

Section 8: Bayesian Inference

Ricky Truong (rickytruong@college.harvard.edu),
Emily Xing (exing@college.harvard.edu)

1 Introduction

1.1 Logistics

Welcome! The goal of our weekly section is to review content from lecture, with emphasis on big-picture intuition, mathematical proof, and hands-on practice.

Specifically, we aim for our sections to be interactive, well-paced, inclusive, and fun! Please do not hesitate to reach out if you have any questions/feedback!

1.2 Office Hours

- Mondays, 7:30–9:30 PM in Adams D-Hall (Ricky).
- Fridays, 11 AM–12 PM in Maxwell-Dworkin 2nd Floor (Emily).
- Saturdays, 10:30–11:30 AM in Cabot D-Hall (Emily).

2 Big Picture

This week we turn to *Bayesian inference*, a fundamentally different philosophy from the frequentist approach we have been using. Most of Harvard Stats faculty (like Joe) are Bayesian; I think you will soon find out why!

In the *frequentist* framework, θ is an unknown, but fixed. We build estimators and confidence intervals whose performance we assess over hypothetical repeated experiments. However, in the *Bayesian* framework, θ is treated as a random variable, and we use probability to quantify our uncertainty about it.

The workflow is:

- *Prior* $\pi(\theta)$: our beliefs about θ before seeing data, encoded as a probability distribution. This is the most controversial part about Bayesian statistics.
- *Likelihood* $L(\theta; y) = f(y | \theta)$: how likely the observed data are under each value of θ . Both Bayesians and frequentists agree on this!
- *Posterior* $\pi(\theta | y)$: our updated beliefs about θ after seeing data. Get familiar with the term, this is the core object of Bayesian inference.

We will also discuss *conjugate priors* (which keep the posterior in a tractable family), *credible intervals* (the Bayesian analog of confidence intervals, with a more direct probability interpretation), and *Bayesian point estimation* (posterior mean, median, and MAP, each optimal under a different loss). We will then discuss *hierarchical models* that give new insight into frequentist problems like the inadmissibility of the MLE (Stein's paradox).

3 Prior to Posterior

Definition 1 (Prior, posterior, marginal likelihood). Consider a parametric model $f(y | \theta)$ with unknown parameter θ . In the Bayesian approach we posit a **joint distribution** for (Y, θ) . The **prior** $\pi(\theta)$ is the marginal distribution of θ , encoding our beliefs before seeing data. The **posterior** $\pi(\theta | y)$ is the conditional distribution of θ given the observed data y . The **marginal likelihood** (or *prior predictive distribution*) is $f(y) = \int f(y | \theta)\pi(\theta) d\theta$.

Theorem 1 (Bayes' rule).

$$\pi(\theta | y) = \frac{L(\theta; y) \pi(\theta)}{f(y)} \propto L(\theta; y) \pi(\theta).$$

We almost always work up to proportionality (normalization is difficult), since $f(y)$ does not depend on θ .

- **Strategy:** Write $\pi(\theta | y) \propto L(\theta; y)\pi(\theta)$, simplify, and pattern-match to a known distribution to identify the posterior without computing $f(y)$.
- **☒ Cromwell's rule.** If $\pi(\theta_0) = 0$ for some θ_0 , then $\pi(\theta_0 | y) = 0$ no matter what the data say. Never assign prior probability of exactly 0 or 1 to something unless it is logically impossible or certain!
- **As $n \rightarrow \infty$:** intuitively, the likelihood dominates the prior (the data overwhelm prior beliefs), and the posterior will concentrate near the MLE. For small n , the prior plays a large role based on weight!

Concept Checker 1. Let $Y | \theta \sim \text{Bern}(\theta)$ and $\theta \sim \text{Beta}(2, 2)$.

1. Write out $\pi(\theta | y) \propto L(\theta; y)\pi(\theta)$ and identify the posterior distribution.
2. What does the $\text{Beta}(2, 2)$ prior encode about our beliefs for θ ?
3. What happens to the posterior as we observe more and more data?

