
Half Title

Solution Manual
Introduction to Bayesian Inference:
A GUIDed tour using R



Title Page

Solution Manual
Introduction to Bayesian Inference:
A GUIDed tour using R

by Andrés Ramírez-Hassan, PhD. Statistical Science.

LOC Page

To my parents, Nancy and Orlando.



Contents

Foreword	ix
Preface	xi
Symbols	xiii
I Foundations: Theory, simulation methods and programming	1
1 Solutions of chapter 1 Basic formal concepts	3
1.1 Solutions of Exercises	3
2 Solutions of chapter 2 Conceptual differences: Bayesian and Frequentist approaches	23
2.1 Solutions of Exercises	23
3 Solutions of chapter 4 Cornerstone models: Conjugate families	29
3.1 Solutions of Exercises	29
Bibliography	45



Foreword





Preface



Symbols

Symbol Description

\neg	Negation symbol.	\mathcal{R}	The Real set.
\propto	Proportional symbol.	\emptyset	Empty set.
\perp	Independence symbol.	$\mathbb{1}$	Indicator function.



Part I

Foundations: Theory, simulation methods and programming



1

Solutions of chapter 1

Basic formal concepts

1.1 Solutions of 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?*

Answer

Set WB and WG equal to the events that the witness said the cab was blue and green, respectively. Set B and G equal to the events that the cabs are blue and green, respectively. We need to calculate $P(B|WB)$, then:

$$\begin{aligned} P(B|WB) &= \frac{P(B, WB)}{P(WB)} \\ &= \frac{P(WB|B) \times P(B)}{P(WB|B) \times P(B) + (1 - P(WB|B)) \times (1 - P(B))} \\ &= \frac{0.8 \times 0.15}{0.8 \times 0.15 + 0.2 \times 0.85} \\ &= 0.41 \end{aligned} \tag{1.1}$$

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?

Answer

Let's name P_i the event *contestant picks door No. i* , H_i the event *host picks*

door No. i , and C_i the event *car is behind door No. i* . Let's assume that the contestant picked door number 1, and the host picked door number 3, then the contestant is interested in the probability of the event $P(C_i|H_3, P_1)$, $i = 2$ or 4 . Then, $P(H_3|C_3, P_1) = 0$, $P(H_3|C_2, P_1) = P(H_3|C_4, P_1) = 1/2$ and $P(H_3|C_1, P_1) = 1/3$. Then,

$$\begin{aligned}
 P(C_i|H_3, P_1) &= \frac{P(C_i, H_3, P_1)}{P(H_3, P_1)} \\
 &= \frac{P(H_3|C_i, P_1)P(C_i|P_1)P(P_1)}{P(H_3|P_1) \times P(P_1)} \\
 &= \frac{P(H_3|C_i, P_1)P(C_i)}{P(H_3|P_1)} \\
 &= \frac{1/2 \times 1/4}{1/3} \\
 &= \frac{3}{8},
 \end{aligned} \tag{1.2}$$

where the third equation uses the fact that C_i and P_i are independent events, and $P(H_3|P_1) = 1/3$ due to this depending just on P_1 (not on C_i).

Therefore, changing the initial decision increases the probability of getting the car from $1/4$ to $3/8$!

Let's check the case with n doors, and assume that the contestant picks the door No. 1, the car is behind the door No. n , and the host, who knows what is behind each door, opens any of the remaining $n - 2$ doors, where there is a goat. The contestant is interested in the probability of the event:

$$\begin{aligned}
 P(C_n|(H_2 \cup \dots \cup H_{n-1}) \cap P_1) &= \frac{P((H_2 \cup H_3 \cup \dots \cup H_{n-1})|C_n \cap P_1)P(C_n|P_1)P(P_1)}{P((H_2 \cup H_3 \cup \dots \cup H_{n-1})|P_1)P(P_1)} \\
 &= \frac{[P(H_2|C_n \cap P_1) + \dots + P(H_{n-1}|C_n \cap P_1)]P(C_n)}{P(H_2|P_1) + P(H_3|P_1) + \dots + P(H_{n-1}|P_1)} \\
 &= \frac{1 \times (\frac{1}{n})}{\frac{1}{n-1} + \frac{1}{n-1} + \dots + \frac{1}{n-1}} \\
 &= \left(\frac{1}{n}\right) \left(\frac{n-1}{n-2}\right).
 \end{aligned} \tag{1.3}$$

In general, the probability of winning the car changing the pick is $\frac{1}{n} \frac{n-1}{n-2}$, while the probability of winning given no change is $\frac{1}{n}$. Given that $\frac{1}{n} \frac{n-1}{n-2} > \frac{1}{n}$ for all $n \geq 3$, where the difference between both probabilities is $\frac{1}{n(n-2)}$. We observe that as the number of doors increases, the difference between the two probabilities becomes zero.

Let's see a code for the general setting,

R code. The Monty Hall Problem

```

set.seed(0101) # Set simulation seed
S <- 100000 # Simulations
Game <- function(opt = 3){
  # opt: number of options. opt > 2
  opts <- 1:opt
  car <- sample(opts, 1) # car location
  guess1 <- sample(opts, 1) # Initial guess

  if(opt == 3 && car != guess1) {
    host <- opts[-c(car, guess1)]
  } else {
    host <- sample(opts[-c(car, guess1)], 1)
  }

  win1 <- guess1 == car # Win given no change

  if(opt == 3) {
    guess2 <- opts[-c(host, guess1)]
  } else {
    guess2 <- sample(opts[-c(host, guess1)], 1)
  }

  win2 <- guess2 == car # Win given change

  return(c(win1, win2))
}
#Win probabilities
Prob <- rowMeans(replicate(S, Game(opt = 4)))
#Winning probabilities no changing door
Prob[1]
0.25151
#Winning probabilities changing door
Prob[2]
0.37267

```

3. Solve the health insurance example using a Gamma prior in the rate parametrization, that is, $\pi(\lambda) = \frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} \lambda^{\alpha_0-1} \exp\{-\lambda\beta_0\}$.

Answer

First, we get the posterior distribution,

$$\pi(\lambda|\mathbf{y}) = \left(\frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} \lambda^{\alpha_0-1} e^{-\lambda\beta_0} \right) \left(\prod_{i=1}^N \frac{\lambda^{y_i} e^{-\lambda}}{y_i!} \right) \quad (1.4)$$

$$\begin{aligned}
&= \left(\frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} \lambda^{\alpha_0-1} e^{-\lambda\beta_0} \right) \left(\frac{\lambda^{\sum_{i=1}^N y_i} e^{-N\lambda}}{\prod_{i=1}^N y_i!} \right) \\
&= \left(\frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} \frac{1}{\prod_{i=1}^N y_i!} \right) \lambda^{\sum_{i=1}^N y_i + \alpha_0 - 1} e^{-\lambda(\beta_0 + N)} \\
&\propto \lambda^{\sum_{i=1}^N y_i + \alpha_0 - 1} e^{-\lambda(\beta_0 + N)}.
\end{aligned} \tag{1.5}$$

The last expression is the kernel of a Gamma distribution with parameters $\alpha_n = \sum_{i=1}^N y_i + \alpha_0$ and $\beta_n = \beta_0 + N$.

Given that $\int_0^\infty \pi(\lambda|\mathbf{y}) d\lambda = 1$, then the constant of proportionality in the last expression is $\Gamma(\alpha_n) / \beta_n^{\alpha_n}$. Therefore the posterior density function $\pi(\lambda|\mathbf{y})$ is $G(\alpha_n, \beta_n)$.

The posterior mean is

$$\begin{aligned}
E[\lambda|\mathbf{y}] &= \frac{\alpha_n}{\beta_n} \\
&= \frac{\sum_{i=1}^N y_i + \alpha_0}{\beta_0 + N} \\
&= \left(\frac{N}{\beta_0 + N} \right) \bar{y} + \left(\frac{\beta_0}{\beta_0 + N} \right) \frac{\alpha_0}{\beta_0} \\
&= w\bar{y} + (1-w)E[\lambda],
\end{aligned} \tag{1.6}$$

where $w = \frac{N}{\beta_0 + N}$, \bar{y} is the sample mean, and $E[\lambda] = \frac{\alpha_0}{\beta_0}$.

The posterior predictive distribution is given by

$$\begin{aligned}
\pi(Y_0|\mathbf{y}) &= \int_0^\infty \frac{\lambda^{y_0} e^{-\lambda}}{y_0!} \pi(\lambda|\mathbf{y}) d\lambda \\
&= \int_0^\infty \left(\frac{\lambda^{y_0} e^{-\lambda}}{y_0!} \right) \left(\frac{\beta_n^{\alpha_n}}{\Gamma(\alpha_n)} \lambda^{\alpha_n-1} e^{-\lambda\beta_n} \right) d\lambda \\
&= \frac{\beta_n^{\alpha_n}}{\Gamma(\alpha_n) y_0!} \int_0^\infty \lambda^{y_0 + \alpha_n - 1} e^{-\lambda(1+\beta_n)} d\lambda \\
&= \frac{\beta_n^{\alpha_n}}{\Gamma(\alpha_n) y_0!} \frac{\Gamma(y_0 + \alpha_n)}{(1 + \beta_n)^{y_0 + \alpha_n}} \\
&= \frac{\Gamma(y_0 + \alpha_n)}{\Gamma(\alpha_n) y_0!} \left(\frac{1}{1 + \beta_n} \right)^{y_0} \left(\frac{\beta_n}{1 + \beta_n} \right)^{\alpha_n}
\end{aligned} \tag{1.7}$$

$$\begin{aligned}
&= \frac{(y_0 + \alpha_n - 1)!}{(\alpha_n - 1)! y_0!} \left(\frac{1}{1 + \beta_n} \right)^{y_0} \left(\frac{\beta_n}{1 + \beta_n} \right)^{\alpha_n} \\
&= \binom{y_0 + \alpha_n - 1}{y_0} \left(\frac{1}{1 + \beta_n} \right)^{y_0} \left(\frac{\beta_n}{1 + \beta_n} \right)^{\alpha_n}.
\end{aligned}$$

Therefore $Y_0|y \sim NB(\alpha_n, p_n)$ where $p_n = \frac{1}{1+\beta_n}$.

