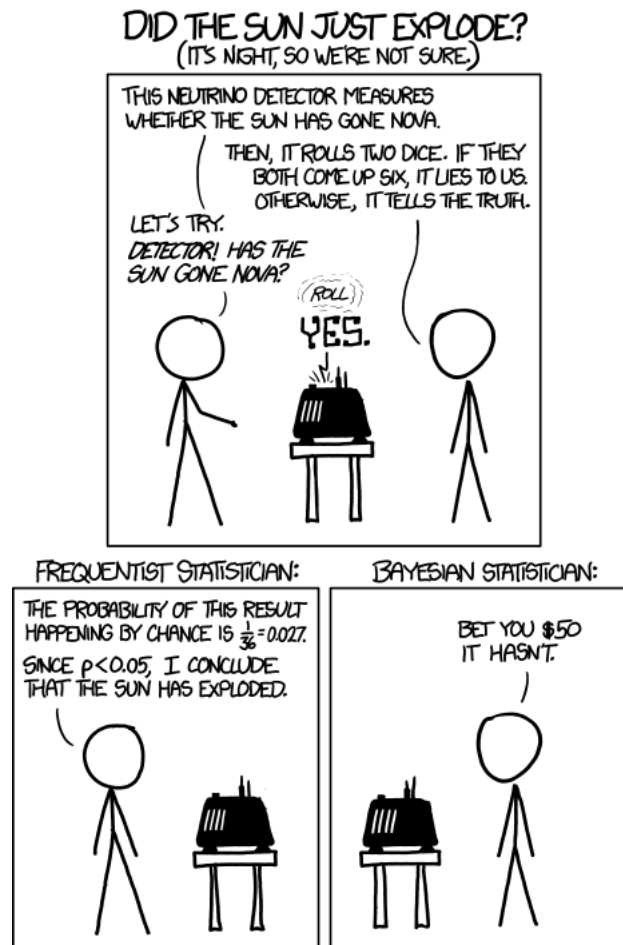


Notes on Bayesian Inference



xkcd.com/1132/

1. What's the Difference?

There are two main approaches to statistical inference, *frequentist* (or classical) methods and *Bayesian* methods. The key differences between the frequentist and Bayesian approaches are as follows:

	frequentist	Bayesian
probability is:	limiting relative frequency	degree of subjective belief
parameter θ is a:	fixed constant	random variable
probability statements are about:	procedures	parameters
frequency guarantees?	yes	no

The field of statistics puts more emphasis on frequentist methods although Bayesian methods certainly have a presence. The machine learning community embraces Bayesian methods more strongly. There are, in fact, many flavors of Bayesian inference. *Subjective Bayesians* interpret probability strictly as personal degrees of belief. *Objective Bayesians* try to find prior distributions that formally express ignorance with the hope that the resulting posterior is, in some sense, objective. *Empirical Bayesians* estimate the prior distribution from the data. *Frequentist Bayesians* are those who use Bayesian methods only when the resulting posterior has good frequency behavior.

2. The Bayesian Method

Let x_1, \dots, x_n be n observations sampled from a probability density $p(x | \theta)$. In this chapter, we write $p(x | \theta)$ if we view θ as a random variable and $p(x | \theta)$ represents the conditional probability density conditioned on θ . In contrast, we write $p_\theta(x)$ if we view θ as a deterministic value. Bayesian inference is usually carried out in the following way.

1. We choose a probability density $\pi(\theta)$ — called the *prior distribution* — that expresses our beliefs about a parameter θ before we see any data.
2. We choose a statistical model $p(x | \theta)$ that reflects our beliefs about x given θ .
3. After observing data $\mathcal{D}_n = \{x_1, \dots, x_n\}$, we update our beliefs and calculate the *posterior distribution* $p(\theta | \mathcal{D}_n)$.

By Bayes' theorem, the posterior distribution can be written as

$$p(\theta | x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n | \theta)\pi(\theta)}{p(x_1, \dots, x_n)} = \frac{\mathcal{L}_n(\theta)\pi(\theta)}{c_n} \propto \mathcal{L}_n(\theta)\pi(\theta) \quad (1)$$

where $\mathcal{L}_n(\theta) = \prod_{i=1}^n p(x_i | \theta)$ is the *likelihood function* and

$$c_n = p(x_1, \dots, x_n) = \int p(x_1, \dots, x_n | \theta)\pi(\theta)d\theta = \int \mathcal{L}_n(\theta)\pi(\theta)d\theta$$

is the normalizing constant, which is also called *evidence*.