Solution

1. $L(\theta; y) = \theta^y(1 - \theta)^{1-y}$ and $\pi(\theta) \propto \theta^{2-1}(1 - \theta)^{2-1} = \theta(1 - \theta)$. Thus

$$\pi(\theta | y) \propto \theta^y(1 - \theta)^{1-y} \cdot \theta(1 - \theta) = \theta^{y+1}(1 - \theta)^{2-y},$$

which is the $\text{Beta}(y+2, 3-y)$ kernel. Since $y \in \{0, 1\}$, the posterior is either $\text{Beta}(2, 3)$ (if $y = 0$) or $\text{Beta}(3, 2)$ (if $y = 1$).

2. $\text{Beta}(2, 2)$ is symmetric about $1/2$ with mean $1/2$. It is like saying we have seen 1 prior success and 1 prior failure, so we mildly believe the coin is fair but are quite uncertain. The Beta is a good choice of prior in these circumstances of symmetry and probability (restricted support).

3. With n i.i.d. $\text{Bern}(\theta)$ observations and $S = \sum y_i$: posterior is $\text{Beta}(2+S, 2+n-S)$. As $n \rightarrow \infty$, the posterior concentrates near the MLE $\hat{\theta} = S/n = \bar{y}$, regardless of the prior. This is a conjugacy we will explore later.

4 Point Estimation

Definition 2 (Posterior mean, median, and mode). Let θ have a continuous posterior density $\pi(\theta | y)$. Then:

$$\text{Posterior mean} = \mathbb{E}[\theta | y] = \int \theta \pi(\theta | y) d\theta$$

$$\text{Posterior median} = Q_{\theta|y}(0.5)$$

$$\text{Posterior mode (MAP)} = \arg \max_{\theta} \pi(\theta | y) = \arg \max_{\theta} \{\log L(\theta; y) + \log \pi(\theta)\}$$

Theorem 2 (Optimal loss functions). • *Squared error loss $(\theta - \hat{\theta})^2$: minimized by the **posterior mean** $\mathbb{E}[\theta | y]$.*

• *Absolute error loss $|\theta - \hat{\theta}|$: minimized by the **posterior median** $Q_{\theta|y}(0.5)$.*

• *0–1 loss (in a limit): minimized by the **posterior mode (MAP)**.*

• **MAP vs. MLE.** MAP = MLE + log $\pi(\theta)$ in the optimization. With a flat (Uniform) prior on a bounded interval, MAP = MLE (why?). With an informative prior, MAP regularizes the MLE toward the prior mean.

Regularization is a good mix between frequentist and Bayesian approaches—if you are curious, it’s worth looking more into here:

• **MAP and LASSO.** With a Laplace prior $\pi(\theta) \propto e^{-d|\theta|}$, the MAP is the LASSO estimator—it thresholds \bar{y} exactly to 0 if $|\bar{y}| < c$. With a Normal prior, the MAP/posterior mean is ridge regression—it *shrinks* \bar{y} toward 0 but never exactly to 0.

• **⚠:** MAP is **not** invariant to reparameterization. If $\hat{\theta}_{\text{MAP}}$ maximizes $\pi(\theta | y)$, then $g(\hat{\theta}_{\text{MAP}})$ does *not* generally maximize $\pi(g(\theta) | y)$.

Concept Checker 2. For the Beta-Binomial model $Y | p \sim \text{Bin}(n, p)$, $p \sim \text{Beta}(a, b)$, write down the posterior mean, median, and MAP explicitly. For large n , what do all three converge to?

Solution

By Beta-Binomial conjugacy, $p | y \sim \text{Beta}(a + y, b + n - y)$, so $a' = a + y$, $b' = b + n - y$, $n' = a' + b' = a + b + n$.

$$\text{Posterior mean} = \frac{a'}{n'} = \frac{a + y}{a + b + n}, \quad \text{Posterior mode (MAP)} = \frac{a' - 1}{n' - 2} = \frac{a + y - 1}{a + b + n - 2}$$

The posterior median has no closed form but $\approx (a + y - 1/3)/(a + b + n - 2/3)$.