To use empirical Bayes, we have the following setting

$$[\hat{\alpha}_0, \hat{\beta}_0] = \arg \max_{\alpha_0, \beta_0} \ln(p(\mathbf{y})),$$

where

$$\begin{aligned}
p(y) &= \int_0^\infty \left(\frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} \lambda^{\alpha_0-1} e^{-\lambda\beta_0} \right) \left(\prod_{i=1}^N \frac{\lambda^{y_i} e^{-\lambda}}{y_i!} \right) d\lambda \quad (1.8) \\
&= \left(\frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0) \prod_{i=1}^N y_i!} \right) \int_0^\infty \lambda^{\sum_{i=1}^N y_i + \alpha_0 - 1} e^{-\lambda(\beta_0 + N)} d\lambda \\
&= \left(\frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0) \prod_{i=1}^N y_i!} \right) \left(\frac{\Gamma(\sum_{i=1}^N y_i + \alpha_0)}{(\beta_0 + N)^{\sum_{i=1}^N y_i + \alpha_0}} \right) \\
&= \frac{\Gamma(\sum_{i=1}^N y_i + \alpha_0)}{\Gamma(\alpha_0) \prod_{i=1}^N y_i!} \left(\frac{1}{\beta_0 + N} \right)^{\sum_{i=1}^N y_i} \left(\frac{\beta_0}{\beta_0 + N} \right)^{\alpha_0}.
\end{aligned}$$

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

```

set.seed(010101)
y <- c(0, 3, 2, 1, 0) # Data
N <- length(y)

# Predictive distribution
ProbBo <- function(y, a0, b0){
  N <- length(y)
  #sample size
  aN <- a0 + sum(y)
  # Posterior shape parameter
  bN <- b0 + N
  # Posterior scale parameter
  p <- 1 / (bN + 1)
  # Probability negative binomial density
  Pr <- 1 - pnbinom(0, size = aN, prob = (1 - p))
  # Probability of visiting the Doctor
  # Observe that in R there is a slightly
  # different parametrization.
  return(Pr)
}

# Using a vague prior:
a0 <- 0.001 # Prior shape parameter
b0 <- 0.001 # Prior scale parameter
PriMeanV <- a0 / b0 # Prior mean
PriVarV <- a0 / b0^2 # Prior variance
Pp <- ProbBo(y, a0 = 0.001, b0 = 0.001)
# This setting is vague prior information.
Pp
0.67

```

*R code. Health insurance, predictive distribution
using empirical Bayes*

```
# Using Empirical Bayes
LogMgLik <- function(theta, y){
  N <- length(y)
  #sample size
  a0 <- theta[1]
  # prior shape hyperparameter
  b0 <- theta[2]
  # prior scale hyperparameter
  aN <- sum(y) + a0
  # posterior shape parameter
  if(a0 <= 0 || b0 <= 0){
    #Avoiding negative values
    lnp <- -Inf
  } else { lnp <- lgamma(aN) - sum(y)*log(b0+N) +
    a0*log(b0/(b0+N)) - lgamma(a0) }
  # log marginal likelihood
  return(-lnp)
}

theta0 <- c(0.01, 0.01)
# Initial values
control <- list(maxit = 1000)
# Number of iterations in optimization
EmpBay <- optim(theta0, LogMgLik, method = "BFGS",
  control = control, hessian = TRUE, y = y)
# Optimization
EmpBay$convergence
# Checking convergence
EmpBay$value # Maximum
a0EB <- EmpBay$par[1]
# Prior shape using empirical Bayes
a0EB
128.383
b0EB <- EmpBay$par[2]
# Prior scale using empirical Bayes
b0EB
106.801

PriMeanEB <- a0EB / b0EB
# Prior mean
PriVarEB <- a0EB / b0EB^2
# Prior variance
PpEB <- ProbBo(y, a0 = a0EB, b0 = b0EB)
# This setting is using empirical Bayes.
PpEB
0.69
```

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, $Y_i \sim \text{Ber}(p)$, 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 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).

Answer

The posterior distribution is given by

$$\begin{aligned} \pi(p|\mathbf{y}) &= \left[\frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0)\Gamma(\beta_0)} p^{\alpha_0-1} (1-p)^{\beta_0-1} \right] \left[\prod_{i=1}^N p^{y_i} (1-p)^{1-y_i} \right] \quad (1.9) \\ &= \frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0)\Gamma(\beta_0)} p^{\sum_{i=1}^N y_i + \alpha_0 - 1} (1-p)^{\beta_0 + N - \sum_{i=1}^N y_i - 1} \\ &\propto p^{\sum_{i=1}^N y_i + \alpha_0 - 1} (1-p)^{\beta_0 + N - \sum_{i=1}^N y_i - 1}. \end{aligned}$$

The last expression is the kernel of a Beta distribution with parameters $\alpha_n = \sum_{i=1}^N y_i + \alpha_0$ and $\beta_n = \beta_0 + N - \sum_{i=1}^N y_i$. Thus, the posterior mean is

$$\begin{aligned} E[p|\mathbf{y}] &= \frac{\alpha_n}{\alpha_n + \beta_n} \\ &= \frac{\sum_{i=1}^N y_i + \alpha_0}{\alpha_0 + \beta_0 + N} \quad (1.10) \\ &= \frac{N\bar{y}}{\alpha_0 + \beta_0 + N} + \frac{\alpha_0}{\alpha_0 + \beta_0 + N} \\ &= \frac{N}{\alpha_0 + \beta_0 + N} (\bar{y}) + \frac{\alpha_0 + \beta_0}{\alpha_0 + \beta_0 + N} \left(\frac{\alpha_0}{\alpha_0 + \beta_0} \right) \\ &= w(\bar{y}) + (1-w)E[p], \end{aligned}$$

where $w = \frac{N}{\alpha_0 + \beta_0 + N}$, \bar{y} is the sample mean, and $E[p] = \frac{\alpha_0}{\alpha_0 + \beta_0}$ is the prior mean.

The posterior predictive distribution of claim the next year is given by

$$\begin{aligned}
\pi(Y_0 = 1|\mathbf{y}) &= \int_0^1 P(Y_0 = 1|\mathbf{y}, p) \pi(p|\mathbf{y}) dp \\
&= \int_0^1 p \times \pi(p|\mathbf{y}) dp \\
&= \mathbb{E}[p|\mathbf{y}] \\
&= \frac{\alpha_n}{\alpha_n + \beta_n}.
\end{aligned} \tag{1.11}$$

To use empirical Bayes, we have the following setting

$$[\hat{\alpha}_0, \hat{\beta}_0] = \arg \max_{\alpha_0, \beta_0} \ln(p(\mathbf{y})),$$

where

$$\begin{aligned}
p(\mathbf{y}) &= \int_0^1 \left[\frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0) \Gamma(\beta_0)} p^{\alpha_0-1} (1-p)^{\beta_0-1} \right] \left[\prod_{i=1}^N (1-p)^{1-y_i} \right] dp \tag{1.12} \\
&= \frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0) \Gamma(\beta_0)} \int_0^1 p^{\sum_{i=1}^N y_i + \alpha_0 - 1} (1-p)^{\beta_0 + N - \sum_{i=1}^N y_i - 1} dp \\
&= \frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0) \Gamma(\beta_0)} \frac{\Gamma\left(\sum_{i=1}^N y_i + \alpha_0\right) \Gamma\left(\beta_0 + N - \sum_{i=1}^N y_i\right)}{\Gamma(\alpha_0 + \beta_0 + N)}.
\end{aligned}$$

R code. Car claim, predictive distribution using vague hyperparameters

```

set.seed(010101)
y <- c(0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0)
# Data
N <- length(y)

#require(TailRank)

# Predictive distribution
ProbBo <- function(y, a0, b0){
  N <- length(y)
  #sample size
  aN <- a0 + sum(y)
  # Posterior shape parameter
  bN <- b0 + N - sum(y)
  # Posterior scale parameter
  pr <- aN / (aN + bN)
  # Probability of a claim the next year
  return(pr)
}

# Using a vague prior:
a0 <- 1 # Prior shape parameter
b0 <- 1 # Prior scale parameter
PriMeanV <- a0 / (a0 + b0)
# Prior mean
PriVarV <- (a0*b0) / (((a0+b0)^2)*(a0+b0+1))
# Prior variance

Pp <- ProbBo(y, a0 = 1, b0 = 1)
# This setting is defining vague prior information.
# The probability of a claim
Pp
0.47

```

R code. Car claim, predictive distribution using empirical Bayes

```

# Using Empirical Bayes
LogMgLik <- function(theta, y){
  N <- length(y)
  #sample size
  a0 <- theta[1]
  # prior shape hyperparameter
  b0 <- theta[2]
  # prior scale hyperparameter
  aN <- sum(y) + a0
  # posterior shape parameter
  if(a0 <= 0 || b0 <= 0){
    #Avoiding negative values
    lnp <- -Inf
  }else{lnp <- lgamma(a0+b0) + lgamma(aN) +
    lgamma(b0+N-sum(y)) -lgamma(a0) -lgamma(b0) -
    lgamma(a0+b0+N)}
  # log marginal likelihood
  return(-lnp)
}

theta0 <- c(0.1, 0.1)
# Initial values
control <- list(maxit = 1000)
# Number of iterations in optimization
EmpBay <- optim(theta0, LogMgLik, method = "BFGS",
  control = control, hessian = TRUE, y = y)
# Optimization
EmpBay$convergence
# Checking convergence
EmpBay$value # Maximum
a0EB <- EmpBay$par[1]
# Prior shape using empirical Bayes
b0EB <- EmpBay$par[2]
# Prior scale using empirical Bayes

PriMeanEB <- a0EB / (a0EB + b0EB)
# Prior mean
PriVarEB <- (a0EB*b0EB) / (((a0EB+b0EB)^2)*(a0EB+b0EB+1))
# Prior variance

PpEB <- ProbBo(y, a0 = a0EB, b0 = b0EB)
# This setting is using empirical Bayes.
PpEB
0.47

```

R code. Car claim, density plots

```

# Density figures
lambda <- seq(0.001, 1, 0.001)
# Values of lambda
VaguePrior <- dbeta(lambda, shape1 = a0, shape2 = b0)
EBPrior <- dbeta(lambda, shape1 = a0EB, shape2 = b0EB)
PosteriorV <- dbeta(lambda, shape1 = a0 + sum(y),
  shape2 = b0 + N - sum(y))
PosteriorEB <- dbeta(lambda, shape1 = a0EB + sum(y),
  shape2 = b0EB + N - sum(y))

# Likelihood function
Likelihood <- function(theta, y){
  LogL <- dbinom(y, 1, theta, log = TRUE)
  # LogL <- dbern(y, theta)
  Lik <- prod(exp(LogL))
  return(Lik)
}

Liks <- sapply(lambda, function(par) {
  Likelihood(par, y = y)})
Sc <- max(PosteriorEB)/max(Liks)
#Scale for displaying in figure
LiksScale <- Liks * Sc

data <- data.frame(cbind(lambda, VaguePrior, EBPrior,
  PosteriorV, PosteriorEB, LiksScale))
#Data frame

require(ggplot2)
# Cool figures
require(latex2exp)
# LaTeX equations in figures
require(ggpubr)
# Multiple figures in one page

fig1 <- ggplot(data = data, aes(lambda, VaguePrior)) +
  geom_line() +
  xlab(TeX("$p$")) + ylab("Density") +
  ggtitle("Prior: Vague Beta")

fig2 <- ggplot(data = data, aes(lambda, EBPrior)) +
  geom_line() +
  xlab(TeX("$p$")) + ylab("Density") +
  ggtitle("Prior: Empirical Bayes Beta")

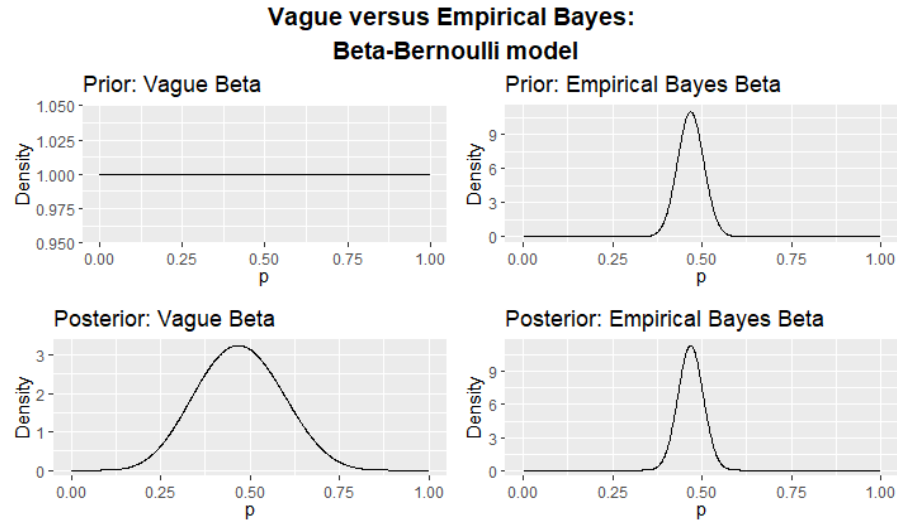
fig3 <- ggplot(data = data, aes(lambda, PosteriorV)) +
  geom_line() +
  xlab(TeX("$p$")) + ylab("Density") +
  ggtitle("Posterior: Vague Beta")

fig4 <- ggplot(data = data, aes(lambda, PosteriorEB)) +
  geom_line() +
  xlab(TeX("$p$")) + ylab("Density") +
  ggtitle("Posterior: Empirical Bayes Beta")

FIG <- ggarrange(fig1, fig2, fig3, fig4,
  ncol = 2, nrow = 2)

annotate_figure(FIG,
  top = text_grob("Vague-versus-Empirical-Bayes:
  Beta-Bernoulli-model", color = "black",
  face = "bold", size = 14))

```

**FIGURE 1.1**

Vague versus Empirical Bayes: Bernoulli-Beta model.

