

MS2505: Bayesian Statistics

Course Project

December 26, 2024

Name

Samuel Jonsson

Email

sajs19@student.bth.se

Person Nr.

19990415-5596

Program

DVAMI19h



1 Analysis problem

The analysis aims to estimate the proportion of emails classified as spam in a given dataset. Using Bayesian methods, the focus is on deriving a reliable estimate of this proportion while accounting for uncertainty. The simplicity of the approach allows for clear insights into the spam classification problem, although it abstracts away complexities like email content or contextual nuances. The objective is to provide a robust probabilistic framework for understanding the distribution of spam within the dataset.

2 Data Description

The dataset selected for this project was obtained from Kaggle and consists of labeled email messages, where each email is classified as either “spam” or “ham”. The dataset contains the following columns:

- **Category:** A binary label indicating whether the email is spam (1) or ham (0).
- **Message:** The text content of the email.

The first 10 rows of the dataset are shown in Table 1. For this analysis, the focus is on the Category column, which serves as the response variable for the Bayesian classification model.

Table 1

<code>mail_data.csv</code>	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 ho. . . 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
...	...

Using a python script, the labels are then converted to numerical values, where spam is represented as 1 and ham as 0. After preprocessing the data, the dataset contains a total of

747	Spam
4825	Ham
5572	Total Emails

These counts are used to update the beta distribution's prior parameters and compute the posterior distribution.

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

3 Model

The model chosen for this analysis is a binomial likelihood model with a beta prior. This approach is particularly suitable for binary classification problems where the response variable is dichotomous, such as classifying emails as either “spam” (1) or “ham” (0).

The binomial likelihood model represents the probability of observing a given number of successes (spam) in a fixed number of trials (emails). The parameter of interest, θ , represents the probability of an email being classified as spam. Using this model allows us to estimate θ while accounting for uncertainty.

To complement the likelihood, we use a $Beta(1, 1)$ prior, which is non-informative and reflects a state of prior ignorance about θ . This choice aligns with common practices in Bayesian analysis, particularly when no strong prior knowledge exists. By combining the prior with the observed data, we compute the posterior distribution of θ , providing a probabilistic estimate of the spam probability.

In addition, Markov Chain Monte Carlo (MCMC) methods are employed to validate the results obtained from the closed-form posterior distribution. MCMC sampling allows us to approximate the posterior distribution of θ when the analytic computation becomes infeasible.

4 Results

This section presents the outcomes of the Bayesian analysis performed using a Beta-binomial model to estimate the probability (θ) that an email is classified as spam. Results from Monte Carlo and MCMC sampling methods are summarized and validated, followed by visualizations to illustrate the posterior distributions and sampling behavior.

4.1 Posterior Analysis

Monte Carlo sampling, performed with 10,000 draws from the posterior distribution of θ , provides a detailed summary of the key statistics in Table 2. The posterior mean estimate for θ is approximately 13.42%, with a narrow 95% credible interval of [0.125, 0.143]. This suggests a high degree of confidence in the estimate and indicates that, on average, 13.42% of emails are classified as spam.

Table 2

Monte Carlo Posterior Analysis Data

Posterior Mean	0.1341994
Posterior SD	0.004576791
95% Credible Interval	[0.1252803, 0.1432051]

4.2 MCMC Validation

To validate the results, MCMC sampling was performed using the `MCMCmetrop1R` function. As shown in Table 3, the posterior mean is consistent at 13.44%, with a 95% credible interval of $[0.125, 0.144]$. These results further confirm the robustness of the estimates and align closely with the Monte Carlo findings.