As $n \rightarrow \infty$: pseudo-counts a, b become negligible, and all three converge to $y/n = \hat{p}_{\text{MLE}}$.

5 Credible Intervals

Definition 3 (Credible interval). Let $0 < \alpha < 1$. A $100(1 - \alpha)\%$ **credible interval** (or **posterior probability interval**) for θ is an interval $[a(y), b(y)]$ such that

$$P(a(y) \leq \theta \leq b(y) \mid y) = 1 - \alpha.$$

The standard choice is the equal-tailed interval $[Q_{\theta|y}(\alpha/2), Q_{\theta|y}(1 - \alpha/2)]$.

- **Direct probability interpretation.** A 95% credible interval means: given the data, there is a 95% probability that θ lies in the interval. This is the statement we *want* to make about confidence intervals, but can't (frequentist intervals are about the procedure, not the parameter).

- **Average frequentist coverage.** A 95% credible interval also has *on average* 95% frequentist coverage, where the average is over both θ and Y . By Adam's law: letting $I = \mathbf{1}(\theta \in C(Y))$,

$$P(I = 1) = \mathbb{E}[I] = \mathbb{E}[\mathbb{E}[I \mid Y]] = \mathbb{E}[P(I = 1 \mid Y)] = \mathbb{E}[0.95] = 0.95.$$

- \otimes : A 95% credible interval is *not* guaranteed to be a 95% confidence interval for a specific fixed θ . Coverage at a *fixed* θ may be higher or lower than 95%, depending on the prior since we include uncertainty in the prior. The averaging above is over the prior distribution of θ , not a fixed θ .

Concept Checker 3. Suppose $\theta \sim \mathcal{N}(0, 2^2)$ and $Y \mid \theta \sim \mathcal{N}(\theta, 1)$. We observe $Y = y$.

1. Find the posterior distribution $\theta \mid y$.
2. Write down a 95% credible interval for θ .
3. If the true $\theta = 1$, is the coverage of this credible interval exactly 95%? Explain intuitively why or why not.

Solution

1. By Normal-Normal conjugacy: $\tau_0^2 = 4$, $\sigma^2 = 1$.

$$\tau_1^{-2} = \sigma^{-2} + \tau_0^{-2} = 1 + \frac{1}{4} = \frac{5}{4}, \quad \tau_1^2 = \frac{4}{5}.$$

$$\mu_1 = \tau_1^2 \left(\frac{y}{\sigma^2} + \frac{\mu_0}{\tau_0^2} \right) = \frac{4}{5}(y + 0) = \frac{4y}{5}.$$

So $\theta \mid y \sim \mathcal{N}\left(\frac{4y}{5}, \frac{4}{5}\right)$.

2. The 95% credible interval is $\left[\frac{4y}{5} \pm 1.96\sqrt{\frac{4}{5}}\right] = \left[\frac{4y}{5} \pm 1.754\right]$.

3. No, not exactly 95% at the fixed value $\theta_0 = 1$. The interval is centered at $4y/5$,

and $Y \mid \theta_0 = 1 \sim \mathcal{N}(1, 1)$, so $4Y/5 \mid \theta_0 = 1 \sim \mathcal{N}(4/5, 16/25)$. The coverage is

$$P(1 \in C(Y) \mid \theta_0 = 1) = P\left(\left|\frac{4Y}{5} - 1\right| \leq 1.754 \mid \theta_0 = 1\right) \approx 0.535$$

(greater than 95% because the prior pulls the interval toward 0, away from the true $\theta = 1$, but the interval is wide enough to still often cover). The 95% guarantee is an *average* over the prior on θ —at $\theta_0 = 1$ it can be higher or lower.

6 Conjugate Priors

Idea. A conjugate prior is one where the posterior stays in the same distributional family as the prior. So we only need to update the parameters, not the family itself—which is very convenient for computation and gives us known results.

Definition 4 (Conjugate prior). A family of priors is **conjugate** for a particular likelihood if choosing a prior in the family always results in a posterior in the same family.