*R code. Car claim, prior, likelihood and posterior
density plots*

```
# Prior, likelihood and posterior:
# Empirical Bayes Binomial-Beta model
dataNew <- data.frame(cbind(rep(lambda, 3),
c(EBPrior, PosteriorEB, Likelihood),
rep(1:3, each = 1000)))
# Data frame

colnames(dataNew) <- c("Lambda", "Density", "Factor")
dataNew$Factor <- factor(dataNew$Factor, levels=c("1", "3",
"2"), labels=c("Prior", "Likelihood", "Posterior"))

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**

Prior, likelihood and posterior: Bernoulli-Beta model.

R code. Car claim, predictive probabilities plots

```
# Predictive distributions
require(TailRank)

PredDen <- function(y, y0, a0, b0){
  N <- length(y)
  aN <- a0 + sum(y) # Posterior shape parameter
  bN <- b0 + N - sum(y) # Posterior scale parameter
  Pr <- aN/(aN+bN)
  Probs <- dbinom(y0, 1, prob = Pr)
  return(Probs)
}
y0 <- 0:1
PredVague <- PredDen(y = y, y0 = y0, a0 = a0, b0 = b0)
PredEB <- PredDen(y = y, y0 = y0, a0 = a0EB, b0 = b0EB)
dataPred <- as.data.frame(cbind(y0, PredVague, PredEB))
colnames(dataPred) <- c("y0", "PredictiveVague",
  "PredictiveEB")

ggplot(data = dataPred) +
  geom_point(aes(y0, PredictiveVague, color = "red")) +
  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(0,1,by=1))
```

**FIGURE 1.3**

Predictive probabilities: Bernoulli-Beta model.

R code. Car claim, Bayesian model average

```

# Posterior odds: Vague vs Empirical Bayes
PO12 <- exp(-LogMgLik(c(a0EB, b0EB), y = y)) /
exp(-LogMgLik(c(a0, b0), y = y))

PostProMEM <- PO12 / (1 + PO12)
# Posterior model probability Empirical Bayes
PostProMEM
0.757
PostProbMV <- 1 - PostProMEM
# Posterior model probability vague prior
PostProbMV
0.242

# Bayesian model average (BMA)
PostMeanEB <- (a0EB + sum(y)) / (a0EB + b0EB + N)
# Posterior mean Empirical Bayes
PostMeanV <- (a0 + sum(y)) / (a0 + b0 + N)
# Posterior mean vague priors
BMAMean <- PostProMEM * PostMeanEB + PostProbMV * PostMeanV
# BMA posterior mean

PostVarEB <- (a0EB + sum(y)) * (b0EB + N - sum(y)) /
((a0EB + b0EB + N)^2 * (a0EB + b0EB + N - 1))
# Posterior variance Empirical Bayes
PostVarV <- (a0 + sum(y)) * (b0 + N - sum(y)) /
((a0 + b0 + N)^2 * (a0 + b0 + N - 1))
# Posterior variance vague prior

BMAVar <- PostProMEM * PostVarEB + PostProbMV * PostVarV
+ PostProMEM * (PostMeanEB - BMAMean)^2 + PostProbMV *
(PostMeanV - BMAMean)^2
# BMA posterior variance

# BMA: Predictive
BMAPred <- PostProMEM * PredEB + PostProbMV * PredVague
dataPredBMA <- as.data.frame(cbind(y0, BMAPred))
colnames(dataPredBMA) <- c("y0", "PredictiveBMA")

ggplot(data = dataPredBMA) +
  geom_point(aes(y0, PredictiveBMA, color = "red")) +
  xlab(TeX("$y_0$")) + ylab("Density") +
  ggtitle("Predictive density: BMA") +
  guides(color = guide_legend(title="BMA")) +
  scale_color_manual(labels = c("Probability"),
    values = c("red"))
+ scale_x_continuous(breaks=seq(0,1,by=1))

```

R code. Car claim, Bayesian updating plots

```

# Bayesian updating
BayUp <- function(y, lambda, a0, b0){
  N <- length(y)
  aN <- a0 + sum(y)
  # Posterior shape parameter
  bN <- b0 + N - sum(y)
  # Posterior scale parameter
  p <- dbeta(lambda, shape1 = aN, shape2 = bN)
  # Posterior density
  return(list(Post = p, a0New = aN, b0New = bN))
}

PostUp <- NULL
for(i in 1:N){
  if(i == 1){
    PostUpi <- BayUp(y[i], lambda, a0 = 1, b0 = 1)
  } else {
    PostUpi <- BayUp(y[i], lambda,
      a0 = PostUpi$a0New, b0 = PostUpi$b0New)
  }
  PostUp <- cbind(PostUp, PostUpi$Post)
}

DataUp <- data.frame(cbind(rep(lambda, 15), c(PostUp),
rep(1:15, each = 1000)))
#Data frame
colnames(DataUp) <- c("Lambda", "Density", "Factor")
DataUp$Factor <- factor(DataUp$Factor, levels=c("1","2",
"3","4","5","6","7","8","9","10","11","12","13","14","15"),
labels=c("Iter_1","Iter_2","Iter_3","Iter_4","Iter_5",
"Iter_6","Iter_7","Iter_8","Iter_9","Iter_10","Iter_11",
"Iter_12","Iter_13","Iter_14","Iter_15"))

