                                                     0.03991
44 XAge
           -0.4180422
                       0.207000
                                 0.532486
                                            0.858221
                                                      1.48632
            -0.0300400 -0.018216 -0.012235 -0.006345
45 XAge2
                                                      0.00497
46 XNatTeam
            0.4384014 0.840123
                                 1.038986 1.234456
47 XGoals
            0.0019513
                       0.005661
                                 0.008176
                                            0.011327
                                                      0.01881
48 XExp
            -0.2320608
                       0.014760
                                 0.139452
                                            0.256663
                                                      0.46053
            -0.0162717 -0.007954 -0.003198
49 XExp2
                                            0.002031
                                                      0.01385
```

7.10 Bayesian bootstrap regression

We implement the Bayesian bootstrap [90] for linear regression models. In particular, the Bayesian bootstrap simulates the posterior distributions assuming that the sample cumulative distribution function (cdf) is the population cdf (this assumption is also implicit in the frequentist bootstrap [30]).

Given $y_i \overset{i.n.d.}{\sim} \mathcal{F}$ where \mathcal{F} does not define a particular parametric family of distributions, i = 1, 2, ..., N, but sets $E(Y_i | \mathbf{x}_i) = \mathbf{x}_i^{\mathsf{T}} \boldsymbol{\beta}$, such that \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 for each y_i where the values of Y that are not observed have zero posterior probability.

The algorithm to implement the Bayesian bootstrap is the following:

Algorithm A9 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: $\boldsymbol{g} = (g_1, g_2, \dots, g_N)$ is the vector of probabilities to attach to $(y_1, \boldsymbol{x}_1^\top), (y_2, \boldsymbol{x}_2^\top), \dots, (y_n, \boldsymbol{x}_N^\top)$ for each Bayesian bootstrap replication.
- 3: Sample $(y_i, \boldsymbol{x}_i^{\top})$ N times with replacement and probabilities g_i , i = 1, 2, ..., N.
- 4: Estimate $\boldsymbol{\beta}$ using ordinary least squares in the model $E(\boldsymbol{Y}|\boldsymbol{X}) = \boldsymbol{X}\boldsymbol{\beta}$, \boldsymbol{y} being an S_1 dimensional vector of realizations of \boldsymbol{Y} , and \boldsymbol{X} an $S_1 \times k$ matrix from the previous stage.*
- 5: Repeat this process S times.
- 6: The distribution of $\beta^{(s_2)}$ is the Bayesian distribution of β .

Example: Simulation exercise

Let's do a simulation exercise to check the performance of the Algorithm A9 to perform inference using the Bayesian bootstrap. The data generating process is given by two regressors that distribute normal standard. The location vector is $\boldsymbol{\beta} = \begin{bmatrix} 1 & 1 \end{bmatrix}^{\mathsf{T}}$, and variance $\sigma^2 = 1$, the sample size is 1,000.

Algorithm A10 shows how to use our GUI to run the Bayesian bootstrap. Our GUI is based on the *bayesboot* command from *bayesboot* package in ${\bf R}$. Exercise 11 asks about using this package to perform inference in this simulation, and compared the results with the ones that we get using our GUI setting S=10000.

The following ${\bf R}$ code shows how to program the Bayesian bootstrap from scratch. We observe from the results that all 95% credible intervals encompass the population parameters, and the posterior means are close to the population parameters.

^{*}Ordinary least squares is the posterior mean of β using Jeffrey's prior in a linear regression.

Algorithm A10 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

```
= ls()); set.seed(010101)
2 N <- 1000; x1 <- runif(N); x2 <- rnorm(N)
3 X <- cbind(x1, x2); k <- dim(X)[2]</pre>
4 B <- rep(1, k+1); sig2 <- 1
        rnorm(N, 0, sig2); y <- cbind(1, X)%*%B + u</pre>
  data <- as.data.frame(cbind(y, X))</pre>
  names(data) <- c("y", "x1", "x2")
8 Reg <- function(d){</pre>
     Reg \leftarrow lm(y x1 + x2, data = d)
     Bhat <- Reg$coef
     return(Bhat)
12 }
13 S <- 10000; alpha <- 1
14 BB <- function(S, df, alpha){</pre>
     Betas <- matrix(NA, S, dim(df)[2])</pre>
     N <- dim(df)[1]</pre>
     pb <- winProgressBar(title = "progress bar", min = 0, max</pre>
       = S, width = 300)
     for(s in 1:S){
       g <- LaplacesDemon::rdirichlet(N, alpha)</pre>
19
       ids <- sample(1:N, size = N, replace = TRUE, prob = g)</pre>
       datas <- df[ids,]
21
       names(datas) <- names(df)</pre>
       Bs <- Reg(d = datas)
23
       Betas[s, ] <- Bs</pre>
       setWinProgressBar(pb, s, title=paste( round(s/S*100, 0),
        "% done"))
     close(pb)
27
     return (Betas)
29 }
30 BBs <- BB(S = S, df = data, alpha = alpha)
31 summary(coda::mcmc(BBs))
```

R. code. Bayesian bootstrap, results

```
Iterations = 1:10000
  Thinning interval = 1
  Number of chains = 1
  Sample size per chain = 10000
  1. Empirical mean and standard deviation for each variable,
  plus standard error of the mean:
                  SD Naive SE Time-series SE
        Mean
              0.06386 0.0006386
       1.1733 0.10888 0.0010888
                                      0.0010201
10 [3,] 1.0137 0.03386 0.0003386
                                      0.0003386
  2. Quantiles for each variable:
        2.5%
                25%
                       50%
                               75% 97.5%
13 var1 0.7926 0.8743 0.9169 0.9599 1.043
14 var2 0.9608 1.0984 1.1739 1.2468 1.389
15 var3 0.9473 0.9910 1.0136 1.0365 1.079
```

7.11 Summary

In this chapter, we present the core univariate regression models and demonstrate how to perform Bayesian inference using Markov Chain Monte Carlo methods. Specifically, we cover a mix of algorithms: Gibbs sampling, Metropolis-Hastings, nested M-H, and M-H-within-Gibbs. These algorithms form the foundation to perform Bayesian inference in more complex settings using cross-sectional data sets.

7.12 Exercises

- 1. Get 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. Show that the posterior conditional distributions of the Gaussian linear model with heteroskedasticity assuming independent priors $\pi(\boldsymbol{\beta}, \sigma^2, \tau) = \pi(\boldsymbol{\beta}) \times \pi(\sigma^2) \times \prod_{i=1}^N \pi(\tau_i)$, where $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \mathbf{B}_0)$, $\sigma^2 \sim IG(\alpha_0/2, \delta_0/2)$ and $\tau_i \sim G(v/2, v/2)$ are $\boldsymbol{\beta}|\sigma^2, \tau, \mathbf{y}, \mathbf{X} \sim N(\boldsymbol{\beta}_n, \mathbf{B}_n)$,

Exercises 151

$$\sigma^{2}|\boldsymbol{\beta}, \boldsymbol{\tau}, \mathbf{y}, \mathbf{X} \sim IG(\alpha_{n}, \delta_{n}) \text{ and } \tau_{i}|\boldsymbol{\beta}, \sigma^{2}, \mathbf{y}, \mathbf{X} \sim G(v_{1n}, v_{2in}), \text{ where } \tau = [\tau_{1}, \dots, \tau_{n}]^{\top}, \mathbf{B}_{n} = (\mathbf{B}_{0}^{-1} + \sigma^{-2}\mathbf{X}^{\top}\mathbf{\Psi}\mathbf{X})^{-1}, \boldsymbol{\beta}_{n} = \mathbf{B}_{n}(\mathbf{B}_{0}^{-1}\boldsymbol{\beta}_{0} + \sigma^{-2}\mathbf{X}^{\top}\mathbf{\Psi}\mathbf{y}), \alpha_{n} = \alpha_{0} + N, \delta_{n} = \delta_{0} + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^{\top}\mathbf{\Psi}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}), v_{1n} = v + 1, v_{2in} = v + \sigma^{-2}(y_{i} - \mathbf{x}_{i}^{\top}\boldsymbol{\beta})^{2}, \text{ and } \mathbf{\Psi} = \text{diagonal } \{\tau_{i}\}.$$

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 of the the market value of soccer players in Europe setting v=5 and same other hyperparameters as the homoscedastic case. Is there any meaningful difference for the coefficient associated with the national team compared to the application in the homoscedastic case?

4. Example: Determinants of hospitalization continues

Program a Gibbs sampling algorithm in the application of determinants of hospitalization.

5. Choice of the fishing mode continues

- •Run the Algorithm A3 of the book to show the results of the Geweke [41], Raftery [82] and Heidelberger [48] tests using our GUI.
- •Use the command rmnpGibbs to do the example of the choice of the fishing mode.

6. Simulation exercise of the multinomial logit model continues

Perform inference in the simulation of the multinomial logit model using the command rmnlIndepMetrop from the bayesm package of \mathbf{R} and using our GUI.

7. Simulation of the ordered probit model

Simulate an ordered probit model where the first regressor distributes N(6,5) and the second distributes G(1,1), the location parameter is $\boldsymbol{\beta} = \begin{bmatrix} 0.5 & -0.25 & 0.5 \end{bmatrix}^{\mathsf{T}}$, and the cutoffs is the vector $\boldsymbol{\alpha} = \begin{bmatrix} 0 & 1 & 2.5 \end{bmatrix}^{\mathsf{T}}$. Program from scratch a Metropolis-within-Gibbs sampling algorithm to perform inference in this simulation.

8. Simulation of the negative binomial model continues

Perform inference in the simulation of the negative binomial model using the bayesm package in \mathbf{R} software.

9. The market value of soccer players in Europe continues

Perform the application of the value of soccer players with left censuring at one million Euros in our GUI using the Algorithm A7, and the hyperparameters of the example.

10. The market value of soccer players in Europe continues

Program from scratch the Gibbs sampling algorithm in the example of the market value of soccer players at the 0.75 quantile.

11. Use the *bayesboot* package to perform inference in the simulation exercise of Section 7.10, and compared the results with the ones that we get using our GUI setting S=10000.

Multivariate models

We describe how to perform Bayesian inference in multivariate response models: multivariate regression, seemingly unrelated regression, instrumental variables, and multivariate probit model. In particular, we show the posterior distributions of the parameters, and perform some applications and simulations. Again, we show how to perform inference in these models using three levels of programming skills: GUI, packages, and programming from scratch the algorithms. Finally, there are some mathematical and computational exercises.

Remember that we can run our GUI typing

```
R code. How to display our graphical user interface

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

in the \mathbf{R} package console or any \mathbf{R} code editor.