We can get a *Bayesian point estimate* by summarizing the center of the posterior. Typically, we use the mean or mode of the posterior distribution. The posterior mean is

$$\bar{\theta}_n = \int \theta p(\theta | \mathcal{D}_n) d\theta = \frac{\int \theta \mathcal{L}_n(\theta) \pi(\theta) d\theta}{\int \mathcal{L}_n(\theta) \pi(\theta) d\theta}. \quad (2)$$

We can also obtain a *Bayesian interval estimate*. For example, for $\alpha \in (0, 1)$, we could find a and b such that

$$\int_{-\infty}^a p(\theta | \mathcal{D}_n) d\theta = \int_b^{\infty} p(\theta | \mathcal{D}_n) d\theta = \alpha/2.$$

Let $C = (a, b)$. Then

$$\mathbb{P}(\theta \in C | \mathcal{D}_n) = \int_a^b p(\theta | \mathcal{D}_n) d\theta = 1 - \alpha$$

so C is a $1 - \alpha$ *Bayesian posterior interval* or *credible interval*. If θ has more than one dimension, the extension is straightforward and we obtain a *credible region*.

Example 2.1. Let $X \sim \text{Bernoulli}(\theta)$ and we have observed data $\mathcal{D}_n = \{x_1, \dots, x_n\}$. Suppose we take the uniform distribution $\pi(\theta) = 1$ as a prior. By Bayes' theorem, the posterior is

$$p(\theta | \mathcal{D}_n) \propto \pi(\theta) \mathcal{L}_n(\theta) = \theta^s (1 - \theta)^{n-s} = \theta^{s+1-1} (1 - \theta)^{n-s+1-1}$$

where $s = \sum_{i=1}^n x_i$ is the number of successes. Recall that a random variable θ on the interval $(0, 1)$ has a Beta distribution with parameters α and β if its density is

$$\pi_{\alpha, \beta}(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}.$$

We see that the posterior distribution for θ is a Beta distribution with parameters $s+1$ and $n-s+1$. That is,

$$p(\theta | \mathcal{D}_n) = \frac{\Gamma(n+2)}{\Gamma(s+1)\Gamma(n-s+1)} \theta^{(s+1)-1} (1 - \theta)^{(n-s+1)-1}.$$

We write this as

$$\theta | \mathcal{D}_n \sim \text{Beta}(s+1, n-s+1).$$

Notice that we have figured out the normalizing constant without actually doing the integral $\int \mathcal{L}_n(\theta) \pi(\theta) d\theta$. Since a density function integrates to one, we see that

$$\int_0^1 \theta^s (1 - \theta)^{n-s} d\theta = \frac{\Gamma(s+1)\Gamma(n-s+1)}{\Gamma(n+2)}. \quad (3)$$

The mean of a $\text{Beta}(\alpha, \beta)$ distribution is $\alpha/(\alpha + \beta)$ so the Bayes posterior estimator is

$$\bar{\theta} = \frac{s+1}{n+2}. \quad (4)$$

It is instructive to rewrite $\bar{\theta}$ as

$$\bar{\theta} = \lambda_n \hat{\theta} + (1 - \lambda_n) \tilde{\theta} \quad (5)$$

where $\hat{\theta} = s/n$ is the maximum likelihood estimate, $\tilde{\theta} = 1/2$ is the prior mean and $\lambda_n = n/(n + 2) \approx 1$. A 95 percent posterior interval can be obtained by numerically finding a and b such that $\int_a^b p(\theta | \mathcal{D}_n) d\theta = .95$.

Suppose that instead of a uniform prior, we use the prior $\theta \sim \text{Beta}(\alpha, \beta)$. If you repeat the calculations above, you will see that $\theta | \mathcal{D}_n \sim \text{Beta}(\alpha + s, \beta + n - s)$. The flat prior is just the special case with $\alpha = \beta = 1$. The posterior mean in this more general case is

$$\bar{\theta} = \frac{\alpha + s}{\alpha + \beta + n} = \left(\frac{n}{\alpha + \beta + n} \right) \hat{\theta} + \left(\frac{\alpha + \beta}{\alpha + \beta + n} \right) \theta_0$$

where $\theta_0 = \alpha/(\alpha + \beta)$ is the prior mean.

An illustration of this example is shown in Figure 1. We use the Bernoulli model to generate $n = 15$ data with parameter $\theta = 0.4$. We observe $s = 7$. Therefore, the maximum likelihood estimate is $\hat{\theta} = 7/15 = 0.47$, which is larger than the true parameter value 0.4. The left plot of Figure 1 adopts a prior $\text{Beta}(4, 6)$ which gives a posterior mode 0.43, while the right plot of Figure 1 adopts a prior $\text{Beta}(4, 2)$ which gives a posterior mode 0.67.