ggplot(data = DataUp, aes_string(x = "Lambda",
y = "Density", group = "Factor")) +
  geom_line(aes(color = Factor)) +
  xlab(TeX("$p$")) + ylab("Density") +
  ggtitle("Bayesian updating:
_Beta-Binomial_model_with_vague_prior") +
  guides(color=guide_legend(title="Update"))

```

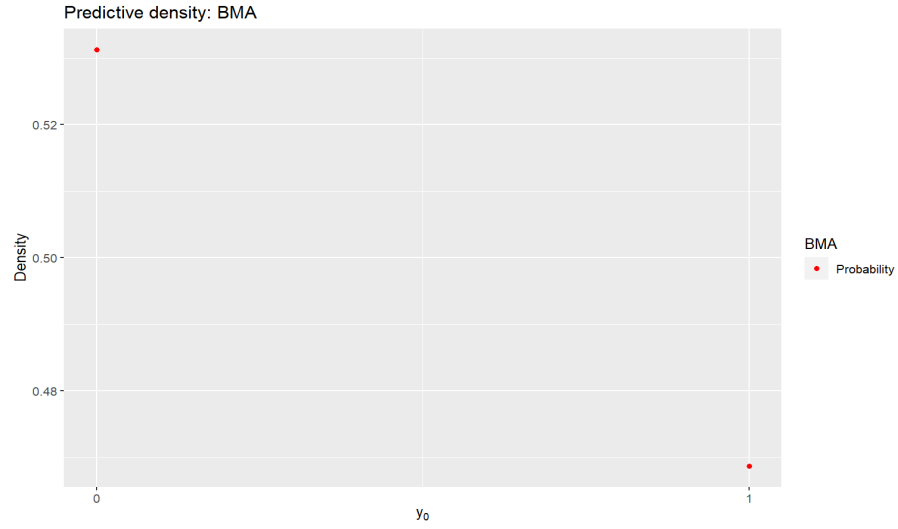
5. Show that given the loss function, $L(\theta, a) = |\theta - a|$, then the optimal decision rule minimizing the risk function, $a^*(\mathbf{y})$, is the median.

sAnswer

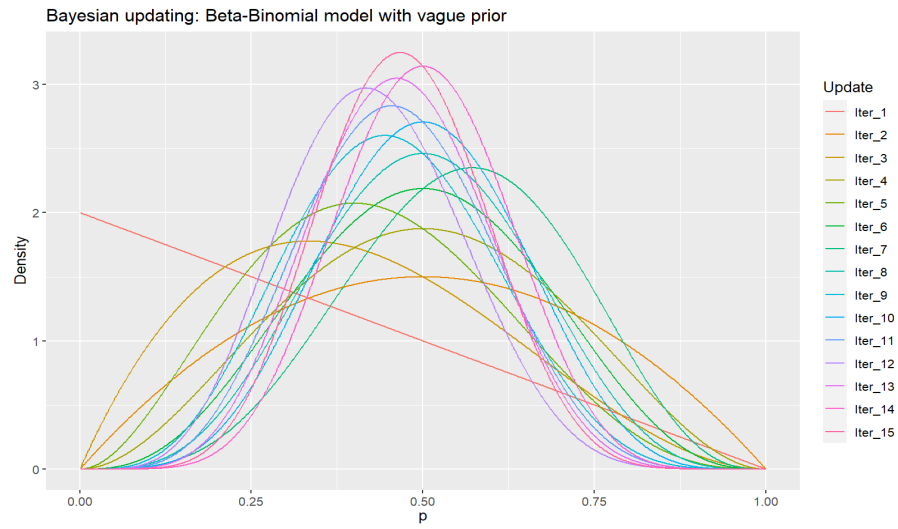
$\int_{\Theta} |\theta - a| \pi(\theta|\mathbf{y}) d\theta = \int_{-\infty}^a (a - \theta) \pi(\theta|\mathbf{y}) d\theta + \int_a^{\infty} (\theta - a) \pi(\theta|\mathbf{y}) d\theta$. Differentiating with respect to a , and equating to zero,

$$\int_{-\infty}^a \pi(\theta|\mathbf{y}) d\theta = \int_a^{\infty} \pi(\theta|\mathbf{y}) d\theta, \quad (1.13)$$

then,

**FIGURE 1.4**

Predictive probabilities: Bernoulli-Beta Bayesian model average.

**FIGURE 1.5**

Predictive probabilities: Bernoulli-Beta Bayesian model updating.

$$2 \int_{-\infty}^a \pi(\theta|\mathbf{y})d\theta = \int_{-\infty}^{\infty} \pi(\theta|\mathbf{y})d\theta = 1, \quad (1.14)$$

that is, $a^*(\mathbf{y})$ is the median.



2

Solutions of chapter 2

Conceptual differences: Bayesian and Frequentist approaches

2.1 Solutions of Exercises

1. Jeffreys-Lindley's paradox

The **Jeffreys-Lindley's paradox** [1, 3] 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 $\mathcal{U}(0, 1)$ under H_1 .
- Show that the p -value for this hypothesis test is equal to 0.023 using the normal approximation, $Y \sim \mathcal{N}(N \times \theta, N \times \theta \times (1 - \theta))$.

Answer

- The marginal likelihood under the null hypothesis is $p(y|H_0) = \binom{n}{y} \theta^y (1 - \theta)^{n-y} \approx 1.95 \times 10^{-4}$ given $\theta = 0.5$ under H_0 , $N = 49,581 + 48,870$ and $y = 49,581$. On the other hand, the marginal likelihood under the alternative hypothesis is

$$\begin{aligned}
 p(y|H_1) &= \int_0^1 \binom{n}{y} \theta^y (1-\theta)^{n-y} d\theta \\
 &= \binom{n}{y} B(y+1, n-k+1) \\
 &= \frac{\Gamma(N+1)}{\Gamma(y+1)\Gamma(n-y+1)} \frac{\Gamma(y+1)\Gamma(N-y+1)}{\Gamma(N+2)} \\
 &= \frac{N!}{(N+1)!} \\
 &= \frac{1}{N+1} \\
 &\approx 1.016 \times 10^{-5}.
 \end{aligned}$$

Then, $PO_{01} = \frac{1.95 \times 10^{-4}}{1.016 \times 10^{-5}} = 19.19$, this implies that the posterior model probability under the null hypothesis is $\pi(H_0|y) = \frac{19.19}{1+19.19} = 0.95$.

- Under the null hypothesis,

$$\begin{aligned}
 p &= 2 \int_{49,581}^{\infty} (2\pi\sigma^2)^{-1/2} \exp\left\{-\frac{1}{2\sigma^2}(y-\mu)^2\right\} dy \\
 &= 0.0235,
 \end{aligned}$$

where $\mu = N \times \theta = 49,225.5$, and $\sigma^2 = N \times \theta \times (1-\theta) = 24,612.75$ under the null hypothesis ($\theta = 0.5$).

Observe that the posterior model probability supports the null hypothesis, whereas the p-value implies rejection of the null hypothesis using a 5% significance level.

Observe that actually this is not a paradox, as we are answering two different questions. The Bayes factor is comparing two models ($\theta = 0.5$ versus $\theta \sim \mathcal{U}(0,1)$), whereas the p-value is checking the compatibility between $\theta = 0.5$ and the sample information. Despite that $\theta = 0.5$ is not compatible with sample information, it is better than the models assuming $\theta \sim \mathcal{U}(0,1)$ as most of these values of θ are far away from the sample mean. Thus, the model under the null is a bad description of the data, but it is better than the model under the alternative hypothesis.¹

¹Observe that there are at least another two issues in this example. First, the prior under the alternative is non-informative, this implies problems for Bayes factors, and second, the prior under the alternative is positive at $\theta = 0.5$, which is the null ([2] propose non-local prior densities in Bayesian hypothesis tests to tackle these issues).

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

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

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

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

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

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.

Answer

$$\begin{aligned} p(\mathbf{y}|\mathcal{M}_1) &= \int_{-\infty}^{\infty} \int_0^{\infty} (2\pi)^{-N/2} \sigma^{-N} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu)^2 \right\} \frac{1}{\sigma} d\sigma d\mu \\ &= (2\pi)^{-N/2} \int_{-\infty}^{\infty} \int_0^{\infty} \sigma^{-(N+1)} \exp \left\{ -\frac{N}{2\sigma^2} \frac{\sum_{i=1}^N (y_i - \mu)^2}{N} \right\} d\sigma d\mu \\ &= (2\pi)^{-N/2} \frac{\Gamma(N/2)}{2} 2^{N/2} \int_{-\infty}^{\infty} \left[\sum_{i=1}^N (y_i - \mu)^2 \right]^{-N/2} d\mu \\ &= (2\pi)^{-N/2} \frac{\Gamma(N/2)}{2} 2^{N/2} \int_{-\infty}^{\infty} \left[\sum_{i=1}^N [(y_i - \bar{y}) - (\mu - \bar{y})]^2 \right]^{-N/2} d\mu \\ &= (2\pi)^{-N/2} \frac{\Gamma(N/2)}{2} 2^{N/2} \int_{-\infty}^{\infty} [\alpha_n \hat{\sigma}^2 + N(\mu - \bar{y})^2]^{-N/2} d\mu \\ &= (2\pi)^{-N/2} \frac{\Gamma(N/2)}{2} 2^{N/2} \left(\frac{\alpha_n \hat{\sigma}^2}{\alpha_n \hat{\sigma}^2} \right)^{-N/2} \int_{-\infty}^{\infty} [\alpha_n \hat{\sigma}^2 + N(\mu - \bar{y})^2]^{-N/2} d\mu \\ &= (2\pi)^{-N/2} \frac{\Gamma(N/2)}{2} 2^{N/2} (\alpha_n \hat{\sigma}^2)^{-N/2} \int_{-\infty}^{\infty} \left[1 + \frac{N(\mu - \bar{y})^2}{\alpha_n \hat{\sigma}^2} \right]^{-N/2} d\mu \\ &= \frac{\pi^{-N/2}}{2} \Gamma(N/2) 2^{N/2} \left(\frac{1}{\alpha_n \hat{\sigma}^2} \right)^{N/2} \left(\frac{N}{\alpha_n \hat{\sigma}^2} \right)^{-1/2} \frac{\Gamma(1/2) \Gamma(\alpha_n/2)}{\Gamma((\alpha_n + 1)/2)}. \end{aligned}$$

The third line takes into account that the integral in the second line is the kernel of an inverted-gamma distribution, and the last line takes into account that the integral in the previous line is the kernel of a student's t distribution [6].

$$\begin{aligned}
 p(\mathbf{y}|\mathcal{M}_0) &= \int_0^\infty (2\pi)^{-N/2} \sigma^{-N} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu_0)^2 \right\} \frac{1}{\sigma} d\sigma \\
 &= (2\pi)^{-N/2} \int_0^\infty \sigma^{-(N+1)} \exp \left\{ -\frac{N}{2\sigma^2} \frac{\sum_{i=1}^N (y_i - \mu_0)^2}{N} \right\} d\sigma \\
 &= (2\pi)^{-N/2} \left[\frac{2}{\Gamma(N/2)} \left(\frac{N}{2} \frac{\sum_{i=1}^N (y_i - \mu_0)^2}{N} \right)^{N/2} \right]^{-1}.
 \end{aligned}$$

The third line takes into account that the integral in the second line is the kernel of an inverted-gamma distribution [6].

Given these results is easy to get PO_{01} .

In addition,

$$\begin{aligned}
 PO_{01} &= \frac{\Gamma((\alpha_n + 1)/2)}{\Gamma(1/2)\Gamma(\alpha_N/2)} (\alpha_n \hat{\sigma}^2/N)^{-1/2} \left[1 + \frac{(\mu_0 - \bar{y})^2}{\alpha_n \hat{\sigma}^2/N} \right]^{-\left(\frac{\alpha_n+1}{2}\right)} \\
 &= \frac{\Gamma((\alpha_n + 1)/2)}{\Gamma(1/2)\Gamma(\alpha_N/2)} (\alpha_n \hat{\sigma}^2/N)^{-1/2} \left[1 + \frac{1}{\alpha_n} \left(\frac{\mu_0 - \bar{y}}{\hat{\sigma}/\sqrt{N}} \right)^2 \right]^{-\left(\frac{\alpha_n+1}{2}\right)} \\
 &= \frac{\Gamma((\alpha_n + 1)/2)}{\Gamma(1/2)\Gamma(\alpha_N/2)} (\alpha_n \hat{\sigma}^2/N)^{-1/2} \left[1 + \frac{1}{\alpha_n} t^2 \right]^{-\left(\frac{\alpha_n+1}{2}\right)},
 \end{aligned}$$

where $t = \frac{\bar{y} - \mu_0}{\hat{\sigma}/\sqrt{N}}$ is the classical statistical test. Then, as t increases then the PO_{01} decreases, both indicating support against the null hypothesis H_0 . $\mu = \mu_0$. However, there are other terms affecting the posterior odds, then, there is no necessary agreement between the classical test statistic and the posterior odds.

3. Using the setting of the **Example: Math test** in subsection 2.6.1 in the book, 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?
- Find the posterior model probability of the null model for each μ_0 .

R code. Example: Math test