8.1 Multivariate regression

A complete presentation of this model is given in Section 4.4. We show here the setting, and the posterior distributions for facility in exposition. In particular, there are M multiply dependent variables which share the same set of regressors, and their stochastic errors are contemporaneously correlated. In particular, $\mathbf{Y} = [\mathbf{y_1}, \mathbf{y_2}, \dots, \mathbf{y_M}]$ is an $N \times M$ matrix that is generated by $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{U}$ where \mathbf{X} is an $N \times K$ matrix of regressors, $\mathbf{B} = [\beta_1 \ \beta_2 \dots \beta_M]$ is a $K \times M$ matrix of parameters, and $\mathbf{U} = [\mu_1 \ \mu_2 \dots \mu_M]$ is a matrix of stochastic random errors such that $\mu_i \sim N(\mathbf{0}, \mathbf{\Sigma}), i = 1, 2, \dots, N$ is each row of \mathbf{U} .

The prior is given by $\pi(B|\Sigma) \sim N(B_0, V_0, \Sigma)$ and $\pi(\Sigma) \sim IW(\Psi_0, \alpha_0)$.

Therefore, the conditional posterior distributions are

$$\boldsymbol{B}|\boldsymbol{\Sigma},\boldsymbol{Y},\boldsymbol{X}\sim N(\boldsymbol{B}_n,\boldsymbol{V}_n,\boldsymbol{\Sigma}),$$

$$\Sigma | Y, X \sim IW(\Psi_n, \alpha_n),$$

where $V_n = (X^\top X + V_0^{-1})^{-1}$, $B_n = V_n(V_0^{-1}B_0 + X^\top X\hat{B})$, $\hat{B} = (X^\top X)^{-1}X^\top Y$, $\Psi_n = \Psi_0 + \mathbf{S} + \mathbf{B}_0^\top V_0^{-1}\mathbf{B}_0 + \hat{\mathbf{B}}^\top X^\top X\hat{\mathbf{B}} - \mathbf{B}_n^\top V_n^{-1}\mathbf{B}_n$, and $\alpha_n = \alpha_0 + N$. We can use a Gibbs sampling algorithm in this model since the conditional posterior distributions are standard.

Example: The effect of institutions on per capita gross domestic product

To illustrate multivariate regression models, we used the dataset provided by [1], who analyzed the effect of property rights on economic growth.

Let's assume that the point of departure is the following *simultaneous* structural economic model:¹

$$\log(\text{pcGDP95}_i) = \beta_1 + \beta_2 \text{PAER}_i + \beta_3 \text{Africa} + \beta_4 \text{Asia} + \beta_5 \text{Other} + u_{1i},$$
(8.1)

$$PAER_i = \alpha_1 + \alpha_2 \log(pcGDP95_i) + \alpha_3 \log(Mort_i) + u_{2i}, \qquad (8.2)$$

where pcGDP95, PAER and Mort are the per capita gross domestic product (GDP) in 1995, the average index of protection against expropriation between 1985 and 1995, and the settler mortality rate during the time of colonization, respectively. Africa, Asia and Other are dummies for continents, with America as the baseline group.

In this model, there is reverse (simultaneous) causality due to the contemporaneous effect of GDP on PAER, and vice verse. Therefore, estimation of the Equations 8.1 and 8.2 without taking into account this phenomenon generates posterior mean estimates that are biased and inconsistent from a sampling (frequentist) point of view. A potential strategy to tackle this issue is to estimate the reduced-form model, that is, a model without simultaneous causality where all endogenous variables are function of exogenous variables. The former variables are determined within the model (log(pcGDP95_i) and PAER in this example), and the latter are determined outside the model (log(Mort_i), Africa, Asia, and Other in this example).

¹This is a model that captures the potential underlying economic relationship between the variables.

 $^{^2}$ Simultaneous causality is the most controversial causation issue from a philosophy of science perspective. The root of the issue is that causation is typically based on the time sequence of cause and effect.

³Observe that $\mathbb{E}[u_1\text{PAER}] \neq 0$, which means failing to meet a necessary requirement to get *unbiased* and *consistent* estimators of β . See exercise 1.

Replacing Equation 8.2 into Equation 8.1, and solving for $\log(pcGDP95)$,

$$\log(\text{pcGDP95}_i) = \pi_1 + \pi_2 \log(\text{Mort}_i) + \pi_3 \text{Africa} + \pi_4 \text{Asia} + \pi_5 \text{Other} + e_{1i}.$$
(8.3)

Then, replacing Equation 8.3 into Equation 8.2, and solving for PAER,

$$PAER_{i} = \gamma_{1} + \gamma_{2} \log(Mort_{i}) + \gamma_{3} Africa + \gamma_{4} Asia + \gamma_{5} Other + e_{2i}, \quad (8.4)$$

where $\pi_2 = \frac{\beta_2 \alpha_3}{1 - \beta_2 \alpha_2}$ and $\gamma_2 = \frac{\alpha_3}{1 - \beta_2 \alpha_2}$ given $\beta_2 \alpha_2 \neq 1$, that is, independent equations (see Exercise 2).

Observe that equations 8.3 and 8.4 have the form of a multivariate regression model where the common set of regressors is $X = [\log(\text{Mort}) \text{ Africa Asia Other}]$ and $Y = [\log(\text{pcGDP95}) \text{ PAER}]$. Thus, we can estimate this model using the setup of this section.

Thus, we estimate in a first stage the parameters from the *reduced-form* model (Equations 8.3 and 8.4), but the main interest is the parameters of the *structural* model (Equations 8.1 and 8.2). Thus, a valid question is if we can recover the *structural* parameters from the *reduced-form* parameters. There are two criteria to respond this question: the order condition, which is necessary, and the rank condition, which is necessary and sufficient.

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 the equation (exclusion restrictions). The equation is exactly identified if the number of excluded variables is M-1, and is over identified if the number of excluded variables is greater than M-1.
- The parameters of equation m in the system are identified if $K K_m \ge M_m 1$, where K_m and M_m are 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 if $K K_m > M_m 1$.

We can see from Equations 8.1 and 8.2 in this example that K=5, M=2, $K_1=4$, $K_2=2$, $M_1=2$ and $M_2=2$. This means that $K-K_1=1=M-1$ and $K-K_2=3>M-1=1$, that is, the order condition says that both equations satisfy the necessary condition of identification, the first equation would be *exactly identified*, and the second equation would be *over identified*. Observe that there is one excluded variable from the first equation, and there are three excluded variables from the second equation.

The rank condition

The rank condition (necessary and sufficient) says that given a *structural* model with M equations (M endogenous variables), an equation is identified if and only if there is at least one determinant different from zero from a $(M-1)\times (M-1)$ matrix built using the excluded variables in the analyzed equation, but included in any other equation of the system.

It is useful to build the *identification matrix* to implement the *rank* condition. Table 8.1 shows this matrix in this example.

TABLE 8.1 Identification matrix.

$\log(\text{pcGDP95})$	PAER	Constant	$\log(Mort)$	Africa	Asia	Other
1	$-\beta_2$	$-\beta_1$	0	$-\beta_3$	$-\beta_4$	$-\beta_5$
- $lpha_2$	1	$-\alpha_1$	$-lpha_3$	0	0	0

The only excluded variable in the log(pcGDP95) equation is log(Mort). Then, there is just one matrix that can be built using the excluded variables from this equation $[-\alpha_3]$ (see column 4 in Table 8.1). Thus, the determinant of this matrix is $-\alpha_3$, and as far as this coefficient is different to zero, that is, that the mortality rate is relevant in the PAER equation ($\alpha_3 \neq 0$), the coefficients in log(pcGDP95) equation are exactly identified. For instance, $\beta_2 = \frac{\pi_2}{\gamma_2}$, which is the effect of property rights on GDP, is exactly identified.

Observe the importance of excluding log(Mort) from the log(pcGDP95) equation, but including log(Mort) in the PAER equation. This is called exclusion restriction, and it is the requirement of having an exogenous source of variability in the PAER equation that helps to identify the log(pcGDP95) equation. Having relevant exogenous sources of variability is a very important aspect in identification, estimation and inference of structural parameters.

Regarding the identification of the *structural* parameters in the PAER equation, there are three potential matrices that can be constructed: $[-\beta_3]$, $[-\beta_4]$ and $[-\beta_5]$ (see columns 5, 6 and 7 in Table 8.1), as far as any of these parameters are relevant in the log(pcGDP95) equation, we achieve identification of the PAER equation. In this case, this equation is *over identified*, that is, there are many ways to find the parameters in this equations. For instance, $\alpha_2 = \gamma_3/\pi_3 = \gamma_4/\pi_4 = \gamma_5/\pi_5$ (see Exercise 2).

In general, trying to recover the *structural* parameters from the *reduced-form* parameters can be challenging due to the requirement of relevant identification restrictions that can be hard to find in some applications.⁴

We set non-informative priors in this example, $\mathbf{B}_0 = [\mathbf{0}_5 \ \mathbf{0}_5]$, $\mathbf{V}_0 = 100\mathbf{I}_K$, $\mathbf{\Psi}_0 = 5\mathbf{I}_2$ and $\alpha_0 = 5.5$ Once our GUI is displayed (see beginning of this

⁴Good text books at introductory level for identification in linear systems are [46, Chap. 19] and [107, Chap. 16].

⁵Observe that we are setting the priors in the *reduced-form* model; this may have unintended consequences for the posterior distributions of the *structural* parameters, which are ultimately the parameters researchers are interested in. See [56, p. 302] for good references in this topic.

chapter), we should follow Algorithm A11 to run multivariate linear models in our GUI (see Chapter 6 for details, particularly how to set the data set):

Algorithm A11 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 ${\bf R}$ 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.