Table 3

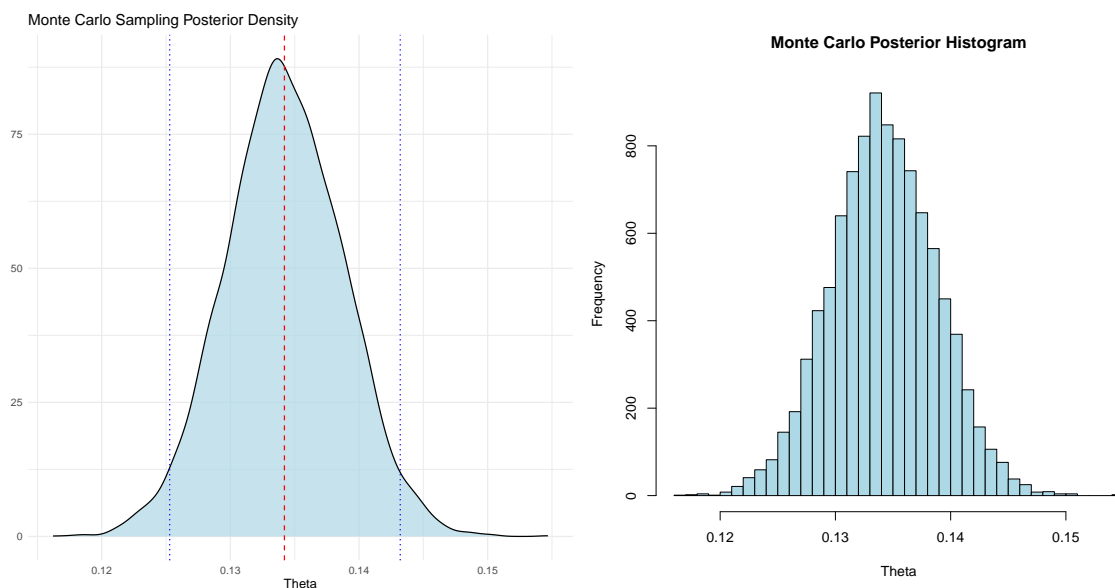
MCMC Posterior Analysis Data

Posterior Mean	0.1344085
Posterior SD	0.004670508
95% Credible Interval	$[0.1254781, 0.1437142]$

The strong agreement between the Monte Carlo and MCMC methods demonstrates the model's validity and the consistency of the Bayesian inference process.

4.3 Visualization

The posterior distributions derived from Monte Carlo and MCMC sampling are visualized in Figures 1 and 2, respectively. Both distributions exhibit a sharp peak around $\theta \approx 0.135$, reflecting the central tendency of the estimates.



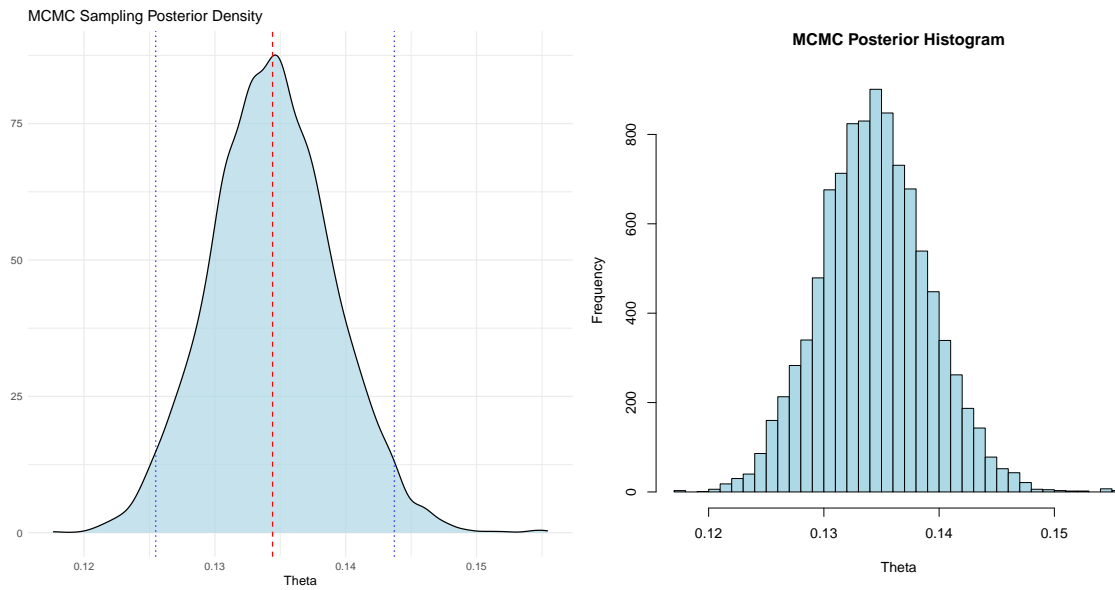
(a) Density plot

(b) Histogram

Figure 1

Posterior density and histogram for Monte Carlo sampling. The sharp peak at $\theta \approx 0.135$ aligns with Table 2.