```
N <- 50 # Sample size
y_bar <- 102 # Sample mean
s2 <- 10 # Sample variance
alpha <- N - 1
serror <- (s2/N)^0.5
y.H0 <- c(100, 100.5, 101, 101.5, 102)
test <- (y.H0 - y_bar)/serror
pval <- 2*pt(test, alpha)
pval
0.0000459 0.0015431 0.0299338 0.2690040 1
# p-values
PO01 <- (gamma(N/2)*((N-1)*serror^2)^(-0.5)*
(1+test^2/alpha)^(-N/2))/(gamma(1/2)*gamma((N-1)/2))
PO01/(1+PO01)
0.0001705 0.0050345 0.0725330 0.3210223 0.4702050
# Posterior model probability of the null hypothesis.
```



3

Solutions of chapter 4 Cornerstone models: Conjugate families

3.1 Solutions of Exercises

1. Write in the canonical form the distribution of the Bernoulli example, and find the mean and variance of the sufficient statistic.

Answer

Given $p(\mathbf{y}|\theta) = (1-\theta)^N \exp \left\{ \sum_{i=1}^N y_i \log \left(\frac{\theta}{1-\theta} \right) \right\}$ where $\eta = \log \frac{\theta}{1-\theta}$ which implies $\theta = \frac{\exp(\eta)}{1+\exp(\eta)}$, then $p(\mathbf{y}|\theta) = \exp \left\{ \sum_{i=1}^N y_i \eta - N \log(1 + \exp(\eta)) \right\}$. Thus $B(\eta) = N \log(1 + \exp(\eta))$, $\nabla(B(\eta)) = N \frac{\exp(\eta)}{1+\exp(\eta)} = N\theta$ and $\nabla^2(B(\eta)) = N \left\{ \frac{\exp(\eta)(1+\exp(\eta))}{(1+\exp(\eta))^2} - \frac{\exp(\eta)\exp(\eta)}{(1+\exp(\eta))^2} \right\} = N\theta(1-\theta)$.

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.

Answer

The density function is

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

Observe that $\sum_{i=1}^N n_i$ is the total sample size of Bernoulli experiments.

Using Theorem 1 in Chapter 4, the prior distribution is

$$\begin{aligned}\pi(\theta) &\propto (1-\theta)^{B_0} \exp \left\{ a_0 \log \left(\frac{\theta}{1-\theta} \right) \right\} \\ &= \theta^{a_0} (1-\theta)^{B_0-a_0} \\ &= \theta^{\alpha_0-1} (1-\theta)^{\beta_0-1},\end{aligned}$$

where $\alpha_0 = a_0 + 1$ and $\beta_0 = B_0 - a_0 + 1$. This is the kernel of a beta distribution. Thus, the posterior distribution is

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

where $\alpha_n = \alpha_0 + \sum_{i=1}^N y_i$ and $\beta_n = \beta_0 + \sum_{i=1}^N n_i - \sum_{i=1}^N y_i$.

The marginal likelihood is

$$\begin{aligned}p(\mathbf{y}) &= \int_0^1 \frac{\theta^{\alpha_0-1} (1-\theta)^{\beta_0-1}}{B(\alpha_0, \beta_0)} \times \prod_{i=1}^N \binom{n_i}{y_i} \theta^{\sum_{i=1}^N y_i} (1-\theta)^{\sum_{i=1}^N n_i - \sum_{i=1}^N y_i} d\theta \\ &= \frac{\prod_{i=1}^N \binom{n_i}{y_i}}{B(\alpha_0, \beta_0)} \int_0^1 \theta^{\alpha_0 + \sum_{i=1}^N y_i - 1} (1-\theta)^{\beta_0 + \sum_{i=1}^N n_i - \sum_{i=1}^N y_i - 1} d\theta \\ &= \frac{\prod_{i=1}^N \binom{n_i}{y_i} B(\alpha_n, \beta_n)}{B(\alpha_0, \beta_0)}.\end{aligned}$$

The third line due to having the kernel of a Beta distribution.

Finally, the predictive distribution is

$$\begin{aligned}p(Y_0|\mathbf{y}) &= \int_0^1 \binom{n_{y_0}}{y_0} \theta^{y_0} (1-\theta)^{n_{y_0}-y_0} \frac{\theta^{\alpha_n-1} (1-\theta)^{\beta_n-1}}{B(\alpha_n, \beta_n)} d\theta \\ &= \frac{\binom{n_{y_0}}{y_0}}{B(\alpha_n, \beta_n)} \int_0^1 \theta^{\alpha_n+y_0-1} (1-\theta)^{\beta_n+n_{y_0}-y_0-1} d\theta \\ &= \binom{n_{y_0}}{y_0} \frac{B(\alpha_n+y_0, \beta_n+n_{y_0}-y_0)}{B(\alpha_n, \beta_n)},\end{aligned}$$

where n_{y_0} is the known size associated with y_0 , and the last line due to having the kernel of a beta distribution. The predictive is a *beta-binomial distribution*.

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.

Answer

We see that the exponential distribution belongs to the exponential family as $p(\mathbf{y}|\lambda) = \prod_{i=1}^N \lambda \exp(-\lambda y_i) = \lambda^N \exp(-\lambda \sum_{i=1}^N y_i)$.

Using the gamma distribution in the rate parametrization, we see that $\pi(\lambda|\mathbf{y}) \propto \lambda^{\alpha_0-1} \exp(-\lambda\beta_0) \times \lambda^N \exp(-\lambda \sum_{i=1}^N y_i) = \lambda^{\alpha_0+N-1} \exp(-\lambda(\beta_0 + \sum_{i=1}^N y_i))$. This is the kernel of a gamma distribution, that is, $\lambda|\mathbf{y} \sim G(\alpha_n, \beta_n)$ where $\alpha_n = \alpha_0 + N$ and $\beta_n = \beta_0 + \sum_{i=1}^N y_i$.

The marginal likelihood is

$$\begin{aligned} p(\mathbf{y}) &= \int_0^\infty \lambda^N \exp\left\{-\lambda \sum_{i=1}^N y_i\right\} \lambda^{\alpha_0-1} \exp\{-\beta_0 \lambda\} \frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} d\lambda \\ &= \frac{\beta_0^{\alpha_0}}{\Gamma(\alpha_0)} \int_0^\infty \lambda^{\alpha_0+N-1} \exp\left\{-\lambda \left(\beta_0 + \sum_{i=1}^N y_i\right)\right\} d\lambda \\ &= \frac{\beta_0^{\alpha_0} \Gamma(\alpha_n)}{\Gamma(\alpha_0) \beta_n^{\alpha_n}}. \end{aligned}$$

Finally, the predictive distribution is

$$\begin{aligned} p(Y_0|\mathbf{y}) &= \int_0^\infty \lambda \exp\{-\lambda y_0\} \lambda^{\alpha_n-1} \exp\{-\beta_n \lambda\} \frac{\beta_n^{\alpha_n}}{\Gamma(\alpha_n)} d\lambda \\ &= \frac{\beta_n^{\alpha_n}}{\Gamma(\alpha_n)} \int_0^\infty \lambda^{\alpha_n+1-1} \exp\{-\lambda(\beta_n + y_0)\} d\lambda \\ &= \frac{\beta_n^{\alpha_n}}{\Gamma(\alpha_n)} \times \frac{\Gamma(\alpha_n + 1)}{(\beta_n + y_0)^{\alpha_n+1}} \\ &= \frac{\alpha_n \beta_n^{\alpha_n}}{(\beta_n + y_0)^{\alpha_n+1}}. \end{aligned}$$

This is a *Lomax distribution*.

4. Given $\mathbf{y} \sim N_N(\mu, \Sigma)$, that is, a *multivariate normal distribution* show that $p(\mathbf{y}|\mu, \Sigma)$ is in the exponential family.

Answer

$$\begin{aligned}
p(\mathbf{y}|\mu, \Sigma) &= (2\pi)^{-N/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mu)^\top \Sigma^{-1} (\mathbf{y} - \mu) \right\} \\
&= (2\pi)^{-N/2} \exp \left\{ -\frac{1}{2} (\mathbf{y}^\top \Sigma^{-1} \mathbf{y} - 2\mathbf{y}^\top \Sigma^{-1} \mu + \mu^\top \Sigma^{-1} \mu + \log(|\Sigma|)) \right\} \\
&= (2\pi)^{-N/2} \exp \left\{ -\frac{1}{2} (tr \{ \mathbf{y}^\top \Sigma^{-1} \mathbf{y} \} - 2\mathbf{y}^\top \Sigma^{-1} \mu + \mu^\top \Sigma^{-1} \mu + \log(|\Sigma|)) \right\} \\
&= (2\pi)^{-N/2} \exp \left\{ -\frac{1}{2} \left(vec(\mathbf{y} \mathbf{y}^\top)^\top vec(\Sigma^{-1}) - 2\mathbf{y}^\top \Sigma^{-1} \mu + \mu^\top \Sigma^{-1} \mu + \log(|\Sigma|) \right) \right\},
\end{aligned}$$

where tr and vec are the trace and vectorization operators, respectively.

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

5. Find the marginal likelihood in the normal/inverse-Wishart model.

Answer

$$\begin{aligned}
p(\mathbf{Y}) &= \int_{\mathcal{R}^p} \int_{\mathcal{S}} (2\pi)^{-pN/2} |\Sigma|^{-N/2} \exp \left\{ -\frac{1}{2} tr[(\mathbf{S} + N(\mu - \hat{\mu})(\mu - \hat{\mu})^\top) \Sigma^{-1}] \right\} \\
&\quad \times (2\pi)^{-p/2} \beta_0^{p/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{\beta_0}{2} tr[(\mu - \mu_0)(\mu - \mu_0)^\top \Sigma^{-1}] \right\} \\
&\quad \times |\Sigma|^{-(\alpha_0 + p + 1)/2} \frac{2^{-\alpha_0 p/2} |\Psi_0|^{\alpha_0/2}}{\Gamma_p(\alpha_0/2)} \exp \left\{ -\frac{1}{2} tr(\Psi_0 \Sigma^{-1}) \right\} d\Sigma d\mu \\
&= \frac{(2\pi)^{-frac{1}{2}(pN+p)} |\Psi_0|^{\alpha_0/2} \beta_0^{p/2} 2^{-\alpha_0 p/2}}{\Gamma_p(\alpha_0/2)} \int_{\mathcal{R}^p} \int_{\mathcal{S}} |\Sigma|^{-\frac{1}{2}(N+1+\alpha_0+p+1)} \\
&\quad \times \exp \left\{ -\frac{1}{2} tr[(\mathbf{S} + N(\mu - \hat{\mu})(\mu - \hat{\mu})^\top + \beta_0(\mu - \mu_0)(\mu - \mu_0)^\top + \Psi_0) \Sigma^{-1}] \right\} d\Sigma d\mu.
\end{aligned}$$

We have in the integral the kernel of an Inverse-Wishart distribution, then

$$\begin{aligned}
p(\mathbf{Y}) &= \frac{\Gamma_p\left(\frac{N+1+\alpha_0}{2}\right) |\Psi_0|^{\alpha_0/2} \beta_0^{p/2}}{\Gamma_p(\alpha_0/2) \pi^{p(N+1)/2}} \\
&\quad \times \int_{\mathcal{R}^p} |\mathbf{S} + \Psi_0 + (N + \beta_0)(\mu - \mu_n)(\mu - \mu_n)^\top \\
&\quad + N\beta_0/(N + \beta_0)(\hat{\mu} - \mu_0)(\hat{\mu} - \mu_0)^\top| d\mu \\
&= \frac{\Gamma_p\left(\frac{N+1+\alpha_0}{2}\right) |\Psi_0|^{\alpha_0/2} \beta_0^{p/2}}{\Gamma_p(\alpha_0/2) \pi^{p(N+1)/2}} \\
&\quad \times \int_{\mathcal{R}^p} |\Psi_n| |1 + \beta_n(\mu - \mu_n)^\top \Psi_n^{-1}(\mu - \mu_n)|^{-\frac{1}{2}(\alpha_n+1)} d\mu \\
&= \frac{\Gamma_p\left(\frac{\alpha_n+1}{2}\right) |\Psi_0|^{\alpha_0/2} \beta_0^{p/2}}{\Gamma_p(\alpha_0/2) \pi^{p(N+1)/2}} |\Psi_n|^{-\frac{1}{2}(\alpha_n+1)} \\
&\quad \times \int_{\mathcal{R}^p} [1 + \beta_n(\mu - \mu_n)^\top \Psi_n^{-1}(\mu - \mu_n)]^{-\frac{1}{2}(\alpha_n+1)} d\mu.
\end{aligned}$$

The last equality uses the definition of Ψ_n , β_n and α_n , and the Sylvester's determinant theorem. Observe that we have the kernel of a multivariate t distribution [4]. Then,

$$\begin{aligned}
p(\mathbf{Y}) &= \frac{\Gamma_p\left(\frac{\alpha_n+1}{2}\right) |\Psi_0|^{\alpha_0/2} \beta_0^{p/2}}{\Gamma_p(\alpha_0/2) \pi^{p(N+1)/2}} |\Psi_n|^{-\frac{1}{2}(\alpha_n+1)} \\
&\quad \times \int_{\mathcal{R}^p} \left[1 + \frac{1}{\alpha_n + 1 - p} (\mu - \mu_n)^\top \left(\frac{\Psi_n}{\beta_n(\alpha_n + 1 - p)}\right)^{-1} (\mu - \mu_n)\right]^{-\frac{1}{2}(\alpha_n+1-p+p)} d\mu \\