R code. The effect of institutions on per capita GDP

```
1 rm(list = ls())
2 set.seed (12345)
3 DataInst <- read.csv("DataApplications/4Institutions.csv",</pre>
      sep = ",", header = TRUE, fileEncoding = "latin1")
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 BO <- 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))
27 summary(coda::mcmc(t(BsCond)))
28 SIGMs <- t(sapply(1:s, function(1) {gdata::lowerTriangle(
      SIGs[,,1], diag=TRUE, byrow=FALSE)}))
29 summary(coda::mcmc(SIGMs))
30 hdiBs \leftarrow HDInterval::hdi(t(BsCond), credMass = 0.95) #
      Highest posterior density credible interval
31 hdiBs
_{\rm 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(beta1)) # Effect of property rights on
      GDP
36 Iterations = 1:10000
37 Thinning interval = 1
38 Number of chains = 1
39 Sample size per chain = 10000
40 1. Empirical mean and standard deviation for each variable,
41 plus standard error of the mean:
                             Naive SE Time-series SE
42 Mean
                   SD
43 0.9796
                 16.8430
                                  0.1684
44 2. Quantiles for each variable:
          25%
                 50%
                         75% 97.5%
46 0.5604 0.7984 0.9677 1.2329 2.8709
```

8.2 Seemingly unrelated regression

In seemingly unrelated regression (SUR) models there are M dependent variables with potentially different regressors such that the stochastic errors are contemporaneously correlated. This is $\mathbf{y}_j = \mathbf{X}_j \mathbf{\beta}_j + \mathbf{\mu}_j$, where \mathbf{y}_j is a N-dimensional vector, \mathbf{X}_j is a matrix of dimension $N \times K_m$ of regressors, $\mathbf{\beta}_j$ is a K_m -dimensional vector of location parameters, and $\mathbf{\mu}_j$ is a N-dimensional vector of stochastic errors, m = 1, 2, ..., M.

vector of stochastic errors, m = 1, 2, ..., M. Setting $\boldsymbol{\mu}_i = [\mu_{i1} \ \mu_{i2} ... \mu_{iM}]^{\top}$ such that $\boldsymbol{\mu}_i \sim N(\mathbf{0}, \boldsymbol{\Sigma})$, and stacking the M equations, we can write $\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\mu}$ where $\boldsymbol{y} = [\boldsymbol{y}_1^{\top} \ \boldsymbol{y}_2^{\top} ... \boldsymbol{y}_M^{\top}]^{\top}$ is a MN-dimensional vector, $\boldsymbol{\beta} = [\boldsymbol{\beta}_1^{\top} \ \boldsymbol{\beta}_2^{\top} ... \boldsymbol{\beta}_M^{\top}]^{\top}$ is a K dimensional vector, $K = \sum_{m=1}^{M} K_m, \boldsymbol{X}$ is an $MN \times K$ block diagonal matrix composed of \boldsymbol{X}_m and $\boldsymbol{\mu} = [\boldsymbol{\mu}_1^{\top} \ \boldsymbol{\mu}_2^{\top} ..., \boldsymbol{\mu}_M^{\top}]^{\top}$ is a MN-dimensional vector of stochastic errors such that $\boldsymbol{\mu} \sim N(\mathbf{0}, \boldsymbol{\Sigma} \otimes \boldsymbol{I}_N)$. Then,

$$p(\boldsymbol{\beta}, \boldsymbol{\Sigma} | \boldsymbol{y}, \boldsymbol{X}) \propto |\boldsymbol{\Sigma}|^{-N/2} \exp \left\{ -\frac{1}{2} (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta})^{\top} (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{I}_N) (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}) \right\}.$$

Using independent priors $\pi(\boldsymbol{\beta}) \sim N(\boldsymbol{\beta}_0, \boldsymbol{B}_0)$ and $\pi(\boldsymbol{\Sigma}^{-1}) \sim W(\alpha_0, \boldsymbol{\Psi}_0)$, the posterior distributions are

$$\boldsymbol{\beta}|\boldsymbol{\Sigma},\boldsymbol{y},\boldsymbol{X}\sim N(\boldsymbol{\beta}_n,\boldsymbol{B}_n),$$

$$\Sigma^{-1}|\boldsymbol{\beta}, \boldsymbol{y}, \boldsymbol{X} \sim W(\alpha_n, \boldsymbol{\Psi}_n),$$

where $\boldsymbol{B}_n = (\boldsymbol{X}^{\top}(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{I}_N)\boldsymbol{X} + \boldsymbol{B}_0^{-1})^{-1}, \boldsymbol{\beta}_n = \boldsymbol{B}_n(\boldsymbol{B}_0^{-1}\boldsymbol{\beta}_0 + \boldsymbol{X}^{\top}(\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{I}_N)\boldsymbol{y}), \ \alpha_n = \alpha_0 + N \text{ and } \boldsymbol{\Psi}_n = (\boldsymbol{\Psi}_0^{-1} + \boldsymbol{U}^{\top}\boldsymbol{U})^{-1}, \text{ where } \boldsymbol{U} \text{ is an } N \times M \text{ matrix whose columns are } \boldsymbol{y}_j - \boldsymbol{X}_j\boldsymbol{\beta}_j.$

Observe that we have standard conditional posteriors, thus, we can employ a Gibbs sampling algorithm to get the posterior draws.

Example: Utility demand

Let's use the dataset *Utilities.csv* to estimate a seemingly unrelated regression model for utilities. We use the same setting as in Exercise 14 in Chapter 4 where we ask to estimate a multivariate regression model omitting households with no consumption in any utility. We see in this exercise that no all regressors are relevant for the demand of electricity, water and gas. Thus, we

estimate the following model:

```
\begin{split} \log(\text{electricity}_i) &= \beta_1 + \beta_2 \log(\text{electricity price}_i) + \beta_3 \log(\text{water price}_i) \\ &+ \beta_4 \log(\text{gas price}_i) + \beta_5 \text{IndSocio1}_i + \beta_6 \text{IndSocio2}_i + \beta_7 \text{Altitude}_i \\ &+ \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) \\ &+ \alpha_4 \log(\text{gas price}_i) + \alpha_5 \text{IndSocio1}_i + \alpha_6 \text{IndSocio2}_i \\ &+ \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) \\ &+ \gamma_4 \log(\text{gas price}_i) + \gamma_5 \text{IndSocio1}_i + \gamma_6 \text{IndSocio2}_i + \gamma_7 \text{Altitude}_i \\ &+ \gamma_8 \text{Nrooms}_i + \gamma_9 \text{HouseholdMem}_i + \mu_{i3}, \end{split}
```

where electricity, water and gas are the monthly consumption of electricity (kWh), water (m³) and gas (m³) of Colombian households. There is information of 2103 households regarding average prices of electricity (USD/kWh), water (USD/m³) and gas (USD/m³), indicators of socioeconomic conditions of the neighborhood where the household is located (IndSocio1 is the lowest and IndSocio3 is the highest), an indicator if the household is located in a municipality that is above 1000 meters above the sea level, the number of rooms in the house, the number of members of the households, and monthly income (USD).

Since there are different sets of regressors in each equation and we suspect correlation between the stochastic errors of the three equations, we should estimate a seemingly unrelated regressions (SUR) model. We expect unobserved correlation in these equations because we are modelling utilities, and in some cases, a single provider handles all three services and issues one bill.

Algorithm A12 shows how to estimate SUR models in our GUI. Our GUI uses the command *rsurGibbs* from the *bayesm* package in **R** software. See Chapter 6 for details, particularly how to set the data set, and templates in our GitHub repository (https://github.com/besmarter/BSTApp) in the folders **DataApp** and **DataSim**.

The following code shows how to program this application using this package. We use 10000 MCMC iterations, $\beta_0 = \mathbf{0}_{27}$, $B_0 = 100 I_{27}$, $\alpha_0 = 5$ and $\Psi = 5 I_3$.

We find that the posterior median estimates of the own-price elasticities of demand of electricity, water and gas are -1.88, -0.36 and -0.62, where there are not 95% credible intervals that encompass 0. This means that a 1% increase in the prices of electricity, water and gas imply a 1.88%, 0.36% and 0.62% decrease in the monthly consumption of these utilities, respectively. In general, there is evidence supporting the relevance of all regressors in these equations,

⁶This is an example where there can be concerns regarding biased and inconsistent posterior mean estimates, for instance, due to reverse causality between quantity and demand. These concerns are valid; although, we are using micro-level data, which implies no demand-supply simultaneity. In addition, the utility providers are operating in regulated

except a few exceptions, and unobserved correlation in the demand of these services supporting the relevance of a SUR model in this application.

Algorithm A12 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

natural monopoly markets, this implies no endogeneity due to searching provider strategies. Finally, we took prices directly from provider records, this avoids price measurement errors [84].

Multivariate models

R code. Utility demand in Colombia

```
1 rm(list = ls())
2 set.seed(010101)
3 library(dplyr)
4 DataUt <- read.csv("DataApplications/Utilities.csv", sep = "
      ,", header = TRUE, fileEncoding = "latin1")
5 DataUtEst <- DataUt %>%
6 filter(Electricity != 0 & Water !=0 & Gas != 0)
7 attach(DataUtEst)
8 y1 <- log(Electricity); y2 <- log(Water); y3 <- log(Gas)</pre>
9 X1 <- cbind(1, LnPriceElect, LnPriceWater, LnPriceGas,
      IndSocio1, IndSocio2, Altitude, Nrooms, HouseholdMem,
      Lnincome)
10 X2 <- cbind(1, LnPriceElect, LnPriceWater, LnPriceGas,
      IndSocio1, IndSocio2, Nrooms, HouseholdMem)
11 X3 <- cbind(1, LnPriceElect, LnPriceWater, LnPriceGas,
      IndSocio1, IndSocio2, Altitude, Nrooms, HouseholdMem)
12 regdata <- NULL
13 regdata[[1]] <- list(y = y1, X = X1); regdata[[2]] <- list(y</pre>
       = y2, X = X2); regdata[[3]] <- list(y = y3, X = X3)
14 M <- length(regdata); K1 <- dim(X1)[2]; K2 <- dim(X2)[2]; K3
       \leftarrow dim(X3)[2]
15 K <- K1 + K2 + K3
16 # Hyperparameters
17 b0 <- rep(0, K); c0 <- 100; B0 <- c0*diag(K); V <- 5*diag(M)
      ; a0 <- M
18 Prior <- list(betabar = b0, A = solve(B0), nu = a0, V = V)
19 #Posterior draws
20 S <- 10000; keep <- 1; Mcmc <- list(R = S, keep = keep)
21 PosteriorDraws <- bayesm::rsurGibbs(Data = list(regdata =</pre>
      regdata), Mcmc = Mcmc, Prior = Prior)
```