Similarly, the MCMC posterior density and histogram (Figure 2) reveal consistent results, reinforcing the robustness of the Bayesian model.



(a) MCMC density plot

(b) MCMC histogram

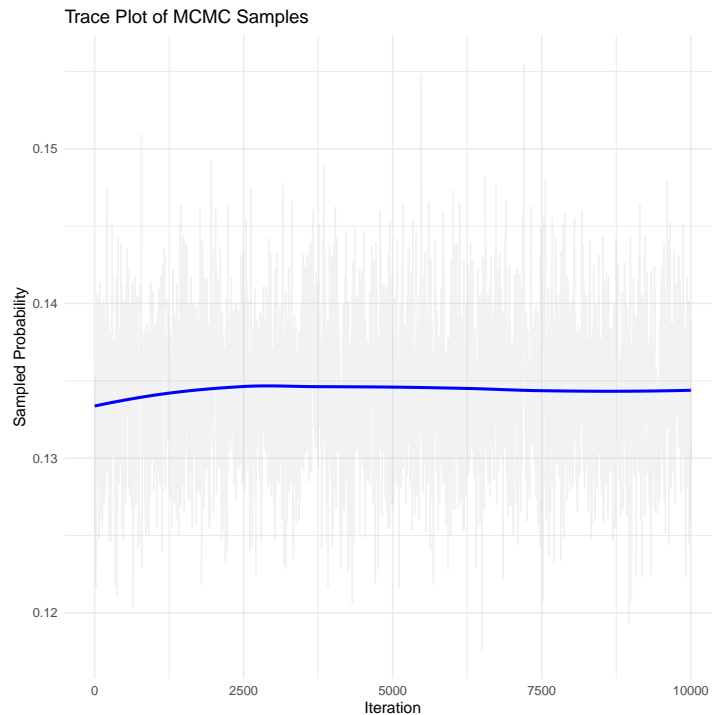
Figure 2

Posterior density and histogram for MCMC sampling, demonstrating consistency with Monte Carlo results and Table 3.

The tails of the posterior distributions provide additional insights into the uncertainty associated with θ . Both Monte Carlo and MCMC density plots show symmetric tails that taper quickly beyond the 95% credible intervals, indicating a low likelihood of extreme values. Sparse population in the outer bins of the histograms further highlights the concentration of posterior mass near the central estimates. Slightly broader tails in the MCMC density plot compared to the Monte Carlo results may suggest minor differences in sampling variability, though these are negligible in the overall context of the model's consistency and reliability.

4.4 MCMC Trace Plot

The trace plot of the MCMC samples (Figure 3) provides further evidence of the sampling process's stability and convergence. The sampled values oscillate consistently around a central range without noticeable trends or significant deviations, confirming that the chain has reached a stationary distribution.

**Figure 3**

MCMC sampling trace plot. The stable fluctuations confirm convergence to the posterior distribution.

This stable behavior, along with minimal autocorrelation and effective mixing, suggests that the MCMC samples are representative of the posterior distribution. While minor clustering is observed in certain regions of the plot, this does not significantly affect the validity of the estimates.

5 Discussion

The Bayesian analysis conducted in this project produced reliable and consistent results for estimating the probability θ that an email is classified as spam. As mentioned in Section 3, the simplicity of the Beta-binomial model provided a clear and computationally efficient approach for this foundational analysis. Both Monte Carlo and MCMC sampling methods yielded posterior mean estimates of approximately 13.4%, with narrow 95% credible intervals, demonstrating a high level of confidence in the central estimate. The visualizations, discussed in Section 4.3, further supported these findings, revealing well-behaved posterior distributions with symmetric tails and rapid tapering, reflecting low uncertainty around the estimates. Additionally, the trace plot in Section 4.4 confirmed convergence of the MCMC sampling process, with sampled values stabilizing without significant deviations.

However, while the Beta-binomial model provides a straightforward and computationally efficient framework, it has inherent limitations. As mentioned in Section 3, the model focuses on binary classification, treating emails as either spam (1) or ham (0). This simplicity, while beneficial for interpretability, abstracts away complexities such as email content or metadata, which could provide richer insights into spam classification. Consequently, this approach restricts the model's ability to generalize to unseen datasets with varying characteristics. Additionally, the choice of a $Beta(1, 1)$ prior assumes no prior knowledge about the data, which, while this neutral stance aligns with the exploratory nature of the analysis, it may not optimally reflect datasets with unusual or heavily skewed properties.