&= \frac{\Gamma_p\left(\frac{\alpha_n+1}{2}\right) \Gamma_p\left(\frac{\alpha_n+1-p}{2}\right) |\Psi_0|^{\alpha_0/2} \beta_0^{p/2} (\alpha_n + 1 - p)^{p/2} \pi^{p/2} |\Psi_n|^{-\frac{1}{2}(\alpha_n+1)}}{\Gamma_p(\alpha_0/2) \pi^{p(N+1)/2} \Gamma_p\left(\frac{\alpha_n+1-p+p}{2}\right) \left(\frac{\Psi_n}{\alpha_n+1-p}\right)^{-1/2}} \\
&= \frac{\Gamma_p\left(\frac{v_n}{2}\right) |\Psi_0|^{\alpha_0/2}}{\Gamma_p\left(\frac{\alpha_0}{2}\right) |\Psi_n|^{\alpha_n/2}} \left(\frac{\beta_0}{\beta_n}\right)^{p/2} (2\pi)^{-Np/2},
\end{aligned}$$

where $v_n = \alpha_n + 1 - p$.

6. Find the posterior predictive distribution in the normal/inverse-Wishart model, and show that $\mathbf{Y}_0|\mathbf{Y} \sim T_{N_0, M}(\alpha_n - M + 1, \mathbf{X}_0 \mathbf{B}_n, \mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{V}_n \mathbf{X}_0^\top, \Psi_n)$ in the multivariate regression linear model.

Answer

$$\begin{aligned}
p(\mathbf{Y}_0|\mathbf{Y}) &\propto \int_{\mathcal{R}^p} \int_S |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} \text{tr}[(\mathbf{y}_0 - \mu)(\mathbf{y}_0 - \mu)^\top \Sigma^{-1}] \right\} \\
&\quad \times |\Sigma|^{-1/2} \exp \left\{ -\frac{\beta_n}{2} \text{tr}[(\mu - \mu_n)(\mu - \mu_n)^\top \Sigma^{-1}] \right\} \\
&\quad \times |\Sigma|^{-(\alpha_n + p + 1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\Psi_n \Sigma^{-1}) \right\} d\Sigma d\mu \\
&\propto \int_{\mathcal{R}^p} |(\mathbf{y}_0 - \mu)(\mathbf{y}_0 - \mu)^\top + (\mu - \mu_n)(\mu - \mu_n)^\top + \Psi_n|^{-(\alpha_n + 2)/2} d\mu.
\end{aligned}$$

The last equality uses that there is the kernel of an Inverse Wishart distribution.

Taking into account that

$$\begin{aligned}
(\mathbf{y}_0 - \mu)(\mathbf{y}_0 - \mu)^\top + (\mu - \mu_n)(\mu - \mu_n)^\top &= (1 + \beta_n) \left(\mu - \frac{(\mathbf{y}_0 + \beta_n \mu_n)}{1 + \beta_n} \right) \left(\mu - \frac{(\mathbf{y}_0 + \beta_n \mu_n)}{1 + \beta_n} \right)^\top \\
&\quad + \frac{\beta_n}{1 + \beta_n} (\mathbf{y}_0 - \mu_n)(\mathbf{y}_0 - \mu_n)^\top.
\end{aligned}$$

Then,

$$\begin{aligned}
p(\mathbf{Y}_0|\mathbf{Y}) &\propto \int_{\mathcal{R}^p} |(\mathbf{y}_0 - \mu)(\mathbf{y}_0 - \mu)^\top + (\mu - \mu_n)(\mu - \mu_n)^\top + \Psi_n|^{-(\alpha_n + 2)/2} d\mu \\
&= \int_{\mathcal{R}^p} \left| (1 + \beta_n) \left(\mu - \frac{(\mathbf{y}_0 + \beta_n \mu_n)}{1 + \beta_n} \right) \left(\mu - \frac{(\mathbf{y}_0 + \beta_n \mu_n)}{1 + \beta_n} \right)^\top \right. \\
&\quad \left. + \frac{\beta_n}{1 + \beta_n} (\mathbf{y}_0 - \mu_n)(\mathbf{y}_0 - \mu_n)^\top + \Psi_n \right|^{-(\alpha_n + 2)/2} d\mu \\
&= \int_{\mathcal{R}^p} \left| \underbrace{\Psi_n + \frac{\beta_n}{1 + \beta_n} (\mathbf{y}_0 - \mu_n)(\mathbf{y}_0 - \mu_n)^\top}_{\Lambda_n} \right| \\
&\quad \left| 1 + (1 + \beta_n) \left(\mu - \frac{(\mathbf{y}_0 + \beta_n \mu_n)}{1 + \beta_n} \right)^\top \frac{1}{\alpha_n + 2 - p} \left(\frac{\Lambda_n}{\alpha_n + 2 - p} \right)^{-1} \left(\mu - \frac{(\mathbf{y}_0 + \beta_n \mu_n)}{1 + \beta_n} \right) \right|^{-(\alpha_n + 2 - p + p)/2} d\mu \\
&\propto \left| \Psi_n + \frac{\beta_n}{1 + \beta_n} (\mathbf{y}_0 - \mu_n)(\mathbf{y}_0 - \mu_n)^\top \right|^{-(\alpha_n + 2)/2} \\
&\quad \times \left| \Psi_n + \frac{\beta_n}{1 + \beta_n} (\mathbf{y}_0 - \mu_n)(\mathbf{y}_0 - \mu_n)^\top \right|^{1/2} \\
&= \left| \Psi_n + \frac{\beta_n}{1 + \beta_n} (\mathbf{y}_0 - \mu_n)(\mathbf{y}_0 - \mu_n)^\top \right|^{-(\alpha_n + 1)/2} \\
&\propto \left[1 + (\mathbf{y}_0 - \mu_n)^\top \frac{1}{\alpha_n + 1 - p} \left(\frac{\Psi_n (1 + \beta_n)}{(\alpha_n + 1 - p) \beta_n} \right)^{-1} (\mathbf{y}_0 - \mu_n) \right]^{-(\alpha_n + 1 - p + p)}.
\end{aligned}$$

The second equality and last line use the Sylvester's determinant theorem, and the second equality uses that there is the kernel of a multivariate t distribution.

Then, we have that the predictive distribution is a multivariate t distribution centered at μ_n , $\alpha_n + 1 - p$ degrees of freedom, and scale matrix $\frac{\Psi_n(1+\beta_n)}{(\alpha_n+1-p)\beta_n}$.

To show the second statement, let's start by the definition of the predictive density to show that $\mathbf{Y}_0|\mathbf{Y} \sim T_{N_0,M}(\alpha_n - M + 1, \mathbf{X}_0\mathbf{B}_n, \mathbf{I}_{N_0} + \mathbf{X}_0\mathbf{V}_n\mathbf{X}_0^\top, \Psi_n)$.

$$\begin{aligned} \pi(\mathbf{Y}_0|\mathbf{Y}) &\propto \int_{\mathcal{S}} \int_{\mathcal{B}} \left\{ |\Sigma|^{-N_0/2} \exp \left\{ -\frac{1}{2} \text{tr}[(\mathbf{Y}_0 - \mathbf{X}_0\mathbf{B})^\top (\mathbf{Y}_0 - \mathbf{X}_0\mathbf{B}) \Sigma^{-1}] \right\} \right. \\ &\quad \times |\Sigma|^{-K/2} \exp \left\{ -\frac{1}{2} \text{tr}[(\mathbf{B} - \mathbf{B}_n)^\top \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) \Sigma^{-1}] \right\} \\ &\quad \times |\Sigma|^{-(\alpha_n+M+1)/2} \exp \left\{ -\frac{1}{2} \text{tr}[\Psi_n \Sigma^{-1}] \right\} \Big\} d\mathbf{B} d\Sigma \\ &= \int_{\mathcal{S}} \int_{\mathcal{B}} \left\{ |\Sigma|^{-(N_0+K+\alpha_n+M+1)/2} \exp \left\{ -\frac{1}{2} \text{tr} [((\mathbf{Y}_0 - \mathbf{X}_0\mathbf{B})^\top (\mathbf{Y}_0 - \mathbf{X}_0\mathbf{B}) \right. \right. \\ &\quad \left. \left. + (\mathbf{B} - \mathbf{B}_n)^\top \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n) + \Psi_n) \Sigma^{-1}] \right\} \right\} d\mathbf{B} d\Sigma. \end{aligned}$$

Setting $\mathbf{M} = (\mathbf{X}_0^\top \mathbf{X}_0 + \mathbf{V}_n^{-1})$, and $\mathbf{B}_* = \mathbf{M}^{-1}(\mathbf{V}_n \mathbf{B}_n + \mathbf{X}_0^\top \mathbf{Y}_0)$, we have that $(\mathbf{B} - \mathbf{B}_*)^\top \mathbf{M} (\mathbf{B} - \mathbf{B}_*) + \mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n + \mathbf{Y}_0^\top \mathbf{Y}_0 - \mathbf{B}_*^\top \mathbf{M} \mathbf{B}_* = (\mathbf{Y}_0 - \mathbf{X}_0\mathbf{B})^\top (\mathbf{Y}_0 - \mathbf{X}_0\mathbf{B}) + (\mathbf{B} - \mathbf{B}_n)^\top \mathbf{V}_n^{-1} (\mathbf{B} - \mathbf{B}_n)$. Then,

$$\begin{aligned} \pi(\mathbf{Y}_0|\mathbf{Y}) &\propto \int_{\mathcal{S}} |\Sigma|^{-(N_0+K+\alpha_n+M+1)/2} \\ &\quad \times \exp \left\{ -\frac{1}{2} \text{tr}[(\Psi_n + \mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n + \mathbf{Y}_0^\top \mathbf{Y}_0 - \mathbf{B}_*^\top \mathbf{M} \mathbf{B}_*) \Sigma^{-1}] \right\} \\ &\quad \times \int_{\mathcal{B}} \exp \left\{ -\frac{1}{2} \text{tr}[(\mathbf{B} - \mathbf{B}_*)^\top \mathbf{M} (\mathbf{B} - \mathbf{B}_*) \Sigma^{-1}] \right\} d\mathbf{B} d\Sigma. \end{aligned}$$

The latter is the kernel of a matrix normal distribution, thus

$$\begin{aligned} \pi(\mathbf{Y}_0|\mathbf{Y}) &\propto \int_{\mathcal{S}} |\Sigma|^{-(N_0+\alpha_n+M+1)/2} \\ &\quad \times \exp \left\{ -\frac{1}{2} \text{tr}[(\Psi_n + \mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n + \mathbf{Y}_0^\top \mathbf{Y}_0 - \mathbf{B}_*^\top \mathbf{M} \mathbf{B}_*) \Sigma^{-1}] \right\} d\Sigma \end{aligned}$$

This is the kernel of an inverse-Wishart distribution, then

$$\pi(\mathbf{Y}_0|\mathbf{Y}) \propto |\Psi_n + \mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n + \mathbf{Y}_0^\top \mathbf{Y}_0 - \mathbf{B}_*^\top \mathbf{M} \mathbf{B}_*|^{-(N_0+\alpha_n)/2}.$$

Setting $\mathbf{C}^{-1} = \mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{V}_n \mathbf{X}_0^\top$ such that $\mathbf{C} = \mathbf{I}_{N_0} - \mathbf{X}_0 (\mathbf{X}_0^\top \mathbf{X}_0 + \mathbf{V}_n^{-1})^{-1} \mathbf{X}_0^\top$ (see footnote 4 in Chapter 4), then $\mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n + \mathbf{Y}_0^\top \mathbf{Y}_0 - \mathbf{B}_*^\top \mathbf{M} \mathbf{B}_* = (\mathbf{Y}_0 - \mathbf{X}_0 \mathbf{B}_n)^\top \mathbf{C} (\mathbf{Y}_0 - \mathbf{X}_0 \mathbf{B}_n)$. This is done following exactly same procedure as deducing the predictive distribution in the linear regression model in the book. Thus,

$$\begin{aligned} \pi(\mathbf{Y}_0 | \mathbf{Y}) &\propto |\boldsymbol{\Psi}_n + (\mathbf{Y}_0 - \mathbf{X}_0 \mathbf{B}_n)^\top \mathbf{C} (\mathbf{Y}_0 - \mathbf{X}_0 \mathbf{B}_n)|^{-(N_0 + \alpha_n)/2} \\ &\propto |\mathbf{I}_{N_0} + \mathbf{C} (\mathbf{Y}_0 - \mathbf{X}_0 \mathbf{B}_n) \boldsymbol{\Psi}^{-1} (\mathbf{Y}_0 - \mathbf{X}_0 \mathbf{B}_n)^\top|^{-(\alpha_n + 1 - M + N_0 + M - 1)/2}. \end{aligned}$$

The second proportionality follows from the Sylvester's theorem. Observe that this is the kernel of a matrix t distribution with $\alpha_n + 1 - M$ degrees of freedom, location $\mathbf{X}_0 \mathbf{B}_n$ and scale matrices $\boldsymbol{\Psi}_n$ and $\mathbf{C}^{-1} = \mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{V}_n \mathbf{X}_0^\top$.

7. Show that $\delta_n = \delta_0 + (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) + (\hat{\beta} - \beta_0)^\top ((\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{B}_0)^{-1} (\hat{\beta} - \beta_0)$ 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.

Answer