R code. Utility demand in Colombia, results

```
1 Bs <- PosteriorDraws[["betadraw"]]</pre>
2 Names <- c("Const", "LnPriceElect", "LnPriceWater", "</pre>
      LnPriceGas", "IndSocio1", "IndSocio2",
3 "Altitude", "Nrooms", "HouseholdMem", "Lnincome", "Const",
  "LnPriceElect", "LnPriceWater", "LnPriceGas", "IndSocio1", "
      IndSocio2",
5 "Nrooms", "HouseholdMem", "Const",
6 "LnPriceElect", "LnPriceWater", "LnPriceGas", "IndSocio1", "
      IndSocio2".
7 "Altitude", "Nrooms", "HouseholdMem")
8 colnames(Bs) <- Names</pre>
9 summary(coda::mcmc(Bs))
summary(PosteriorDraws[["Sigmadraw"]])
11 2. Quantiles for each variable:
            2.5%
                       25%
                                          75%
                0.44452 1.03120 1.342407
                                            1.65192 2.25376
13 Const
14 LnPriceElect -2.39679 -2.06328 -1.882706 -1.70369 -1.36996
15 LnPriceWater -0.44221 -0.38678 -0.356850 -0.32669 -0.26969
               -0.21655 -0.13777 -0.098191 -0.05902 0.01872
16 LnPriceGas
17 IndSocio1
               -0.87630 -0.78653 -0.737701 -0.68840 -0.59675
               -0.24601 -0.18286 -0.151440 -0.11896 -0.05681
18 IndSocio2
                -0.27080 -0.23838 -0.220742 -0.20385 -0.17259
19 Altitude
                                            0.07835 0.09422
                0.04596 0.06178 0.070023
20 Nrooms
21 HouseholdMem 0.06600 0.07994
                                   0.086857
                                             0.09411
                                                      0.10785
22 Lnincome
                0.03836 0.05421
                                   0.062957
                                             0.07165
                                                      0.08717
                0.88957
                         1.73496
23 Const
                                  2.169638
                                             2.62170
                                                      3.47216
24 LnPriceElect -0.81956 -0.31624 -0.054075
                                             0.21132
25 LnPriceWater -0.49559 -0.40995 -0.364248
                                            -0.32026 -0.23639
26 LnPriceGas
                0.06075 0.16754 0.226690
                                            0.28570
                                                      0.39476
27 IndSocio1
               -0.64203 -0.50302 -0.427819 -0.35226 -0.21315
                -0.50401 -0.40949 -0.359821 -0.31063 -0.21199
28 IndSocio2
                0.05688 0.08023
                                   0.093139
                                             0.10555
29 Nrooms
                                                      0.12968
30 HouseholdMem 0.10041 0.12065 0.131506
                                             0.14260
                                                      0.16314
31 Const
                -2.28569 -1.58566 -1.220078 -0.84612 -0.14787
32 LnPriceElect -2.42484 -2.01228 -1.797269 -1.57889 -1.16396
33 LnPriceWater -0.10684 -0.03923 -0.004088
                                            0.03153 0.09905
34 LnPriceGas
                -0.76526 -0.67445 -0.625899 -0.57734 -0.48125
                -0.91381 -0.80243 -0.744909 -0.68577 -0.57341
35 IndSocio1
               -0.31791 -0.24388 -0.203300 -0.16415 -0.09012
36 IndSocio2
                         0.29099
37 Altitude
                0.24896
                                   0.311668
                                             0.33256
                                                      0.37278
38 Nrooms
                0.06050
                         0.07921
                                   0.089386
                                             0.09943
                                                       0.11793
39 HouseholdMem 0.14467
                         0.16144
                                   0.170024
                                             0.17843
                                                       0.19431
40 summary(coda::mcmc(PosteriorDraws[["Sigmadraw"]]))
41 2. Quantiles for each variable:
          2.5%
                   25%
                            50%
                                    75% 97.5%
43 var1 0.19912 0.20822 0.21332 0.21863 0.2290
44 var2 0.08183 0.09284 0.09870 0.10475 0.1160
45 var3 0.05121 0.05973 0.06426 0.06882 0.0781
46 var4 0.08183 0.09284 0.09870 0.10475 0.1160
47 var5 0.47763 0.49934 0.51131 0.52387 0.5493
48 var6 0.07318 0.08653 0.09351 0.10079 0.1145
49 var7 0.05121 0.05973 0.06426 0.06882 0.0781
50 var8 0.07318 0.08653 0.09351 0.10079 0.1145
51 var9 0.29523 0.30900 0.31654 0.32428 0.3397
```

We ask in the Exercise 5 to run this application using our GUI and the information in the dataset *Utilities.csv*. Observe that this file should be modified to agree the structure that requires our GUI (see the dataset *5Institutions.csv* in the folder *DataApp* of our GitHub repository - https://github.com/besmarter/BSTApp- for a template). In addition, we ask to program from scratch the Gibbs sampler algorithm in this application.

8.3 Instrumental variable

This inferential approach is used when there are *endogeneity* issues, that is, the stochastic error is not independent of the regressors, this in turn generates bias in posterior mean estimates when we use an inferential approach that does not take this issue into. *Endogeneity* can be caused by reverse causality, omitting relevant correlated variables, or measurement error in the regressors.⁷

Let's specify the dependent variable as a linear function of one endogenous regressor and some exogenous regressors. That is, $y_i = \boldsymbol{x}_{ei}^{\top} \boldsymbol{\beta}_1 + \boldsymbol{\beta}_s x_{si} + \mu_i$ where $x_{si} = \boldsymbol{x}_{ei}^{\top} \boldsymbol{\gamma}_1 + \boldsymbol{z}_i^{\top} \boldsymbol{\gamma}_2 + v_i$, x_s is the variable which generates the endogeneity issues ($\mathbb{E}[\mu|x_s] \neq 0$), x_e are K_1 exogenous regressors ($\mathbb{E}[\mu|x_e] = \mathbf{0}$), and \boldsymbol{z} are K_2 instruments, that is, regressors that drive x_s ($\mathbb{E}[x_s \boldsymbol{z}] \neq \mathbf{0}$), but do not have a direct effect on y ($\mathbb{E}[y\boldsymbol{z}|x_s] = \mathbf{0}$). The equation of y is called the *structural equation*, and is the equation that the reseracher is interested in.

Assuming $(u_i, v_i)^{\top} \stackrel{i.i.d.}{\sim} N(0, \Sigma), \Sigma = [\sigma_{lm}], l, m = 1, 2$, the likelihood function is

$$p(\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\Sigma} | \boldsymbol{y}, \boldsymbol{X}, \boldsymbol{Z}) = \frac{1}{(2\pi)^{\frac{N}{2}} |\boldsymbol{\Sigma}|^{\frac{N}{2}}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^{N} (y_i - \boldsymbol{x}_i^{\top} \boldsymbol{\beta}, x_{si} - \boldsymbol{w}_i^{\top} \boldsymbol{\gamma}) \boldsymbol{\Sigma}^{-1} \begin{pmatrix} y_i - \boldsymbol{x}_i^{\top} \boldsymbol{\beta} \\ x_{si} - \boldsymbol{w}_i^{\top} \boldsymbol{\gamma} \end{pmatrix} \right\},$$

where
$$\boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1^\top, \boldsymbol{\beta}_s \end{bmatrix}^\top$$
, $\boldsymbol{\gamma} = \begin{bmatrix} \boldsymbol{\gamma}_1^\top, \boldsymbol{\gamma}_2^\top \end{bmatrix}^\top$, $\boldsymbol{x}_i = \begin{bmatrix} \boldsymbol{x}_{ei}^\top, x_{si} \end{bmatrix}^\top$ and $\boldsymbol{w}_i = \begin{bmatrix} \boldsymbol{x}_{ei}^\top, \boldsymbol{z}_i^\top \end{bmatrix}^\top$. We get standard conditional posterior densities using the following inde-

We get standard conditional posterior densities using the following independent priors $\gamma \sim N(\gamma_0, G_0)$, $\beta \sim N(\beta_0, B_0)$ and $\Sigma^{-1} \sim W(\alpha_0, \Psi_0)$. In particular,

$$\begin{split} \boldsymbol{\beta}|\boldsymbol{\gamma}, \boldsymbol{\Sigma}, \boldsymbol{y}, \boldsymbol{X}, \boldsymbol{Z} \sim N(\boldsymbol{\beta}_n, \boldsymbol{B}_n) \\ \boldsymbol{\gamma}|\boldsymbol{\beta}, \boldsymbol{\Sigma}, \boldsymbol{y}, \boldsymbol{X}, \boldsymbol{Z} \sim N(\boldsymbol{\gamma}_n, \boldsymbol{G}_n) \\ \boldsymbol{\Sigma}^{-1}|\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{y}, \boldsymbol{X}, \boldsymbol{Z} \sim W(\boldsymbol{\alpha}_n, \boldsymbol{\Psi}_n) \end{split}$$
 where $\boldsymbol{\beta}_n = \boldsymbol{B}_n \left(\boldsymbol{B}_0^{-1} \boldsymbol{\beta}_0 + \omega_1^{-1} \sum_{i=1}^N \left[\boldsymbol{x}_i \left(y_i - \frac{\sigma_{12}(\boldsymbol{x}_{si} - \boldsymbol{w}_i^{\top} \boldsymbol{\gamma})}{\sigma_{22}} \right) \right] \right), \ \boldsymbol{B}_n = \boldsymbol{B}_n \left(\boldsymbol{B}_0^{-1} \boldsymbol{\beta}_0 + \omega_1^{-1} \sum_{i=1}^N \left[\boldsymbol{x}_i \left(y_i - \frac{\sigma_{12}(\boldsymbol{x}_{si} - \boldsymbol{w}_i^{\top} \boldsymbol{\gamma})}{\sigma_{22}} \right) \right] \right), \ \boldsymbol{B}_n = \boldsymbol{B}_n \left(\boldsymbol{B}_0^{-1} \boldsymbol{\beta}_0 + \boldsymbol{\omega}_1^{-1} \sum_{i=1}^N \left[\boldsymbol{x}_i \left(\boldsymbol{y}_i - \frac{\sigma_{12}(\boldsymbol{x}_{si} - \boldsymbol{w}_i^{\top} \boldsymbol{\gamma})}{\sigma_{22}} \right) \right] \right), \ \boldsymbol{B}_n = \boldsymbol{B}_n \left(\boldsymbol{B}_0^{-1} \boldsymbol{\beta}_0 + \boldsymbol{\omega}_1^{-1} \sum_{i=1}^N \left[\boldsymbol{x}_i \left(\boldsymbol{y}_i - \frac{\sigma_{12}(\boldsymbol{x}_{si} - \boldsymbol{w}_i^{\top} \boldsymbol{\gamma})}{\sigma_{22}} \right) \right] \right)$

 $^{^7{\}rm See}$ [107, Chap. 15] for an introductory treatment of instrumental variable in the Frequentist inferential approach.

$$\begin{aligned} &(\omega_{1}^{-1} \sum_{i=1}^{N} \boldsymbol{x}_{i} \boldsymbol{x}_{i}^{\top} + \boldsymbol{B}_{0}^{-1})^{-1}, \ \omega_{1} = \sigma_{11} - \sigma_{12}^{2} / \sigma_{22}, \ \boldsymbol{G}_{n} = (\omega_{2}^{-1} \sum_{i=1}^{N} \boldsymbol{w}_{i} \boldsymbol{w}_{i}^{\top} + \boldsymbol{G}_{0}^{-1})^{-1}, \ \boldsymbol{\gamma}_{n} = \boldsymbol{G}_{n} \left(\boldsymbol{G}_{0}^{-1} \boldsymbol{\gamma}_{0} + \omega_{2}^{-1} \sum_{i=1}^{N} \left[\boldsymbol{w}_{i} \left(\boldsymbol{x}_{si} - \frac{\sigma_{12} (y_{i} - \boldsymbol{x}_{i}^{\top} \boldsymbol{\beta})}{\sigma_{11}} \right) \right] \right), \ \omega_{2} = \sigma_{22} - \sigma_{12}^{2} / \sigma_{11}, \ \boldsymbol{\Psi}_{n} = \left[\boldsymbol{\Psi}_{0}^{-1} + \sum_{i=1}^{N} \left(\boldsymbol{y}_{i} - \boldsymbol{x}_{i}^{\top} \boldsymbol{\beta} \right) \left(\boldsymbol{y}_{i} - \boldsymbol{x}_{i}^{\top} \boldsymbol{\beta}, \boldsymbol{x}_{si} - \boldsymbol{w}_{i}^{\top} \boldsymbol{\gamma} \right) \right]^{-1}, \\ \boldsymbol{\alpha}_{n} = \boldsymbol{\alpha}_{0} + N, \ \text{and} \ \boldsymbol{\sigma}_{lj} \ \text{are the elements of} \ \boldsymbol{\Sigma}. \end{aligned}$$