To address these challenges, several improvements could be explored. First, as noted in Section 2, the dataset contains the email text content, and incorporating this into the model

could enhance its predictive capabilities, allowing it to better generalize to diverse or unknown datasets. For instance, incorporating natural language processing techniques to analyze the textual content could provide additional predictors beyond the binary classification framework. Additionally, extending the model to include hierarchical structures, as a refinement of the foundational framework described in Section 3, could capture variations across subpopulations or different email sources, offering a more nuanced understanding of θ .

The choice of prior distribution, discussed in Section 3, is another area for potential improvement. Rather than using a neutral $Beta(1, 1)$ prior, future analyses could explore data-informed priors that align with specific dataset characteristics. For instance, if prior information about spam prevalence is available from other sources, it could be incorporated to improve the robustness and sensitivity of the posterior analysis. The nature of the problem, however, would make a generally representative prior quite challenging, as the number of spam emails in an unknown dataset would be difficult to determine without more contextual knowledge.

On the computational side, the MCMC sampling methods, while effective, could benefit from the adoption of advanced techniques such as Hamiltonian Monte Carlo or adaptive MCMC methods. As noted in Section 4.2, MCMC sampling produced slightly broader tails compared to Monte Carlo sampling, suggesting some variability in the results. Enhanced sampling techniques could reduce this variability while improving computational efficiency.

In summary, while the simplicity of the beta-binomial model proved effective for this initial analysis, it highlights a trade-off between clarity and flexibility. The straightforward approach ensured interpretable results and computational feasibility. However, the challenges discussed here underscore the need for a more comprehensive framework in future analyses. Building on the current findings, incorporating richer features, hierarchical structures, and optimized priors would enhance the model's applicability to real-world spam detection scenarios. This iterative refinement process is essential for balancing simplicity with the complexities of dynamic and diverse datasets.