Taking into account that

$$\begin{aligned} \delta^* &= \delta_0 + \mathbf{y}^\top \mathbf{y} + \beta_0^\top \mathbf{B}_0^{-1} \beta_0 - \beta_n^\top \mathbf{B}_n^{-1} \beta_n \\ &= \delta_0 + \mathbf{y}^\top \mathbf{y} + \beta_0^\top \mathbf{B}_0^{-1} \beta_0 - (\mathbf{B}_0^{-1} \beta_0 + \mathbf{X}^\top \mathbf{X} \hat{\beta})^\top \mathbf{B}_n (\mathbf{B}_0^{-1} \beta_0 + \mathbf{X}^\top \mathbf{X} \hat{\beta}) \\ &= \delta_0 + \mathbf{y}^\top \mathbf{y} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \mathbf{B}_n \mathbf{X}^\top \mathbf{X} \hat{\beta} - 2\hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \mathbf{B}_n \mathbf{B}_0^{-1} \beta_0 + \beta_0^\top (\mathbf{B}_0^{-1} - \mathbf{B}_0^{-1} \mathbf{B}_n \mathbf{B}_0^{-1}) \beta_0 \\ &\quad - \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta} + \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta} \\ &= \delta_0 + \mathbf{y}^\top \mathbf{y} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta} + \hat{\beta}^\top (\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X} \mathbf{B}_n \mathbf{X}^\top \mathbf{X}) \hat{\beta} \\ &\quad - 2\hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \mathbf{B}_n \mathbf{B}_0^{-1} \beta_0 + \beta_0^\top (\mathbf{B}_0^{-1} - \mathbf{B}_0^{-1} \mathbf{B}_n \mathbf{B}_0^{-1}) \beta_0. \end{aligned}$$

Observe that

$$\begin{aligned} (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) &= \mathbf{y}^\top \mathbf{y} - 2\hat{\beta}^\top \mathbf{X}^\top \mathbf{y} + \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta} \\ &= \mathbf{y}^\top \mathbf{y} - 2\hat{\beta}^\top \mathbf{X}^\top (\mathbf{X}\hat{\beta} + \hat{\mu}) + \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta} \\ &= \mathbf{y}^\top \mathbf{y} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{X} \hat{\beta}, \end{aligned}$$

where $\mathbf{y} = \mathbf{X}\hat{\beta} + \hat{\mu}$, and $\mathbf{X}^\top \hat{\mu} = 0$.

The following matrix identities are useful [5]:

$$(\mathbf{D} + \mathbf{E})^{-1} = \mathbf{D}^{-1} - \mathbf{D}^{-1} (\mathbf{D}^{-1} + \mathbf{E}^{-1})^{-1} \mathbf{D}^{-1},$$

and

$$(\mathbf{D} + \mathbf{E})^{-1} = \mathbf{D}^{-1} (\mathbf{E}^{-1} + \mathbf{D}^{-1}) \mathbf{E}^{-1}.$$

Using these identities,

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

Then,

$$\begin{aligned} \delta^* &= \delta_0 + (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) + \hat{\beta}^\top [(\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{B}_0]^{-1} \hat{\beta} \\ &\quad - 2\hat{\beta}^\top [(\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{B}_0]^{-1} \beta_0 + \beta_0^\top [(\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{B}_0]^{-1} \beta_0 \\ &= \delta_0 + (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) \\ &\quad + (\hat{\beta} - \beta_0)^\top [(\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{B}_0]^{-1} (\hat{\beta} - \beta_0). \end{aligned}$$

In a similar way for the second part,

$$\begin{aligned} (\mathbf{V}_0 + (\mathbf{X}^\top \mathbf{X})^{-1})^{-1} &= \mathbf{V}_0^{-1} - \mathbf{V}_0^{-1}(\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{V}_0^{-1} \\ &= \mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X}(\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \\ &= \mathbf{X}^\top \mathbf{X}((\mathbf{X}^\top \mathbf{X})^{-1} + \mathbf{V}_0)^{-1} \mathbf{V}_0^{-1}, \end{aligned}$$

we use these results and some algebra to show that $\mathbf{B}_0^\top \mathbf{V}_0^{-1} \mathbf{B}_0 + \hat{\mathbf{B}}^\top \mathbf{X}^\top \mathbf{X} \hat{\mathbf{B}} - \mathbf{B}_n^\top \mathbf{V}_n^{-1} \mathbf{B}_n = (\hat{\mathbf{B}} - \mathbf{B}_0)^\top \mathbf{V}_n (\hat{\mathbf{B}} - \mathbf{B}_0)$ taking into account that $\mathbf{V}_n = (\mathbf{V}_0^{-1} + \mathbf{X}^\top \mathbf{X})^{-1}$ and $\hat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$.

8. Show that in the linear regression model $\beta_n^\top (\mathbf{B}_n^{-1} - \mathbf{B}_n^{-1} \mathbf{M}^{-1} \mathbf{B}_n^{-1}) \beta_n = \beta_{**}^\top \mathbf{C} \beta_{**}$ and $\beta_{**} = \mathbf{X}_0 \beta_n$.

Answer

Taking into account that $(\mathbf{A} + \mathbf{B})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \mathbf{A}^{-1}$ [5], then we observe that $(\mathbf{B}_n^{-1} - \mathbf{B}_n^{-1} \mathbf{M}^{-1} \mathbf{B}_n^{-1}) = (\mathbf{B}_n + (\mathbf{X}_0^\top \mathbf{X}_0)^{-1})^{-1}$, where $(\mathbf{B}_n + (\mathbf{X}_0^\top \mathbf{X}_0)^{-1})^{-1} = \mathbf{X}_0^\top \mathbf{X}_0 - \mathbf{X}_0^\top \mathbf{X}_0 (\mathbf{B}_n^{-1} + \mathbf{X}_0^\top \mathbf{X}_0)^{-1} \mathbf{X}_0^\top \mathbf{X}_0 = \mathbf{X}_0^\top \mathbf{X}_0 - \mathbf{X}_0^\top \mathbf{X}_0 \mathbf{M}^{-1} \mathbf{X}_0^\top \mathbf{X}_0$, thus

$$\begin{aligned} \beta_n^\top (\mathbf{B}_n^{-1} - \mathbf{B}_n^{-1} \mathbf{M}^{-1} \mathbf{B}_n^{-1}) \beta_n &= \beta_n^\top (\mathbf{X}_0^\top \mathbf{X}_0 - \mathbf{X}_0^\top \mathbf{X}_0 \mathbf{M}^{-1} \mathbf{X}_0^\top \mathbf{X}_0) \beta_n \\ &= \beta_n^\top \mathbf{X}_0^\top (\mathbf{I}_{N_0} - \mathbf{X}_0 \mathbf{M}^{-1} \mathbf{X}_0^\top) \mathbf{X}_0 \beta_n \\ &= \beta_n^\top \mathbf{X}_0^\top \mathbf{C} \mathbf{X}_0 \beta_n \\ &= \beta_{**}^\top \mathbf{C} \beta_{**}. \end{aligned}$$

Let's show that $\beta_{**} = \mathbf{X}_0 \beta_n$,

$$\begin{aligned}
\beta_{**} &= \mathbf{C}^{-1} \mathbf{X}_0 \mathbf{M}^{-1} \mathbf{B}_n^{-1} \beta_n \\
&= (\mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{B}_n \mathbf{X}_0^\top) \mathbf{X}_0 \mathbf{M}^{-1} \mathbf{B}_n^{-1} \beta_n \\
&= (\mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{B}_n \mathbf{X}_0^\top) \mathbf{X}_0 (\mathbf{B}_n - \mathbf{B}_n ((\mathbf{X}_0^\top \mathbf{X}_0)^{-1} + \mathbf{B}_n)^{-1} \mathbf{B}_n) \mathbf{B}_n^{-1} \beta_n \\
&= (\mathbf{I}_{N_0} + \mathbf{X}_0 \mathbf{B}_n \mathbf{X}_0^\top) (\mathbf{X}_0 \beta_n - \mathbf{X}_0 \mathbf{B}_n ((\mathbf{X}_0^\top \mathbf{X}_0)^{-1} + \mathbf{B}_n)^{-1} \beta_n) \\
&= \mathbf{X}_0 \beta_n - \mathbf{X}_0 \mathbf{B}_n ((\mathbf{X}_0^\top \mathbf{X}_0)^{-1} + \mathbf{B}_n)^{-1} \beta_n + \mathbf{X}_0 \mathbf{B}_n \mathbf{X}_0^\top \mathbf{X}_0 \beta_n \\
&\quad - \mathbf{X}_0 \mathbf{B}_n \mathbf{X}_0^\top \mathbf{X}_0 \mathbf{B}_n ((\mathbf{X}_0^\top \mathbf{X}_0)^{-1} + \mathbf{B}_n)^{-1} \beta_n \\
&= \mathbf{X}_0 \beta_n - \mathbf{X}_0 \mathbf{B}_n [((\mathbf{X}_0^\top \mathbf{X}_0)^{-1} + \mathbf{B}_n)^{-1} - \mathbf{X}_0^\top \mathbf{X}_0 + \mathbf{X}_0^\top \mathbf{X}_0 \mathbf{B}_n ((\mathbf{X}_0^\top \mathbf{X}_0)^{-1} + \mathbf{B}_n)^{-1}] \beta_n.
\end{aligned}$$

Using that $(\mathbf{A} + \mathbf{B})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{A} + \mathbf{B})^{-1}$, we observe that the expression in brackets is equal to $\mathbf{0}$, then we have the result.

9. Show that $(\mathbf{Y} - \mathbf{X}\mathbf{B})^\top (\mathbf{Y} - \mathbf{X}\mathbf{B}) = \mathbf{S} + (\mathbf{B} - \hat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \hat{\mathbf{B}})$ where $\mathbf{S} = (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^\top (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$, $\hat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ in the multivariate regression model.

Answer

$$\begin{aligned}
(\mathbf{Y} - \mathbf{X}\mathbf{B})^\top (\mathbf{Y} - \mathbf{X}\mathbf{B}) &= (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}} + \mathbf{X}\hat{\mathbf{B}} - \mathbf{X}\mathbf{B})^\top (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}} + \mathbf{X}\hat{\mathbf{B}} - \mathbf{X}\mathbf{B}) \\
&= (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^\top (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}) + 2(\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^\top (\mathbf{X}\hat{\mathbf{B}} - \mathbf{X}\mathbf{B}) \\
&\quad + (\mathbf{X}\mathbf{B} - \mathbf{X}\hat{\mathbf{B}})^\top (\mathbf{X}\mathbf{B} - \mathbf{X}\hat{\mathbf{B}}) \\
&= \mathbf{S} + (\mathbf{B} - \hat{\mathbf{B}})^\top \mathbf{X}^\top \mathbf{X} (\mathbf{B} - \hat{\mathbf{B}}),
\end{aligned}$$

given that $(\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^\top (\mathbf{X}\hat{\mathbf{B}} - \mathbf{X}\mathbf{B}) = \hat{\mathbf{U}}^\top \mathbf{X} (\hat{\mathbf{B}} - \mathbf{B})$, using that $\hat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ which implies $\mathbf{X}^\top \mathbf{X} \hat{\mathbf{B}} = \mathbf{X}^\top \mathbf{Y} = \mathbf{X}^\top \mathbf{X} \hat{\mathbf{B}} + \mathbf{X}^\top \hat{\mathbf{U}}$, then $\mathbf{X}^\top \hat{\mathbf{U}} = \mathbf{0}$.

10. 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 supporters of Hillary vs Trump?
 - What is the probability that sampling one hundred independent individuals 44, 40 and 16 support Hillary, Trump and other, respectively?

Answer

R code. Multinomial-Dirichlet model: Polling 2016 USA presidential race