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 $[\mu_i \ v_i]^{\top} \sim N(\mathbf{0}, \mathbf{\Sigma}), \ \mathbf{\Sigma} = [\sigma_{lj}]$ 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 $\mathbb{E}[\mu|x_s] = \mathbb{E}[\mu|v] = \frac{\sigma_{12}}{\sigma_{22}}v \neq 0$ given $\sigma_{12} \neq 0$ and $\mathbb{E}[\mu|z] = 0$. Let's set all location parameters equal to 1, and $\sigma_{11} = \sigma_{22} = 1$, $\sigma_{12} = 0.8$, and $z \sim N(0,1)$. We know from the large sampling properties of the posterior mean that this converge to the maximum likelihood estimator (see Section 1.1, and [61, 103]), which in this setting is $\hat{\beta}_2 = \frac{\widehat{Cov}(x_s, y)}{\widehat{Var}(x_s)}$ which converges in probability to $\beta_2 + \frac{\sigma_{12}}{\sigma_{22}Var(x_s)} = \beta_2 + \frac{\sigma_{12}}{\sigma_{22}(\gamma_2^2Var(z) + \sigma_{22})} = 1.4$, that is, the asymptotic bias when using the posterior mean of a linear regression without taking into account endogeneity is 0.4 in this example.

We assess the sampling performance of Bayesian "estimators" simulating this setting 100 times. The following code shows how to do this using a linear model without taking into account the *endogeneity* issue (see Section 7.1), and implementing the variable instrumental model. We use $\boldsymbol{B}_0 = 1000\boldsymbol{I}_2$, $\boldsymbol{\beta}_0 = \boldsymbol{0}_2$, and the parameters of the inverse gamma distribution equal to 0.0005. In the case of the instrumental variable setting, we set $\boldsymbol{\gamma}_0 = \boldsymbol{0}_2$, $\boldsymbol{G}_0 = 1000\boldsymbol{I}_2$ $\alpha_0 = 3$ and $\boldsymbol{\Psi}_0 = 3\boldsymbol{I}_2$ in addition.

R code. Simulation exercise, sampling properties ordinary and instrumental models

```
1 rm(list = ls()); set.seed(010101)
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 <-</pre>
6 U <- replicate(S, MASS::mvrnorm(n = N, mu = rep(0, 2), SIGMA
      ))
7 x <- G[1] + G[2]*z + U[,2,]; y <- B[1] + B[2]*x + U[,1,]
8 # Hyperparameters
9 d0 <- 0.001/2; a0 <- 0.001/2
10 b0 <- rep(0, k); c0 <- 1000; B0 <- c0*diag(k)
11 B0i <- solve(B0); g0 <- rep(0, 2)</pre>
12 GO <- 1000*diag(2); GOi <- solve(GO)
13 nu <- 3; Psi0 <- nu*diag(2)
_{14} # MCMC parameters
15 mcmc <- 5000; burnin <- 1000
16 tot <- mcmc + burnin; thin <- 1
17 # Gibbs sampling
18 Gibbs <- function(x, y){
    Data \leftarrow list(y = y, x = x, w = w, z = Z)
    Mcmc \leftarrow list(R = mcmc, keep = thin, nprint = 0)
    Prior <- list(md = g0, Ad = G0i, mbg = b0, Abg = B0i, nu =
21
       nu, V = Psi0)
    RestIV <- bayesm::rivGibbs(Data = Data, Mcmc = Mcmc, Prior</pre>
        = Prior)
    PostBIV <- mean(RestIV[["betadraw"]])</pre>
23
    ResLM <- MCMCpack::MCMCregress(y \tilde{x} + w - 1, b0 = b0, B0
24
      = B0i, c0 = a0, d0 = d0)
    PostB <- mean(ResLM[,1]); Res <- c(PostB,PostBIV)</pre>
25
    return (Res)
26
27 }
28 PosteriorMeans <- sapply(1:S, function(s) {Gibbs(x = x[,s],
      y = y[,s])
29 rowMeans(PosteriorMeans)
30 Model <- c(replicate(S, "Ordinary"), replicate(S, "</pre>
      Instrumental"))
31 postmeans <- c(t(PosteriorMeans))</pre>
32 df <- data.frame(postmeans, Model, stringsAsFactors = FALSE)
33 library(ggplot2); library(latex2exp)
34 histExo \leftarrow ggplot(df, aes(x = postmeans, fill = Model)) +
      geom_histogram(bins = 40, position = "identity", color =
"black", alpha = 0.5) + labs(title = "Overlayed")
      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")
35 histExo
```

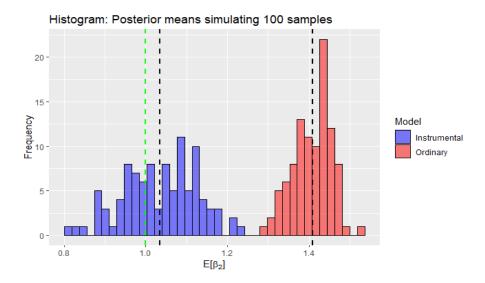


FIGURE 8.1 Histogram of posterior means: Ordinary and instrumental models.

Figure 8.1 displays the histograms of the posterior means of β_2 using the ordinary model without taking endogeneity into account, and the instrumental variable model. In one hand, the mean of the posterior means of the ordinary model is 1.41 (black dashed line in red histogram), this implies a bias equal to 0.41, which is very close to the population bias (0.40). On the other hand, the mean of the posterior means of the instrumental variable model is 1.04 (black dashed line in blue histogram), which is close to the population value of $\beta_2 = 1$ (green dashed line).

We also see that the histogram of the posterior means of the ordinary model is less disperse, that is, this "estimator" is more efficient, which is a well-known result in the Frequentist inferential approach comparing ordinary least squares and two-stage least squares (see [106, Chap. 5]).

Two very relevant aspects in the instrumental variables literature are the weakness and exogeneity of the instruments. The former refers how strong is the relationship between the instruments and the endogeneous regressors, and the latter refers to the independence of the instruments of the stochastic error in the structural equation. We ask in Exercise 6 to use the previous code as a baseline to study this two aspects. Observe the link between the weakness and exogeneity of the instrument, and the exclusion restrictions ($\mathbb{E}[x_s \mathbf{z}] \neq \mathbf{0}$ and $\mathbb{E}[y\mathbf{z}|x_s] = \mathbf{0}$). This is the point of departure of [24] who propose to assess the plausibility of the exclusion restrictions defining plausible exogeneity as having prior information that the effect of the instrument in the structural equation is near zero, but perhaps not exactly zero.

Algorithm A13 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 A13 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

8.4 Multivariate probit model

In the multivariate probit model [29], the response variable $y_{il} = \{0, 1\}$ indicates that individual i makes binary choices regarding no mutually exclusive alternatives l, i = 1, 2, ..., N, l = 1, 2, ..., L. In particular,

$$y_{il} = \begin{cases} 0, \ y_{il}^* \le 0 \\ 1, \ y_{il}^* > 0 \end{cases},$$

where $\boldsymbol{y}_i^* = \boldsymbol{X}_i \boldsymbol{\beta} + \boldsymbol{\mu}_i \overset{i.i.d.}{\sim} N(\boldsymbol{0}, \boldsymbol{\Sigma}), \, \boldsymbol{y}_i^*$ is an unobserved latent L-dimensional vector, \boldsymbol{X}_i is an $L \times K$ design matrix of regressors, $K = L \times k, \, k$ is the number of regressors, and $\boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1^\top \ \boldsymbol{\beta}_2^\top \dots \boldsymbol{\beta}_k^\top \end{bmatrix}^\top$, where the $\boldsymbol{\beta}_j$ make up an L-dimensional vector of coefficients, $j = 1, 2, \dots, k$. We simultaneously take

into account the alternative-varying regressors (alternative attributes) and alternative-invariant regressors (individual characteristics).

The likelihood function in this model is $p(\boldsymbol{\beta}, \boldsymbol{\Sigma}|\boldsymbol{y}, \boldsymbol{X}) = \prod_{i=1}^{N} \prod_{l=1}^{L} p_{il}^{y_{il}}$ where $p_{il} = p(y_{il}^* \geq 0)$. Observe that $p(y_{il}^* \geq 0) = p(\lambda_{ll}y_{il}^* \geq 0)$, $\lambda_{ll} > 0$. This generates identification issues because just the correlation matrix can be identified, same case as the univariate probit model where the variance of the model is fixed to 1. We follow the post processing strategy proposed by [29] to get identified parameters, that is, $\tilde{\boldsymbol{\beta}} = vec\{\boldsymbol{\Lambda}\boldsymbol{B}\}$ and the correlation matrix $\boldsymbol{R} = \boldsymbol{\Lambda}\boldsymbol{\Sigma}\boldsymbol{\Lambda}$, where $\boldsymbol{\Lambda} = diag\{\sigma_{ll}\}^{-1/2}$ and $\boldsymbol{B} = [\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_k]$.⁸