6.1 Beta–Binomial

Theorem 3 (Beta-Binomial conjugacy). If $p \sim \text{Beta}(a, b)$ and $Y \mid p \sim \text{Bin}(n, p)$, then

$$p \mid (Y = y) \sim \text{Beta}(a + y, b + n - y).$$

Posterior mean: $\frac{a + y}{a + b + n}$. Interpret $a - 1$ as prior successes, $b - 1$ as prior failures.

6.2 Gamma–Poisson

Theorem 4 (Gamma-Poisson conjugacy). If $\lambda \sim \text{Gamma}(r_0, b_0)$ (rate b_0) and $Y_1, \dots, Y_n \mid \lambda \stackrel{i.i.d.}{\sim} \text{Pois}(\lambda)$, then with $S = \sum y_i$:

$$\lambda \mid y \sim \text{Gamma}(r_0 + S, b_0 + n).$$

Posterior mean: $\frac{r_0 + S}{b_0 + n}$. Predictive: $\tilde{Y} \mid y \sim \text{NBin}\left(r_0 + S, \frac{b_0 + n}{b_0 + n + 1}\right)$.

6.3 Normal–Normal

Theorem 5 (Normal-Normal conjugacy, general sample size). Let $Y_1, \dots, Y_n \stackrel{i.i.d.}{\sim} \mathcal{N}(\mu, \sigma^2)$ with σ^2 known, and prior $\mu \sim \mathcal{N}(\mu_0, \tau_0^2)$. Then

$$\mu \mid y \sim \mathcal{N}(\mu_n, \tau_n^2), \quad \tau_n^{-2} = n\sigma^{-2} + \tau_0^{-2}, \quad \mu_n = \tau_n^2 \left(\frac{n\bar{y}}{\sigma^2} + \frac{\mu_0}{\tau_0^2} \right).$$

Writing $b_n = \tau_n^2 / \tau_0^2 = \sigma^2 / (\sigma^2 + n\tau_0^2)$ (shrinkage factor):

$$\mu_n = (1 - b_n)\bar{y} + b_n\mu_0.$$

The posterior mean is a **precision-weighted average** of the sample mean and the prior mean. As $n \rightarrow \infty$, $b_n \rightarrow 0$ and $\mu_n \rightarrow \bar{y}$.

• **Intuition.** The posterior mean compromises between \bar{y} (the data's estimate) and μ_0 (the prior's estimate), weighted by their relative precisions. More data \Rightarrow less shrinkage toward the prior.

• **Normal-Normal predictive.** By Adam's and Eve's laws:

$$\tilde{Y} | y \sim \mathcal{N}(\mu_n, \sigma^2 + \tau_n^2).$$

The extra τ_n^2 reflects *parameter uncertainty*—even if we knew μ exactly, \tilde{Y} would have variance σ^2 ; our uncertainty about μ adds τ_n^2 on top.

Theorem 6 (Normal with heteroskedasticity). If $Y_j | \mu \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mu, \sigma_j^2)$ and $\mu \sim \mathcal{N}(m_0, \tau_0^2)$, then

$$\tau_n^{-2} = \tau_0^{-2} + \sum_j \sigma_j^{-2}, \quad \mu_n = \tau_n^2 \left(\tau_0^{-2} m_0 + \sum_j \sigma_j^{-2} y_j \right).$$

Observations with smaller variances receive more weight. Reduces to the standard case when all $\sigma_j^2 = \sigma^2$.

Theorem 7 (Bayesian linear regression). If $Y_j | (X = x, \theta) \stackrel{\text{ind.}}{\sim} \mathcal{N}(\theta x_j, \sigma^2)$ and $\theta | X \sim \mathcal{N}(m_0, \tau_0^2)$, then $\theta | y, x \sim \mathcal{N}(m_n, \tau_n^2)$ where

$$\tau_n^{-2} = \tau_0^{-2} + \sigma^{-2} \sum_j x_j^2, \quad m_n = \tau_n^2 \left(\tau_0^{-2} m_0 + \sigma^{-2} \sum_j x_j y_j \right).$$

Concept Checker 4. A manufacturer claims their product weights are $\mathcal{N}(\theta, 100)$ grams (so $\sigma^2 = 100$). Your prior for θ is $\mathcal{N}(200, 400)$. You weigh $n = 25$ items and find $\bar{y} = 190$.

