

MS2505: Bayesian Statistics

Course Project

December 12, 2024

Name	Samuel Jonsson
E-Mail	sajs19@student.bth.se
Person Nr.	19990415-5596
Program	DVAMI19h



1 Setup

- Describe the data and the analysis problem.
- Choose and describe the modeling approach (e.g., non-hierarchical or hierarchical model).
- Justify your prior choice.
- Perform posterior predictive checks.

1.1 Analysis problem

1.2 Data Selection

Describe the data and the analysis problem.

The dataset selected is a datasets containing a list of emails, as well as a label marking each email as "spam" or "ham" (spam or not spam). The first 10 rows of the dataset looks as follows:

Table 1

mail_data.csv dataset first 10 rows

Category	Message
ham	"Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat..."
ham	Ok lar... Joking wif u oni...
spam	Free entry in 2 a wkly comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entry question(std txt rate)T&C's apply 08452810075over18's
ham	U dun say so early hor... U c already then say...
ham	"Nah I don't think he goes to usf, he lives around here though"
spam	"FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like some fun you up for it still? Tb ok! XxX std chgs to send, £1.50 to rcv"
ham	Even my brother is not like to speak with me. They treat me like aids patent.
ham	As per your request 'Melle Melle (Oru Minnaminunginte Nurungu Vettam)' has been set as your callertune for all Callers. Press *9 to copy your friends Callertune
...	...

Then, using a Python script, the labels were converted to 1 if it was "spam" and 0 if it was "ham", for easier analysis.

1.3 Model

- Choose and describe the modeling approach (e.g., non-hierarchical or hierarchical model).
- Justify your prior choice.

The model chosen was a binomial likelihood model with a beta prior. As the goal is to analyse the probability of an email being spam, the fallout will be binary (either it is spam or it is not). Hence, a binomial likelihood, where I want to find the parameter θ in a dataset of fixed size with a set number of "successes" and "fails" (spam and ham), is appropriate.

Additionally, as I do not have any prior knowledge in regards to this distribution, a non-informative prior is the most suited option, and as $Beta(1, 1)$ is a common prior used with binomial likelihood functions, I chose it for this problem.

1.4 Prior checks

Perform posterior predictive checks.

2 Results

Include diagnostics to assess model convergence and adequacy.

3 Discussion

Discuss results, problems encountered, and possible improvements.

A R Code

Listing 1

Project R code

```
1 # Load required libraries
2 library(bayesplot)
3 library(rstanarm)
4 library(ggplot2)
5 library(brms)
6
7 # Set a global random seed for reproducibility
8 set.seed(123)
9
10 # Create the directory if it doesn't exist
11 if (!dir.exists("Project/logs")) {
12   dir.create("Project/logs", recursive = TRUE)
13 }
14
15 # Create the directory for figures if it doesn't exist
16 if (!dir.exists("figures")) {
17   dir.create("figures", recursive = TRUE)
18 }
19
20 # Specify the log file path
21 log_file <- "Project/logs/R_output.log"
22
23 # Open the sink to redirect output
24 sink(log_file)
25
```

```
26 # Read the data
27 mail_data <- read.csv("Project/data/mail_data_bin.csv")
28
29 # Set metadata
30 alpha_prior <- 1
31 beta_prior <- 1
32 total_emails <- nrow(mail_data)
33 spam_count <- sum(mail_data$Category == 1)
34 ham_count <- sum(mail_data$Category == 0)
35
36 # Fit the model with a vague prior
37 fit_prior <- brm(
38   Category ~ 1,
39   data = mail_data,
40   family = bernoulli(),
41   prior = prior(beta(1, 1),
42     class = "Intercept"
43   )
44 )
45
46 # Prior predictive check
47 pdf("figures/prior_predictive_check.pdf")
48 pp_check(fit_prior, type = "hist") +
49   ggtitle("Prior Predictive Check for Email Spam Model")
50 dev.off()
51
52 # Fit a robust model with a more informative prior
53 fit_robust <- brm(
54   Category ~ 1,
55   data = mail_data,
56   family = bernoulli(),
57   prior = prior(beta(2, 2), class = "Intercept")
58 )
59
60 # Posterior predictive check
61 pdf("figures/posterior_predictive_check.pdf")
62 pp_check(fit_robust, type = "dens_overlay") +
63   ggtitle("Posterior Predictive Check for Robust Email Spam Model")
64 dev.off()
65
66 # Compute posterior parameters for P(spam)
67 alpha_post <- alpha_prior + spam_count
68 beta_post <- beta_prior + ham_count
69
70 # Posterior probability of an email being spam
71 posterior_spam <- alpha_post / (alpha_post + beta_post)
72
73 # Print results
74 cat("Prior: Beta(", alpha_prior, ",", beta_prior, ")\n")
75 cat("Spam Count:", spam_count, "\n")
76 cat("Ham Count:", ham_count, "\n")
```

```
77 cat("Posterior: Beta(", alpha_post, ",", beta_post, ")\n")
78 cat("P(spam):", posterior_spam, "\n")
79
80 # Posterior Predictive Checking
81 cat("\n--- Posterior Predictive Checking ---\n")
82
83 # Simulate posterior predictive samples
84 num_samples <- 1000
85 posterior_samples <- rbeta(num_samples, alpha_post, beta_post)
86
87 observed_counts <- spam_count / total_emails
88
89 # Generate density overlay
90 y <- rep(observed_counts, num_samples)
91 yrep <- matrix(posterior_samples, nrow = num_samples)
92
93 # Save density overlay plot
94 # pdf("../figures/ppc_density_overlay.pdf")
95 # ppc_dens_overlay(y = y, yrep = yrep) +
96 #   ggtitle("Posterior Predictive Check: Density Overlay")
97 # dev.off()
98
99 # Generate and save histogram of posterior samples
100 pdf("figures/ppc_histogram.pdf")
101 hist(posterior_samples,
102     breaks = 30, col = "blue", border = "white",
103     main = "Posterior Predictive Distribution", xlab = "P(spam)"
104 )
105 dev.off()
106
107 # Sensitivity Analysis
108 cat("\n--- Sensitivity Analysis ---\n")
109 sensitivity_results <- data.frame()
110 alpha_values <- seq(0.5, 2, by = 0.5)
111 beta_values <- seq(0.5, 2, by = 0.5)
112
113 for (alpha in alpha_values) {
114   for (beta in beta_values) {
115     alpha_post_temp <- alpha + spam_count
116     beta_post_temp <- beta + ham_count
117     posterior_temp <- alpha_post_temp / (alpha_post_temp +
118         beta_post_temp)
119     sensitivity_results <- rbind(
120         sensitivity_results,
121         data.frame(alpha, beta, posterior_temp)
122     )
123   }
124 }
125
126 # Display the sensitivity analysis results
127 cat("Sensitivity Analysis Results:\n")
```

```
127 print(sensitivity_results)
128
129 # Save sensitivity results to CSV
130 write.csv(sensitivity_results,
131           "Project/logs/sensitivity_analysis.csv",
132           row.names = FALSE
133 )
134
135 # Flush the output and close the sink
136 flush.console()
137 sink()
```