We assume independent priors, $\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \boldsymbol{B}_0)$ and $\boldsymbol{\Sigma}^{-1} \sim W(\alpha_0, \boldsymbol{\Psi}_0)$. We can employ Gibbs sampling in this model because this is a standard Bayesian linear regression model when data augmentation in \boldsymbol{y}^* is used. The posterior conditional distributions are

$$\boldsymbol{\beta}|\boldsymbol{\Sigma},\boldsymbol{w}\sim N(\boldsymbol{\beta}_n,\boldsymbol{B}_n),$$

$$\boldsymbol{\Sigma}^{-1}|\boldsymbol{\beta},\boldsymbol{w}\sim W(\alpha_n,\boldsymbol{\Psi}_n),$$

$$\boldsymbol{y}_{il}^*|\boldsymbol{y}_{i,-l}^*,\boldsymbol{\beta},\boldsymbol{\Sigma}^{-1},\boldsymbol{y}_i\sim TN_{I_{il}}(m_{il},\tau_{ll}^2)$$
 where $\boldsymbol{B}_n=(\boldsymbol{B}_0^{-1}+\boldsymbol{X}^{*\top}\boldsymbol{X}^*)^{-1},~\boldsymbol{\beta}_n=\boldsymbol{B}_n(\boldsymbol{B}_0^{-1}\boldsymbol{\beta}_0+\boldsymbol{X}^{*\top}\boldsymbol{y}^{**}),~\boldsymbol{\Sigma}^{-1}=$
$$\boldsymbol{C}^{\top}\boldsymbol{C},~\boldsymbol{X}_i^{*\top}=~\boldsymbol{C}^{\top}\boldsymbol{X}_i,~\boldsymbol{y}_i^{**}=~\boldsymbol{C}^{\top}\boldsymbol{y}_i^*,~\boldsymbol{X}^*=\begin{bmatrix}\boldsymbol{X}_1^*\\\boldsymbol{X}_2^*\\\vdots\\\boldsymbol{X}_N^*\end{bmatrix},~\boldsymbol{\alpha}_n=\boldsymbol{\alpha}_0+\boldsymbol{N},$$

$$\boldsymbol{\Psi}_n=(\boldsymbol{\Psi}_0+\sum_{i=1}^N(\boldsymbol{y}_i^*-\boldsymbol{X}_i\boldsymbol{\beta})^{\top}(\boldsymbol{y}_i^*-\boldsymbol{X}_i\boldsymbol{\beta}))^{-1},~\boldsymbol{I}_{il}=\begin{cases}\boldsymbol{y}_{il}^*>0,~\boldsymbol{y}_{il}=1\\\boldsymbol{y}_{il}^*\leq 0,~\boldsymbol{y}_{il}=0\end{cases},$$

$$\boldsymbol{m}_{il}=\boldsymbol{x}_{il}^{\top}\boldsymbol{\beta}+\boldsymbol{f}^{\top}(\boldsymbol{y}_{i,-l}^*-\boldsymbol{X}_{i,-l}\boldsymbol{\beta}),~\boldsymbol{y}_{i,-l}^*\text{ is an }L-1\text{ dimensional vector of all components of }\boldsymbol{y}_i^*\text{ excluding }\boldsymbol{y}_{il}^*,~\boldsymbol{x}_{il}\text{ is the }l\text{-th row of }\boldsymbol{X}_i,~\boldsymbol{X}_{i,-l}\text{ is }\boldsymbol{X}_i\text{ after deleting the }l\text{-th row},~\boldsymbol{\tau}_{ll}^2=1/\boldsymbol{\sigma}^{ll},~\text{and }\boldsymbol{f}=-\boldsymbol{\sigma}^{ll}\boldsymbol{\omega}_{l,-l},~\boldsymbol{\sigma}^{jl}\text{ is the }jl\text{-th element of }\boldsymbol{\Sigma}^{-1},~\boldsymbol{\Sigma}^{-1}=\begin{bmatrix}\boldsymbol{\omega}_1^{\top}~\boldsymbol{\omega}_2^{\top};\boldsymbol{\omega}_L^{\top},&\boldsymbol{\omega}_{l,-l}^{\top}\text{ is the }l\text{-th row of }\boldsymbol{\Sigma}^{-1}\text{ extracting the }l\text{-th element.}$$

Example: Self selection in hospitalization due to a subsidized health care program in Medellín

We use the dataset 7HealthMed.csv where the dependent variable is equal to y = [Hosp SHI]' where Hosp is equal to 1 if an individual was hospitalized in 2007, 0 otherwise, and SHI is equal to 1 if the individual had subsidized health insurance that year, and 0 otherwise. Recall that our application in binary response models was to uncover the determinants of hospitalization in Medellín (Colombia), where one of the regressors was a binary indicator of being in a subsidized health care program. We can use a bivariate probit

⁸In a Bayesian setting, we can have a non identified model; however, the posterior of the model parameters exists given a proper prior distribution [29].

model if we suspect there is a dependence regarding the decisions involving these two variables. We would expect a priori that being in a subsidized health care program would imply a higher probability of being hospitalized *ceteris paribus*. However, if an individual expects to be hospitalized in the future, and the factors that drive this decision are unobserved to the econometrician, we would have a feedback effect from being hospitalized on being in a subsidized health care program.

We took into account 9 regressors: a constant, female, age, squared age, self perception of health status taking as reference bad (four categories), and the proportion of the individual's age spent living in her/his neighborhood (PTL). The last variable tries to take into account the social capital that can affect being in the subsidized health insurance program, as the target population is identified by the local government [85]. We have 12975 individuals chosen two options (subsidized regime and hospitalization).

The Algorithm A14 shows how to run a multivariate probit model in our GUI.

Algorithm A14 Instrumental variable 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 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

We set 20,000 MCMC iterations plus 1,000 iterations as burn-in, and a thinning parameter equal to 5. This implies an effective length of the posterior chains equal to 4,000 draws. We also used default values for the hyperparameters of the prior distributions. In general, the convergence diagnostics seem

good, except that there is a high level of autocorrelation for the posterior chain of the correlation between the two equations, as indicated by the dependence factors, and the trace and autocorrelation plots. Observe that the tests of [41] and [48] have NaN values for the elements (1, 1) and (2, 2) of the covariance matrix, as these parameters were set equal to 1 due to identification restrictions. This also means just two values for the dependence factors, which are actually the same due to symmetry.

R code. Self selection in hospitalization

```
1 rm(list = ls()); set.seed(010101)
2 Data <- read.csv("DataApplications/7HealthMed.csv", sep = ",</pre>
       ", header = TRUE, fileEncoding = "latin1")
3 attach(Data); str(Data)
4 p <- 2; nd <- 9; N <- length(y)/p; y <- y
_5 Xd <- as.matrix(Data[seq(1, p*N, 2),3:11])
6 XcreateMP <-function(p,nxs,nind,Data){</pre>
    pandterm = function(message) {
      stop(message, call. = FALSE)
    if (missing(nxs))
    pandterm("requires number of regressors: include intercept
        if required")
    if (missing(nind))
12
    pandterm("requires number of units (individuals)")
13
    if (missing(Data))
    pandterm("requires dataset")
     if (nrow(Data)!=nind*2)
    pandterm("check dataset! number of units times number
       alternatives should be equal to dataset rows")
    XXDat <-array(0,c(p,1+nxs,nind))</pre>
    XX<-array(0,c(p,nxs*p,nind))</pre>
19
    YY <- array (0, c(p, 1, nind))
    is<- seq(p,nind*p,p)</pre>
21
    cis<- seq(nxs,nxs*p+1,nxs)</pre>
    for(i in is){
23
24
       j <- which (i == is)
       XXDat[,,j] <-as.matrix(Data[c((i-(p-1)):i),-1])</pre>
25
      YY[,,j]<-XXDat[,1,j]
26
       for(1 in 1:p){
         XX[1,((cis[1]-(nxs-1)):cis[1]),j]<-XXDat[1,-1,j]</pre>
28
29
30
31
    return(list(y=YY,X=XX))
32 }
33 Dat <- XcreateMP(p = p, nxs = nd, nind = N, Data = Data)
34 y <- NULL; X <- NULL
35 for(i in 1:dim(Dat$y)[3]){
36 y <-c (y, Dat $ y [,,i])
    X<-rbind(X,Dat$X[,,i])</pre>
38 }
39 DataMP = list(p=p, y=y, X=X)
_{40} # Hyperparameters
41 k \leftarrow dim(X)[2]; b0 \leftarrow rep(0, k); c0 \leftarrow 1000; B0 \leftarrow c0*diag(k)
42 B0i <- solve(B0); a0 <- p - 1 + 3; Psi0 <- a0*diag(p)
43 Prior <- list(betabar = b0, A = B0i, nu = a0, V = Psi0)
_{44} # MCMC parameters
45 mcmc <- 100000; thin <- 5
46 Mcmc <- list(R = mcmc, keep = thin)
```

R code. Self selection in hospitalization, results

```
Results <- bayesm::rmvpGibbs(Data = DataMP, Mcmc = Mcmc,
      Prior = Prior)
  betatilde <- Results$betadraw / sqrt(Results$sigmadraw[,1])
  attributes(betatilde)$class <- "bayesm.mat"
  summary(coda::mcmc(betatilde))
  Quantiles for each variable:
  2.5%
              25%
                          50%
                                     75%
                                              97.5%
         -5.5137018
                    -1.946e+00
                               -1.298e+00 -6.494e-01
         0.0327632
                    9.315e-02
                                1.252e-01
                                           1.572e-01
                                                      0.2178014
  var2
         -0.0078208 -3.120e-03
                               -6.648e-04
  var3
                                            1.798e-03 0.0065977
                    1.462e-05
                                4.411e-05
  var4
         -0.0000417
                                            7.356e-05 0.0001297
         -3.8462383
                    -2.747e-01
                                3.677e-01
                                            1.015e+00
                                                      4.5937307
  var5
  var6
         -4.3469133 -7.901e-01 -1.459e-01
                                            4.952e-01
                                                      4.0485322
         -5.0522003 -1.496e+00 -8.543e-01
                                           -2.101e-01 3.3375816
13 var7
        -4.9152373 -1.360e+00 -7.058e-01 -7.392e-02 3.4898306
14 var8
         -0.1838265 -1.097e-01 -6.940e-02 -2.936e-02 0.0472931
  var9
  var10
         -2.8060844
                     3.396e-01
                                1.169e+00
                                            2.732e+00
  var11
         0.2280206
                    5.567e-01
                                9.842e-01
                                            2.293e+00 3.3507157
         -0.0574587 -2.311e-02 -9.777e-03
                                           -3.976e-03 0.0033351
18 var12
         0.0001256
                    3.125e-04
                                5.964e-04
                                           1.323e-03 0.0020878
         -3.0896704
                    2.245e-01
                                9.712e-01
                                            2.419e+00 6.1875147
20 var14
         -3.8584533
                    -1.963e-01
                                3.769e-01
                                            1.285e+00 4.8754045
        -4.5826124 -8.322e-01 -1.249e-01
                                            4.637e-01 3.9548010
22 var16
         -4.5922641 -9.026e-01 -1.634e-01
                                            4.225e-01 3.9128193
         0.1511020 3.760e-01 7.023e-01
                                           1.561e+00 2.4146495
25 sigmadraw <- Results$sigmadraw / Results$sigmadraw[,1]
26 attributes(sigmadraw)$class = "bayesm.var'
27 summary(coda::mcmc(sigmadraw))
28 Quantiles for each variable:
29 2.5%
           25%
                      50%
                               75%
                                       97.5%
        1.0000
                1.0000 1.000000
                                   1.00000
                                              1.0000
30 var1
        -0.4574 -0.0828 -0.007985
                                              0.3883
                                   0.05506
32 var3 -0.4574 -0.0828 -0.007985
                                   0.05506
                                              0.3883
        0.5711
                3.3214 9.789460 56.30992 110.2846
```

The results suggest that only female is relevant to explain hospitalization. The 95% credible interval is (3.13e-02, 0.22). Observe that only 3.11% of the sample has been hospitalized. Probit models are not well designed for this kind of dataset, but our main purpose is to illustrate the use of our GUI. On the other hand, the results suggest that age, squared age, and the proportion of age spent living in the neighborhood are statistically relevant to explain enrollment in the subsidized program. Their 95% credible intervals are (2.63e-01, 3.27e-01), (-8.05e-03, -2.05e-03) and (0.15, 0.27), respectively. The latter result seems to support the social capital hypothesis. Lastly, the 95% credible interval for the correlation between the two binary equations is (-0.07,

0.06), suggesting that there is no self selection regarding these two decisions (hospitalization and subsidized insurance). So, it seems that it is better to estimate univariate binary models for each of these dependent variables, for the sake of parsimony.