1. Find the posterior distribution of θ .
2. Interpret the shrinkage factor b_n : how much does the posterior mean shrink toward the prior?
3. Find the posterior predictive distribution for a new item's weight.

Solution

1. $\sigma^2 = 100$, $\tau_0^2 = 400$, $n = 25$, $\bar{y} = 190$, $\mu_0 = 200$.

$$\tau_n^{-2} = \frac{25}{100} + \frac{1}{400} = 0.25 + 0.0025 = 0.2525, \quad \tau_n^2 = \frac{1}{0.2525} \approx 3.96.$$

$$\mu_n = 3.96 \left(\frac{25 \times 190}{100} + \frac{200}{400} \right) = 3.96(47.5 + 0.5) = 3.96 \times 48 = 190.08.$$

So $\theta | y \approx \mathcal{N}(190.1, 3.96)$.

2. $b_n = \sigma^2 / (\sigma^2 + n\tau_0^2) = 100 / (100 + 25 \times 400) = 100 / 10100 \approx 0.0099$. The posterior mean is pulled only about 1% toward the prior mean of 200, because the data are plentiful and precise. Almost all weight goes to \bar{y} .

3. $\tilde{Y} | y \sim \mathcal{N}(\mu_n, \sigma^2 + \tau_n^2) \approx \mathcal{N}(190.1, 103.96)$. The extra variance is small because the posterior uncertainty about μ is already small.

7 Practice Problems

Problem 1 (Normal posterior predictive). A random sample of n students is drawn from a population whose weights are $\mathcal{N}(\theta, 400)$ (so $\sigma = 20$) with unknown mean θ . The sample mean is $\bar{y} = 150$. Use prior $\theta \sim \mathcal{N}(180, 1600)$ (so $\tau_0 = 40$).

- Find the posterior distribution $\theta \mid y$ as a function of n .
- Find the posterior predictive distribution $\tilde{y} \mid y$ for a new student's weight, justifying the parameters using Adam's and Eve's laws.
- For $n = 10$: give a 95% posterior interval for θ and a 95% posterior predictive interval for \tilde{y} .

Solution

(a) By Normal-Normal conjugacy with $\sigma^2 = 400$, $\tau_0^2 = 1600$, $\mu_0 = 180$, $\bar{y} = 150$:

$$\tau_n^{-2} = \frac{n}{400} + \frac{1}{1600} = \frac{4n+1}{1600}, \quad \tau_n^2 = \frac{1600}{4n+1}.$$

$$\mu_n = \tau_n^2 \left(\frac{n \cdot 150}{400} + \frac{180}{1600} \right) = \frac{1600}{4n+1} \cdot \frac{600n+180}{1600} = \frac{600n+180}{4n+1} = 150 + \frac{30}{4n+1}.$$

So $\theta \mid y \sim \mathcal{N}\left(150 + \frac{30}{4n+1}, \frac{1600}{4n+1}\right)$.

(b) The posterior predictive $\tilde{y} \mid y$ is Normal (since $(Y_1, \dots, Y_n, \tilde{Y}, \mu)$ is jointly Normal, so any conditional is Normal). By Adam's law:

$$\mathbb{E}[\tilde{y} \mid y] = \mathbb{E}[\mathbb{E}[\tilde{y} \mid y, \theta] \mid y] = \mathbb{E}[\theta \mid y] = \mu_n.$$

By Eve's law:

$$\text{Var}(\tilde{y} \mid y) = \mathbb{E}[\text{Var}(\tilde{y} \mid y, \theta) \mid y] + \text{Var}(\mathbb{E}[\tilde{y} \mid y, \theta] \mid y) = \mathbb{E}[\sigma^2 \mid y] + \text{Var}(\theta \mid y) = \sigma^2 + \tau_n^2.$$

So $\tilde{y} \mid y \sim \mathcal{N}\left(\mu_n, 400 + \frac{1600}{4n+1}\right)$.