```

set.seed(010101)
# Multinomial-Dirichlet example:
# Polling 2016 USA presidential race
y <- c(411, 373, 149)
# Clinton, Trump, Other
# Public Policy Polling September 27-28,
# 2016 five-way race

alpha0 <- rep(1, 3)
# Hyperparameters: non-informative distribution
alphan <- alpha0 + y
S <- 100000
# Sample draws of posterior
thetas <- MCMCpack::rdirichlet(S, alphan)
colnames(thetas) <- c("Clinton", "Trump", "Other")
head(thetas)

```

	Clinton	Trump	Other
[1,]	0.4211346	0.4188607	0.1600046
[2,]	0.4244207	0.4224523	0.1531270
[3,]	0.4349268	0.3843953	0.1806779
[4,]	0.4533499	0.4005530	0.1460972
[5,]	0.4381799	0.3968502	0.1649699
[6,]	0.4436852	0.3971321	0.1591827

```

dif <- thetas[,1] - thetas[,2]
# Difference of shares Hillary vs Trump
data <- data.frame(dif)
names(data) <- c("Difference")
library(ggplot2)
p <- ggplot(data) +
  geom_histogram(aes(x = Difference), binwidth = 0.01) +
  geom_vline(xintercept=0.0, lwd=1, colour="red") +
  ggtitle("Percentage difference Clinton vs Trump  
2016 presidential race") +
  xlab("Percentage Difference") + ylab("")

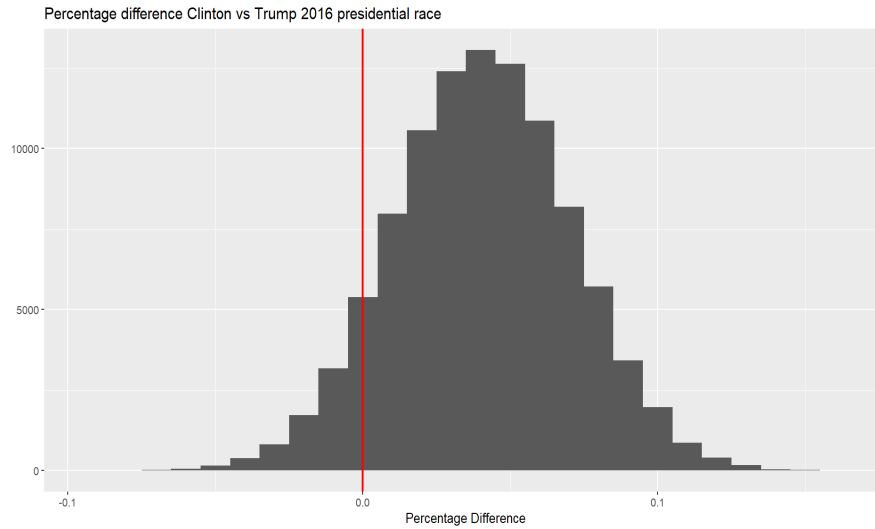
difmcmc <- coda::mcmc(dif)
# Declaring a MCMC object
summary(difmcmc)

Iterations = 1:1e+05
Thinning interval = 1
Number of chains = 1
Sample size per chain = 1e+05

1. Empirical mean and standard deviation for each
variable, plus standard error of the mean:

```

Mean	SD	Naive SE	Time-series SE
4.062e-02	2.996e-02	9.474e-05	
9.474e-05			

**FIGURE 3.1**

Percentage difference: Hillary Clinton vs Donald Trump, five-way race.

There is a 95% probability that the percentage difference between Hillary and Trump according to this poll is (-1.8%, 9.9%). The probability of Hillary having more supporters is 91.3%

***R code. Multinomial-Dirichlet model: Polling 2016
USA presidential race***

2. Quantiles for each variable:

```

2.5%      25%      50%      75%      97.5%
-0.01817  0.02033  0.04058  0.06089  0.09923
CW <- mean(difmcmc>0)
CW
0.91339

# Predictive distribution by simulation
y0 <- c(44, 40, 16)

Pred <- apply(thetas, 1, function(p) {
  rmultinom(1, size = sum(y0), prob = p)})
sum(sapply(1:S, function(s) {
  sum(Pred[,s] == y0) == 3}))/S
0.00825

# Predictive distribution by analytical expression
PredY0 <- function(y0){
  n <- sum(y0)
  Res1 <- sum(sapply(1:length(y), function(l){
    lgamma(alphan[l]+y0[l]) -
    lgamma(alphan[l]) - lfactorial(y0[l]))))
  Res <- lfactorial(n)+lgamma(sum(alphan))
  -lgamma(sum(alphan)+n) + Res1
  return(exp(Res))
}
PredY0(y0)
0.00850

```

The probability that from one hundred random selected people 44 support Hillary, 40 support Trump and 16 support other candidate is 0.85%.

11. 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 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 95% confidence and credible intervals for μ .
- What is the posterior probability that $\mu > 103$?

Answer

R code. Math test example continues

```

set.seed(010101)
N <- 50
# Sample size
muhat <- 102
# Sample mean
sig2hat <- 10
# Sample variance

# Hyperparameters
mu0 <- 100
beta0 <- 1
delta0 <- 0.001
alpha0 <- 0.001

S <- 100000
# Posterior draws
alphan <- alpha0 + N
deltan <- sig2hat*(N - 1) + delta0 +
  beta0*N/(beta0 + N)*(muhat - mu0)^2
sig2Post <- invgamma::rinvgamma(S, shape = alphan,
  rate = deltan)
summary(sig2Post)
betan <- beta0 + N
mun <- (beta0*mu0 + N*muhat)/betan
muPost <- sapply(sig2Post, function(s2){rnorm(1, mun,
  sd = (s2/betan)^0.5)})

muPostq <- quantile(muPost, c(0.025, 0.5, 0.975))
muPostq
      2.5%      50%      97.5%
101.0929 101.9625 102.8311
cutoff <- 103
PmuPostcutoff <- mean(muPost > cutoff)
PmuPostcutoff
0.00994
# Using Student's t
muPost_t <- ((deltan/(alphan*betan))^0.5)*rt(S, alphan)
  + mun
c1 <- rgb(173,216,230,max = 255, alpha = 50,
  names = "lt.blue")
c2 <- rgb(255,192,203, max = 255, alpha = 50,
  names = "lt.pink")
hist(muPost, main = "Histogram:_Posterior_mean",
  xlab = "Posterior_mean", col = c2)
hist(muPost_t, main = "Histogram:_Posterior_mean",
  xlab = "Posterior_mean", add = T, col = c1)
muPost_tq <- quantile(muPost_t, c(0.025, 0.5, 0.975))
muPost_tq
      2.5%      50%      97.5%
101.0837 101.9608 102.8435
PmuPost_t cutoff <- mean(muPost_t > cutoff)
PmuPost_t cutoff
0.01087

```

We perform our calculations using the posterior conditional distribution, and the posterior marginal distribution. Both procedures give similar results as we can observe from Figure 3.2.

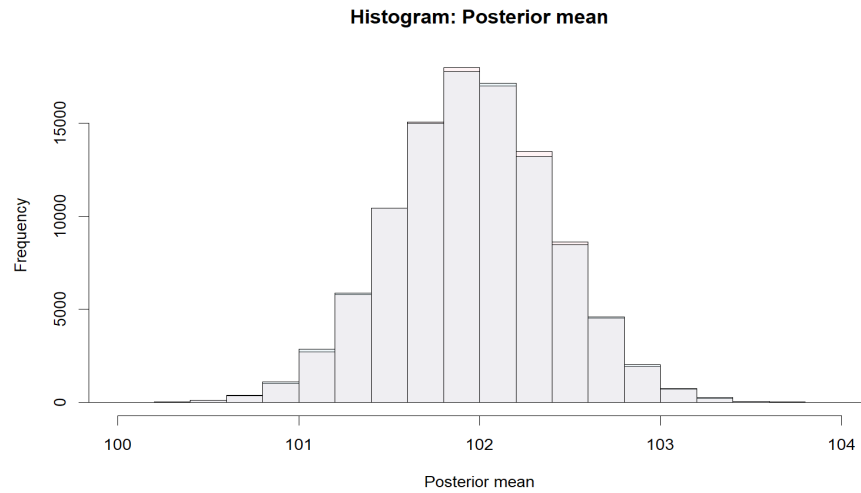


FIGURE 3.2

Histogram using the posterior conditional distribution and the posterior marginal distribution

We have that the 95% credible interval is $(101.08, 102.84)$, and the probability of having a value greater than 103 is 1.09%.



Bibliography

- [1] H. Jeffreys. *Theory of Probability*. Oxford University Press, London, 1961.
- [2] Valen E Johnson and David Rossell. On the use of non-local prior densities in bayesian hypothesis tests. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 72(2):143–170, 2010.
- [3] Dennis V Lindley. A statistical paradox. *Biometrika*, 44(1/2):187–192, 1957.
- [4] Kevin P Murphy. Conjugate bayesian analysis of the gaussian distribution. *def*, 1(2 σ 2):16, 2007.
- [5] A. F. M. Smith. A General Bayesian Linear Model. *Journal of the Royal Statistical Society. Series B (Methodological)*., 35(1):67–75, 1973.
- [6] A. Zellner. *Introduction to Bayesian inference in econometrics*. John Wiley & Sons Inc., 1996.