8.5 Summary

8.6 Exercises

- 1. Show that $\mathbb{E}[u_1 \text{PAER}] = \frac{\alpha_1}{1-\beta_1 \alpha_1} \sigma_1^2$ assuming that $\mathbb{E}[u_1 u_2] = 0$ where $Var(u_1) = \sigma_1^2$ in the effect of institutions on per capita GDP.
- 2. Show that $\beta_1 = \pi_1/\gamma_1$ in the effect of institutions on per capita GDP.
- 3. The effect of institutions on per capita gross domestic product 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

Given the structural demand-supply model:

$$q_i^d = \beta_1 + \beta_2 p_i + \beta_3 y_i + \beta_4 p c_i + \beta_5 p s_i + u_{i1}$$

$$q_i^s = \alpha_1 + \alpha_2 p_i + \alpha_3 e r_i + u_{i2},$$

where q^d is demand, q^s is supply, p, y, pc, ps and er are price, income, complementary price, substitute price, and exchange rate. Complementary and substitute prices are prices of a complementary and substitute goods of q. Assume that $\boldsymbol{\beta} = \begin{bmatrix} 5 & -0.5 & 0.8 & -0.4 & 0.7 \end{bmatrix}^{\top}$, $\boldsymbol{\alpha} = \begin{bmatrix} -2 & 0.5 & -0.4 \end{bmatrix}^{\top}$, $u_1 \sim N(0, 0.5^2)$ and $u_2 \sim N(0, 0.5^2)$. In addition, assume that $y \sim N(10, 1)$, $pc \sim N(5, 1)$, $ps \sim N(5, 1)$ and $tc \sim N(15, 1)$.

- •Find the reduce-form model using that in equilibrium demand and supply are equal, that is, $q^d = q^s$. This condition defines the observable quantity (q).
- •Simulate p and q from the reduce-form equations.
- Preform inference of the *reduce-form* model using the *rmultireg* command from the *bayesm* package.

Exercises 175

•Use the posterior draws of the *reduce-form* parameters to perform inference of the *structural* parameters. Any issue? Hint: Are all *structural* parameters exactly identified?

5. Utility demand continues

- •Run the **Utility demand** application using our GUI and the information in the dataset *Utilities.csv*. Hint: This file should be modified to agree the structure that requires our GUI (see the dataset *5Institutions.csv* in the folder *DataApp* of our GitHub repository -https://github.com/besmarter/BSTApp- for a template).
- Program from scratch the Gibbs sampler algorithm in this application.

6. Simulation exercise of instrumental variables continues I

- (a) Use the setting of the simulation exercise of instrumental variables to analyze what happens when the instrument is weak, for instance, setting $\gamma_2 = 0.2$, and compare the performance of the posterior means of the ordinary and instrumental models.
- (b) Perform a simulation that helps to analyze how the degree of exogeneity of the instrument affects the performance of the posterior mean of the instrumental variable model.

7. Simulation exercise of instrumental variables continues II

Program from scratch the Gibbs sampling algorithm of the instrumental model for the simulation exercise of the instrumental variables.

8. The effect of institutions on per capita gross domestic product continues II

Estimate the structural Equation 8.1 using the instrumental variable model where the instrument of PAER is $\log(Mort)$. Compare the effect of property rights on per capita GDP of this model with the effect estimated in the example of the effect of institutions on per capita gross domestic product. Use the file 6Institutions.csv to do this exercise in our GUI, and set $B_0 = 100I_5$, $\beta_0 = 0_5$, $\gamma_0 = 0_2$, $G_0 = 100I_2$ $\alpha_0 = 3$ and $\Psi_0 = 3I_2$. The MCMC iterations, burn-in and thinning parameters are 50000, 1000 and 5, respectively.

Time series models

Panel data models

Bayesian model average

- 11.1 Calculating the marginal likelihood
- 11.1.1 Savage-Dickey density ratio
- 11.1.2 Gelfand-Dey method
- 11.1.3 Chib's methods

Part III

Advanced methods: Theory, applications and programming

Hierarchical models

12.1 Finite mixtures

12.2 Direchlet processes

Causal inference

Machine learning

14.1 Cross validation and Bayes facto	14.1	\mathbf{Cross}	validation	and	Bayes	factor
---------------------------------------	------	------------------	------------	-----	-------	--------

14.2 Regularization

14.3 Bayesian additive regression trees

14.4 Gaussian processes

Further topics

15.1	Approximate	Bayesian	computation

- 15.2 Variational Bayes
- 15.3 Integrated nested Laplace approximations
- 15.4 Bayesian exponential tilted empirical likelihood

- [1] 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.
- [2] 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.
- [3] R. Baath. Package bayesboot, 2018.
- [4] M. Bayarri and J. Berger. P-values for composite null models. *Journal of American Statistical Association*, 95:1127–1142, 2000.
- [5] M. J. Bayarri and J. Berger. The interplay of bayesian and frequentist analysis. *Statistical Science*, 19(1):58–80, 2004.
- [6] 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.
- [7] 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.
- [8] 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.
- [9] J. Berger. Statistical Decision Theory and Bayesian Analysis. Springer, third edition edition, 1993.
- [10] J. Berger. The case for objective bayesian analysis. Bayesian Analysis, 1(3):385-402, 2006.
- [11] James O Berger. Statistical decision theory and Bayesian analysis. Springer Science & Business Media, 2013.
- [12] J. Bernardo and A. Smith. Bayesian Theory. Wiley, Chichester, 1994.

[13] 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.

- [14] G. E. P. Box. Science and statistics. *Journal of the American Statistical Association*, 71:791–799, 1976.
- [15] George EP Box. Robustness in the strategy of scientific model building. In *Robustness in statistics*, pages 201–236. Elsevier, 1979.
- [16] Colin Cameron and Pravin Trivedi. *Microeconometrics: Methods and Applications*. Cambridge, 2005.
- [17] George Casella and Roger Berger. Statistical inference. CRC Press, 2024.
- [18] W. Chang. Web Application Framework for R: Package shiny. R Studio, 2018.
- [19] V. Chernozhukov and H. Hong. An MCMC approach to classical estimation. *Journal of Econometrics*, 115:293–346, 2003.
- [20] S. Chib. Bayes inference in the Tobit censored regression model. *Journal of Econometrics*, 51:79–99, 1992.
- [21] Siddhartha Chib. Marginal likelihood from the gibbs output. *Journal* of the american statistical association, 90(432):1313–1321, 1995.
- [22] Siddhartha Chib and Ivan Jeliazkov. Marginal likelihood from the metropolis-hastings output. *Journal of the American Statistical Association*, 96(453):270–281, 2001.
- [23] M. Clyde and E. George. Model uncertatinty. Statistical Science, 19(1):81–94, 2004.
- [24] T. Conley, C. Hansen, and P. Rossi. Plausibly exogenous. *The Review of Economics and Statistics*, 94(1):260–272, 2012.
- [25] A. P. Dawid, M. Musio, and S. E. Fienberg. From statistical evidence to evidence of causality. *Bayesian Analysis*, 11(3):725–752, 2016.
- [26] 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.
- [27] M. H. DeGroot. *Probability and statistics*. Addison-Wesley Publishing Co., London, 1975.
- [28] Persi Diaconis and Donald Ylvisaker. Conjugate priors for exponential families. *The Annals of statistics*, 7(2):269–281, 1979.

[29] Y. D. Edwards and G. M. Allenby. Multivariate analysis of multiple response data. *Journal of Marketing Research*, 40:321–334, 2003.

- [30] B. Efron. Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, 7:1–26, 1979.
- [31] Bradley Efron and Trevor Hastie. Computer age statistical inference, volume 5. Cambridge University Press, 2016.
- [32] R. Fisher. Statistical Methods for Research Workers. Hafner, New York, 13th edition, 1958.
- [33] A. E. Gelfand and A. F. M. Smith. Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Associ*ation, 85:398–409, 1990.
- [34] 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.
- [35] 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.
- [36] A. Gelman, X. Meng, and H. Stern. Posterior predictive assessment of model fitness via realized discrepancies. *Statistica sinica*, pages 733–760, 1996.
- [37] Andrew Gelman, John B Carlin, Hal S Stern, David Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. Chapman and Hall/CRC, 2021.
- [38] 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.
- [39] Andrew Gelman and Guido Imbens. Why ask why? forward causal inference and reverse causal questions. Technical report, National Bureau of Economic Research, 2013.
- [40] 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.
- [41] J. Geweke. Bayesian Statistics, chapter Evaluating the accuracy of sampling-based approaches to calculating posterior moments. Clarendon Press, Oxford, UK., 1992.
- [42] John Geweke. Contemporary Bayesian econometrics and statistics, volume 537. John Wiley & Sons, 2005.

[43] I. J. Good. The bayes/non bayes compromise: A brief review. *Journal of the American Statistical Association*, 87(419):597–606, September 1992.

- [44] S. N. Goodman. Toward evidence-based medical statistics. 1: The P value fallacy. *Annals of internal medicine*, 130(12):995–1004, 1999.
- [45] Edward Greenberg. *Introduction to Bayesian econometrics*. Cambridge University Press, 2012.
- [46] Damodar N. Gujarati and Dawn C. Porter. *Basic Econometrics*. McGraw-Hill Education, New York, NY, 5th edition, 2009.
- [47] W. Hastings. Monte Carlo sampling methods using Markov chains and their application. *Biometrika*, 57:97–109, 1970.
- [48] P. Heidelberger and P. D. Welch. Simulation run length control in the presence of an initial transient. *Operations Research*, 31(6):1109–1144, 1983.
- [49] H. Jeffreys. Some test of significance, treated by the theory of probability. *Proceedings of the Cambridge philosophy society*, 31:203–222, 1935.
- [50] H. Jeffreys. *Theory of Probability*. Oxford University Press, London, 1961.
- [51] M. I. Jordan, Z. Ghahramani, T. Jaakkola, and L. Saul. Introduction to variational methods for graphical models. *Machine Learning*, 37:183–233, 1999.
- [52] Daniel Kahneman. Thinking, fast and slow. Macmillan, 2011.
- [53] 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.
- [54] R. Kass. Statistical inference: the big picture. Statistical science, 26(1):1–9, 2011.
- [55] Robert E. Kass and Adrian E. Raftery. Bayes factorss. Journal of American Statistical Association, 90(430):773-795, 1995.
- [56] Gary M Koop. Bayesian econometrics. John Wiley & Sons Inc., 2003.
- [57] Hideo Kozumi and Genya Kobayashi. Gibbs sampling methods for Bayesian quantile regression. *Journal of Statistical Computation and Simulation*, 81(11):1565—1578, 2011.
- [58] Tony Lancaster. An introduction to modern Bayesian econometrics. Blackwell Oxford, 2004.