6 Appendix

A R Code

Listing 1

Project R code

```
1 # =====
2 # Bayesian Analysis for Email Spam Classification
3 # =====
4
5 # Load Required Libraries
6 required_packages <- c("ggplot2", "dplyr", "MCMCpack", "coda", "
  tidyrr")
7 installed_packages <- rownames(installed.packages())
8 for (pkg in required_packages) {
9   if (!pkg %in% installed_packages) install.packages(pkg)
10  library(pkg, character.only = TRUE)
11 }
12
13 # Set a global random seed for reproducibility
14 set.seed(123)
15
16 # Output Logs
17 if (!dir.exists("logs")) dir.create("logs")
18 results_log_file <- "logs/combined_results.log"
19
20 # =====
21 # Load Data
22 # =====
23
24 data_path <- "data/mail_data_bin.csv"
25 if (!file.exists(data_path)) stop("Data file not found!")
26 mail_data <- read.csv(data_path)
27
28 # Ensure the response variable is binary and properly coded
29 if (!all(mail_data$Category %in% c(0, 1))) {
30   stop("The response variable 'Category' must be binary (0 or 1).")
31 }
32
33 # Metadata
34 spam_count <- sum(mail_data$Category == 1)
35 ham_count <- sum(mail_data$Category == 0)
36 total_emails <- nrow(mail_data)
37
38
39 sink(results_log_file)
40 cat("--- Data Metadata ---\n")
41 cat("Spam Count:", spam_count, "\n")
42 cat("Ham Count:", ham_count, "\n")
43 cat("Total Emails:", total_emails, "\n")
```



```
44 sink()
45
46 # =====
47 # Beta Posterior Analysis
48 # =====
49
50 # Prior parameters
51 prior_alpha <- 1
52 prior_beta <- 1
53
54 # Posterior parameters
55 posterior_alpha <- prior_alpha + spam_count
56 posterior_beta <- prior_beta + ham_count
57
58 # Monte Carlo Sampling
59 n_samples <- 10000
60 beta_samples <- rbeta(n_samples, posterior_alpha, posterior_beta)
61
62 # Summary statistics
63 beta_mean <- mean(beta_samples)
64 beta_sd <- sd(beta_samples)
65 beta_ci <- quantile(beta_samples, c(0.025, 0.975))
66
67 sink(results_log_file, append = TRUE)
68 cat("\n--- Beta Posterior Analysis ---\n")
69 cat("Posterior Mean:", beta_mean, "\n")
70 cat("Posterior SD:", beta_sd, "\n")
71 cat("95% Credible Interval:", beta_ci, "\n")
72 sink()
73
74 # =====
75 # MCMC Sampling
76 # =====
77
78 # Define log-posterior function
79 log_posterior <- function(params) {
80   theta <- params[1]
81   log_prior <- dbeta(theta, 1, 1, log = TRUE)
82   log_likelihood <- sum(dbinom(mail_data$Category,
83                               size = 1,
84                               prob = theta,
85                               log = TRUE))
86   return(log_prior + log_likelihood)
87 }
88
89 # Run MCMC sampling
90 mcmc_results <- MCMCmetrop1R(
91   fun = log_posterior,
92   theta.init = 0.5,
93   burnin = 1000,
94   mcmc = n_samples,
```

```
95   thin = 1,
96   verbose = 0
97 )
98
99 # Extract posterior samples
100 mcmc_samples <- as.vector(mcmc_results)
101 mcmc_mean <- mean(mcmc_samples)
102 mcmc_sd <- sd(mcmc_samples)
103 mcmc_ci <- quantile(mcmc_samples, c(0.025, 0.975))
104
105 sink(results_log_file, append = TRUE)
106 cat("\n--- MCMC Sampling Analysis ---\n")
107 cat("Posterior Mean:", mcmc_mean, "\n")
108 cat("Posterior SD:", mcmc_sd, "\n")
109 cat("95% Credible Interval:", mcmc_ci, "\n")
110 sink()
111
112 # =====
113 # Diagnostics and Visualization
114 # =====
115
116 if (!dir.exists("figures")) dir.create("figures")
117
118 # 1. Beta Density Plot
119 pdf("figures/beta_posterior_density_plot.pdf")
120 ggplot(data = data.frame(samples = beta_samples), aes(x = samples))
121   +
122   geom_density(fill = "lightblue", alpha = 0.7) +
123   geom_vline(xintercept = beta_mean, color = "red", linetype = "
     dashed") +
124   geom_vline(xintercept = beta_ci, color = "blue", linetype = "
     dotted") +
125   labs(title = "Monte Carlo Sampling Posterior Density", x = "Theta
     ", y = "") +
126   theme_minimal()
127 dev.off()
128
129 # 2. MCMC Density Plot
130 pdf("figures/mcmc_posterior_density_plot.pdf")
131 ggplot(data = data.frame(samples = mcmc_samples), aes(x = samples))
132   +
133   geom_density(fill = "lightblue", alpha = 0.7) +
134   geom_vline(xintercept = mcmc_mean, color = "red", linetype = "
     dashed") +
135   geom_vline(xintercept = mcmc_ci, color = "blue", linetype = "
     dotted") +
136   labs(title = "MCMC Sampling Posterior Density", x = "Theta", y =
     "") +
137   theme_minimal()
138 dev.off()
```

```
138 # 3. Beta Posterior Histogram
139 pdf("figures/beta_posterior_histogram.pdf")
140 hist(beta_samples,
141       breaks = 30, col = "lightblue", border = "black",
142       xlab = "Theta", main = "Monte Carlo Posterior Histogram"
143 )
144 dev.off()
145
146 # 4. MCMC Posterior Histogram
147 pdf("figures/mcmc_posterior_histogram.pdf")
148 hist(mcmc_samples,
149       breaks = 30, col = "lightblue", border = "black",
150       xlab = "Theta", main = "MCMC Posterior Histogram"
151 )
152 dev.off()
153
154 # 5. Trace Plot for MCMC
155 pdf("figures/mcmc_trace_plot.pdf")
156 ggplot(
157   data.frame(Iteration = 1:n_samples, Sample = mcmc_samples),
158   aes(x = Iteration, y = Sample)
159 ) +
160   geom_line(alpha = 0.2, color = "gray") +
161   geom_smooth(color = "blue", method = "loess", se = FALSE) +
162   labs(title = "Trace Plot of MCMC Samples",
163        x = "Iteration",
164        y = "Sampled Probability") +
165   theme_minimal()
166 dev.off()
167
168 # =====
169 # Completion
170 # =====
171
172 cat("Analysis complete. Check 'logs' and 'figures' directories for
    results.\n")
```