(c) For $n = 10$: $\mu_{10} = 150 + 30/41 = 6180/41 \approx 150.73$, $\tau_{10}^2 = 1600/41 \approx 39.02$.

95% posterior interval for θ : $\mu_{10} \pm 1.96\sqrt{\tau_{10}^2} \approx 150.73 \pm 1.96 \times 6.25 \approx [138.5, 163.0]$.

For \tilde{y} : $\sigma^2 + \tau_{10}^2 = 400 + 39.02 = 439.02$, so 95% predictive interval: $150.73 \pm 1.96\sqrt{439.02} \approx 150.73 \pm 41.1 \approx [109.7, 191.8]$.

Note: the predictive interval is much wider than the posterior interval, since it must also account for the $\sigma^2 = 400$ sampling variance of each individual.

Problem 2 (Posterior as a compromise). Let Y be the number of heads in n coin flips with unknown probability θ .

- With a $\text{Unif}(0, 1)$ prior, derive the prior predictive distribution $P(Y = k)$ for each $k = 0, \dots, n$.

- (b) With $\theta \sim \text{Beta}(\alpha, \beta)$ and y heads observed, show algebraically that the posterior mean always lies strictly between the prior mean $\frac{\alpha}{\alpha+\beta}$ and the observed frequency $\frac{y}{n}$.
- (c) Show that if the prior is $\text{Unif}(0, 1)$, the posterior variance is always less than the prior variance.
- (d) Give an example of a $\text{Beta}(\alpha, \beta)$ prior and data (n, y) where the posterior variance *exceeds* the prior variance.

Solution

(a) $\text{Unif}(0, 1) = \text{Beta}(1, 1)$.

$$P(Y = k) = \int_0^1 \binom{n}{k} \theta^k (1-\theta)^{n-k} d\theta = \binom{n}{k} B(k+1, n-k+1) = \binom{n}{k} \cdot \frac{k! (n-k)!}{(n+1)!} = \frac{1}{n+1}.$$

The prior predictive is $\text{Unif}\{0, 1, \dots, n\}$ —all outcomes equally likely. This is known as Bayes-Laplace.

(b) By Beta-Binomial conjugacy, $\theta \mid y \sim \text{Beta}(\alpha + y, \beta + n - y)$. The posterior mean is

$$\mathbb{E}[\theta \mid y] = \frac{\alpha + y}{\alpha + \beta + n} = \frac{(\alpha + \beta) \cdot \frac{\alpha}{\alpha + \beta} + n \cdot \frac{y}{n}}{\alpha + \beta + n},$$

which is a weighted average (with positive weights $\alpha + \beta$ and n) of the prior mean $\frac{\alpha}{\alpha + \beta}$ and the sample frequency $\frac{y}{n}$. A weighted average of two distinct values lies strictly between them.

(c) With $\text{Unif}(0, 1) = \text{Beta}(1, 1)$: prior variance = $\frac{1 \cdot 1}{2^2 \cdot 3} = \frac{1}{12}$. Posterior is $\text{Beta}(1 + y, 1 + n - y)$, with variance

$$\text{Var}(\theta \mid y) = \frac{(1 + y)(1 + n - y)}{(2 + n)^2(3 + n)} = \frac{1 + y}{2 + n} \cdot \frac{1 + n - y}{2 + n} \cdot \frac{1}{3 + n} \leq \frac{1}{4} \cdot \frac{1}{3 + n} < \frac{1}{12}.$$

The first inequality uses $ab \leq (a + b)^2/4$; the last uses $n \geq 1$.

(d) Take $\alpha = 1$, $\beta = 12$ (prior mean ≈ 0.077 , prior variance = $\frac{1 \cdot 12}{13^2 \cdot 14} \approx 0.00507$). With $n = 6$, $y = 3$: posterior $\text{Beta}(4, 15)$, posterior variance = $\frac{4 \cdot 15}{19^2 \cdot 20} \approx 0.00831 > 0.00507$. This is possible because the data ($y/n = 0.5$) are far from the prior mean, causing the posterior to be spread across a wider range than the prior. (For large n this cannot happen—on average the posterior variance is always smaller.)