- [59] P. Laplace. Théorie Analytique des Probabilités. Courcier, 1812.
- [60] Pierre Simon Laplace. Mémoire sur la probabilité de causes par les évenements. Mémoire de l'académie royale des sciences, 1774.
- [61] E.L. Lehmann and George Casella. *Theory of Point Estimation*. Springer, second edition edition, 2003.
- [62] Alex Lenkoski, Anna Karl, and Andreas Neudecker. *Package ivbma*, 2013.
- [63] D. V. Lindley. The philosophy of statistics. The Statistician, 49(3):293–337, 2000.
- [64] D. V. Lindley and L. D. Phillips. Inference for a Bernoulli process (a Bayesian view). *American Statistician*, 30:112–119, 1976.
- [65] Dennis V Lindley. A statistical paradox. *Biometrika*, 44(1/2):187–192, 1957.
- [66] 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.
- [67] Andrew D. Martin, Kevin M. Quinn, and Jong Hee Park. Package MCMCpack, 2018.
- [68] R. McCulloch and P. Rossi. An exact likelihood analysis of the multinomial probit model. *Journal of Econometrics*, 64:207–240, 1994.
- [69] 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.
- [70] 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.
- [71] 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.
- [72] 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.
- [73] Agostino Nobile. Comment: Bayesian multinomial probit models with a normalization constraint. *Journal of Econometrics*, 99(2):335–345, 2000.

[74] G. Parmigiani and L. Inoue. Decision theory principles and approaches. John Wiley & Sons, 2008.

- [75] 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.
- [76] Giovanni Petris, Sonia Petrone, and Patrizia Campagnoli. Dynamic linear models. In *Dynamic Linear Models with R*, pages 31–84. Springer, 2009.
- [77] Martyn Plummer, Nicky Best, Kate Cowles, Karen Vines, Deepayan Sarkar, Douglas Bates, Russell Almond, and Arni Magnusson. Output Analysis and Diagnostics for MCMC, 2016.
- [78] R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, 2021.
- [79] R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, 2023.
- [80] A. Raftery. Bayesian model selection in social research. Sociological Methodology, 25:111–163, 1995.
- [81] Adrian Raftery, Jennifer Hoeting, Chris Volinsky, Ian Painter, and Ka Yee Yeung. *Package BMA*, 2012.
- [82] 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.
- [83] 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.
- [84] 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.
- [85] R. Ramírez-Hassan, A. Guerra-Urzola. Bayesian treatment effects due to a subsidized health program: The case of preventive health care utilization in medellín (Colombia). *Empirical Economics*, Forthcoming, 2019.
- [86] 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.

[87] A. Ramírez-Hassan and M. Graciano-Londoño. A guided tour of Bayesian regression. Technical report, Universidad EAFIT, 2020.

- [88] P. Rossi. Package bayesm, 2017.
- [89] Peter E Rossi, Greg M Allenby, and Rob McCulloch. *Bayesian statistics and marketing*. John Wiley & Sons, 2012.
- [90] Donnald B. Rubin. The Bayesian bootstrap. *The Annals of Statistics*, 9(1):130–134, 1981.
- [91] L. J. Savage. The foundations of statistics. John Wiley & Sons, Inc., New York, 1954.
- [92] Robert Schlaifer and Howard Raiffa. Applied statistical decision theory. Wiley New York, 1961.
- [93] 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.
- [94] Steve Selvin. A problem in probability (letter to the editor). The American Statistician, 11(1):67–71, 1975.
- [95] Steve Selvin. A problem in probability (letter to the editor). The American Statistician, 11(3):131–134, 1975.
- [96] 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.
- [97] A. F. M. Smith. A General Bayesian Linear Model. *Journal of the Royal Statistical Society. Series B (Methodological).*, 35(1):67–75, 1973.
- [98] Stan Development Team. shinystan: Interactive visual and numerical diagnostics and posterior analysis for Bayesian models., 2017. R package version 2.3.0.
- [99] Stephen Stigler. Richard price, the first bayesian. Statistical Science, 33(1):117–125, 2018.
- [100] 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.
- [101] Luke Tierney. Markov chains for exploring posterior distributions. *The Annals of Statistics*, pages 1701–1728, 1994.
- [102] 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.

[103] Aad W Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.

- [104] Ronald L. Wasserstein and Nicole A. Lazar. The ASA's statement on p-values: context, process and purpose. *The American Statistician*, 2016.
- [105] P. Woodward. BugsXLA: Bayes for the common man. *Journal of Statistical Software*, 14(5):1–18, 2005.
- [106] Jeffrey M Wooldridge. Econometric analysis of cross section and panel data. MIT press, 2010.
- [107] Jeffrey M. Wooldridge. Introductory Econometrics: A Modern Approach. Cengage Learning, Boston, MA, 6th edition, 2016.
- [108] A. Zellner. Introduction to Bayesian inference in econometrics. John Wiley & Sons Inc., 1996.
- [109] S. Ziliak. Guinnessometrics; the Economic Foundation of student's t. Journal of Economic Perspectives, 22(4):199–216, 2008.

TABLE 15.1 Libraries and commands in BEsmarter GUI.

Univariate models				
Model	Library	Command	Reference	
Normal	MCMCpack	MCMCregress	[67]	
Logit	MCMCpack	MCMClogit	[67]	
Probit	bayesm	rbprobitGibbs	[88]	
Multinomial(Mixed) Probit	bayesm	rmnpGibbs	[88]	
Multinomial (Mixed) Logit	bayesm	rmnlIndepMetrop	[88]	
Ordered Probit	bayesm	rordprobitGibbs	[88]	
Negative Binomial(Poisson)	bayesm	rnegbinRw	[88]	
Tobit	MCMCpack	MCMCtobit	[67]	
Quantile	MCMCpack	MCMCquantreg	[67]	
	tivariate mo			
Model	Library	Command	Reference	
Multivariate	bayesm	rmultireg	[88]	
Seemingly Unrelated Regression	bayesm	rsurGibbs	[88]	
Instrumental Variable	bayesm	rivGibbs	[88]	
Bivariate Probit	bayesm	rmvpGibbs	[88]	
Hierarchical longitudinal models				
Model	Library	Command	Reference	
Normal	MCMCpack	MCMChregress	[67]	
Logit	MCMCpack	MCMChlogit	[67]	
Poisson	MCMCpack	MCMChpoisson	[67]	
	esian Bootst			
Model	Library	Command	Reference	
Bayesian bootstrap	bayesboot	bayesboot	[3]	
	an model ave			
Model	Library	Command	Reference	
Normal (BIC)	BMA	bic.glm	[81]	
Normal (MC ³)	BMA	MC3.REG	[81]	
Normal (instrumental variables)	ivbma	ivbma	[62]	
Logit (BIC)	BMA	bic.glm	[81]	
Gamma (BIC)	BMA	bic.glm	[81]	
Poisson (BIC)	BMA	bic.glm	[81]	
	Diagnostics			
Diagnostic	Library	Command	Reference	
Trace plot	coda	traceplot	[77]	
Autocorrelation plot	coda	autocorr.plot	[77]	
Geweke test	coda	geweke.diag	[77]	
Raftery & Lewis test	coda	raftery.diag	[77]	
Heidelberger & Welch test	coda	heidel.diag	[77]	

TABLE 15.2 Datasets templates in folder *DataSim*.

Univariate models			
Model	Data set file	Data set simulation	
Normal	11SimNormalmodel.csv	11SimNormal.R	
Logit	12SimLogitmodel.csv	12SimLogit	
Probit	13SimProbitmodel.csv	13SimProbit.R	
Multinomial(Mixed) Probit	14SimMultProbmodel.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	
Multivariate models			
Model	Data set file	Data set simulation	
Multivariate	21SimMultivariate.csv	21SimMultReg.R	
Seemingly Unrelated Regression	22SimSUR.csv	22SimSUR.R	
Instrumental Variable	23SimIV.csv	23SimIV.R	
Bivariate Probit	24SimMultProbit.csv	24SimMultProbit.R	
Hierarchical longitudinal models			
Model	Data set file	Data set simulation	
Normal	31SimLogitudinalNormal.csv	31SimLogitudinalNormal.R	
Logit	32SimLogitudinalLogit.csv	32SimLogitudinalLogit.R	
Poisson	33SimLogitudinalPoisson.csv	33SimLogitudinalPoisson.R	
Bayesian Bootstrap			
Model	Data set file	Data set simulation	
Bayesian bootstrap	41SimBootstrapmodel.csv	41SimBootstrapmodel.R	
Bayesian model averaging			
Model	Data set file	Data set simulation	
Normal (BIC)	511SimNormalBMA.csv	511SimNormalBMA.R	
Normal (MC ³)	512SimNormalBMA.csv	512SimNormalBMA.R	
Normal (instrumental variables)	513SimNormalBMAivYXW.csv	513SimNormalBMAiv.R	
	513SimNormalBMAivZ.csv		
Logit (BIC)	52SimLogitBMA.csv	52SimLogitBMA.R	
Gamma (BIC)	53SimGammaBMA.csv	53SimGammaBMA.R	
Poisson (BIC)	53SimPoissonBMA.csv	53SimPoissonBMA.R	

TABLE 15.3 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	2 Health Med.csv	MedVisPrevOr	
Negative Binomial(Poisson)	2 Health Med.csv	MedVisPrev	
Tobit	1ValueFootballPlayers.csv	log(ValueCens)	
Quantile	1ValueFootballPlayers.csv	$\log(\text{Value})$	
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 = [Hosp SHI]'	
Hierarchical longitudinal models			
Model	Data set file	Dependent variable	
Normal	8PublicCap.csv	$\log(\text{gsp})$	
Logit	9VisitDoc.csv	DocVis	
Poisson	9VisitDoc.csv	DocNum	
Bayesian Bootstrap			
Model	Data set file	Dependent variable	
Bayesian bootstrap	1ValueFootballPlayers.csv	log(Value)	
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	avghhi and avglgdpcap	
	12ExportDiversificationHHIInstr.csv		
Logit (BIC)	13InternetMed.csv	internet	
Gamma (BIC)	14ValueFootballPlayers.csv	log market value	
Poisson (BIC)	15Fertile2.csv	ceb	