Problem 3 (Bayesian persuasion). Let $R \sim \text{Unif}(0, 1)$ be the probability of a successful persuasion. Leo will stop resisting once he experiences s successes; let S_{total} be the total number of attempts needed.

- (a) Show that $\mathbb{E}[S_{\text{total}}]$ is infinite.
- (b) After n total attempts with exactly s successes, find the posterior distribution of R .

Solution

(a) Given $R = p$, the number of failures before the s -th success is $\text{NBin}(s, p)$, so $S_{\text{total}} - s \mid R = p \sim \text{NBin}(s, p)$ and $\mathbb{E}[S_{\text{total}} \mid R = p] = s/p$. By LOTE:

$$\mathbb{E}[S_{\text{total}}] = \int_0^1 \frac{s}{p} \cdot 1 \, dp = s \int_0^1 p^{-1} \, dp,$$

which diverges. So the expected number of attempts is infinite.

(b) After n attempts with s successes and $n - s$ failures, the likelihood is $\propto p^s(1-p)^{n-s}$. With $\text{Unif}(0, 1) = \text{Beta}(1, 1)$ prior:

$$\pi(R \mid S_{\text{total}} = n) \propto p^s(1-p)^{n-s} \cdot 1 = p^{(s+1)-1}(1-p)^{(n-s+1)-1},$$

so $R \mid S_{\text{total}} = n \sim \text{Beta}(s + 1, n - s + 1)$.

Problem 4 (Censored and uncensored data). Suppose $Y \mid \theta \sim \text{Expo}(\theta)$ with conjugate prior $\theta \sim \text{Gamma}(\alpha, \beta)$.

- We observe $Y \geq 100$ (but not the exact value). Find $\pi(\theta \mid Y \geq 100)$, and the posterior mean and variance.
- Now suppose we are told $Y = 100$ exactly. Find $\pi(\theta \mid Y = 100)$, and the posterior mean and variance.
- Why is the posterior variance *higher* in (b), even though more information was given?

Solution

(a) $P(Y \geq 100 \mid \theta) = e^{-100\theta}$. So

$$\pi(\theta \mid Y \geq 100) \propto \pi(\theta) P(Y \geq 100 \mid \theta) = \theta^{\alpha-1} e^{-\beta\theta} \cdot e^{-100\theta} = \theta^{\alpha-1} e^{-(\beta+100)\theta},$$

giving $\theta \mid Y \geq 100 \sim \text{Gamma}(\alpha, \beta+100)$. Posterior mean = $\frac{\alpha}{\beta+100}$; variance = $\frac{\alpha}{(\beta+100)^2}$.

(b) $f(Y = 100 \mid \theta) = \theta e^{-100\theta}$. So

$$\pi(\theta \mid Y = 100) \propto \theta^{\alpha-1} e^{-\beta\theta} \cdot \theta e^{-100\theta} = \theta^{(\alpha+1)-1} e^{-(\beta+100)\theta},$$

giving $\theta \mid Y = 100 \sim \text{Gamma}(\alpha + 1, \beta + 100)$. Posterior mean = $\frac{\alpha+1}{\beta+100}$; variance = $\frac{\alpha+1}{(\beta+100)^2}$.

(c) Posterior variance = $\frac{\alpha+1}{(\beta+100)^2} > \frac{\alpha}{(\beta+100)^2}$.

The key insight is that Eve's law says $\text{Var}(\theta) \leq \mathbb{E}[\text{Var}(\theta \mid Y)] + \text{Var}(\mathbb{E}[\theta \mid Y])$, so the posterior variance is *on average* smaller than the prior variance—but not necessarily for every single observation. Observing $Y = 100$ exactly tells us the rate is moderate: it is unlikely to be very high (which would make $Y = 100$ very improbable) or very low (same reason). This pushes probability mass toward an intermediate range, which can actually *spread* the distribution more than the one-sided constraint $Y \geq 100$ did.

Censored data tells us θ is small (rate is low $\Rightarrow Y$ is large), concentrating mass; exact data tells us θ is moderate, which is less restrictive.