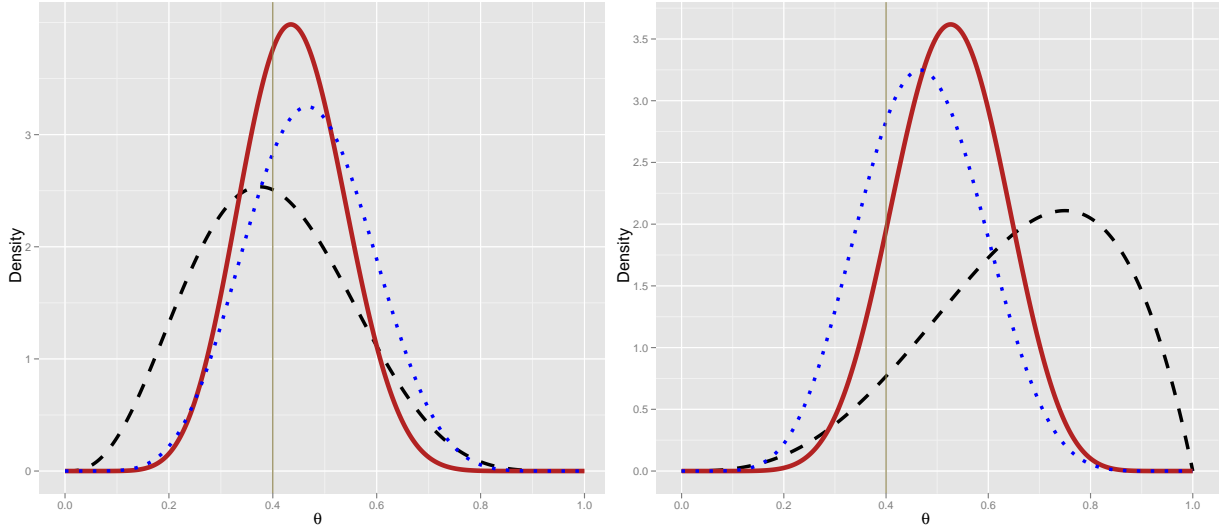


Figure 1: Illustration of Bayesian inference on Bernoulli data with two priors. The three curves are prior distribution (black-dashed), likelihood function (blue-dotted), and the posterior distribution (red-solid). The true parameter value $\theta = 0.4$ is indicated by the vertical line.

Example 2.2. Let $\theta = (\theta_1, \dots, \theta_K)$ be a K -dimensional parameter ($K > 1$). The multinomial model with a Dirichlet prior is a generalization of the Bernoulli model and Beta prior of the previous example. The Dirichlet distribution for K outcomes is the exponential family distribution on

the $K - 1$ dimensional probability simplex¹ Δ_K given by

$$\pi_\alpha(\theta) = \frac{\Gamma(\sum_{j=1}^K \alpha_j)}{\prod_{j=1}^K \Gamma(\alpha_j)} \prod_{j=1}^K \theta_j^{\alpha_j-1} \quad (6)$$

where $\alpha = (\alpha_1, \dots, \alpha_K) \in \mathbb{R}_+^K$ is a non-negative vector of scaling coefficients, which are the parameters of the model. We can think of the sample space of the multinomial with K outcomes as the set of vertices of the K -dimensional hypercube \mathbb{H}_K , made up of vectors with exactly one 1 and the remaining elements 0:

$$x = \underbrace{(0, 0, \dots, 0, 1, 0, \dots, 0)}_{K \text{ places}}. \quad (7)$$

Let $x_i = (x_{i1}, \dots, x_{iK}) \in \mathbb{H}_K$. If

$$\theta \sim \text{Dirichlet}(\alpha) \quad \text{and} \quad x_i | \theta \sim \text{Multinomial}(\theta) \quad \text{for } i = 1, 2, \dots, n \quad (8)$$

then the posterior satisfies:

$$p(\theta | x_1, \dots, x_n) \propto \mathcal{L}_n(\theta) \pi(\theta) \propto \prod_{i=1}^n \prod_{j=1}^K \theta_j^{x_{ij}} \prod_{j=1}^K \theta_j^{\alpha_j-1} = \prod_{j=1}^K \theta_j^{\sum_{i=1}^n x_{ij} + \alpha_j - 1}. \quad (9)$$

We see that the posterior is also a Dirichlet distribution:

$$\theta | x_1, x_2, \dots, x_n \sim \text{Dirichlet}(\alpha + n\bar{x}) \quad (10)$$

where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \in \Delta_K$.

Since the mean of a Dirichlet distribution $\pi_\alpha(\theta)$ is given by

$$\mathbb{E}(\theta) = \left(\frac{\alpha_1}{\sum_{i=1}^K \alpha_i}, \dots, \frac{\alpha_K}{\sum_{i=1}^K \alpha_i} \right), \quad (11)$$

the posterior mean of a multinomial with Dirichlet prior is

$$\mathbb{E}(\theta | x_1, \dots, x_n) = \left(\frac{\alpha_1 + \sum_{i=1}^n x_{i1}}{\sum_{i=1}^K \alpha_i + n}, \dots, \frac{\alpha_K + \sum_{i=1}^n x_{iK}}{\sum_{i=1}^K \alpha_i + n} \right). \quad (12)$$

¹ The probability simplex Δ_K is defined as

$$\Delta_K = \left\{ \theta = (\theta_1, \dots, \theta_K) \in \mathbb{R}^K \mid \theta_i \geq 0 \text{ for all } i \text{ and } \sum_{i=1}^K \theta_i = 1 \right\}.$$

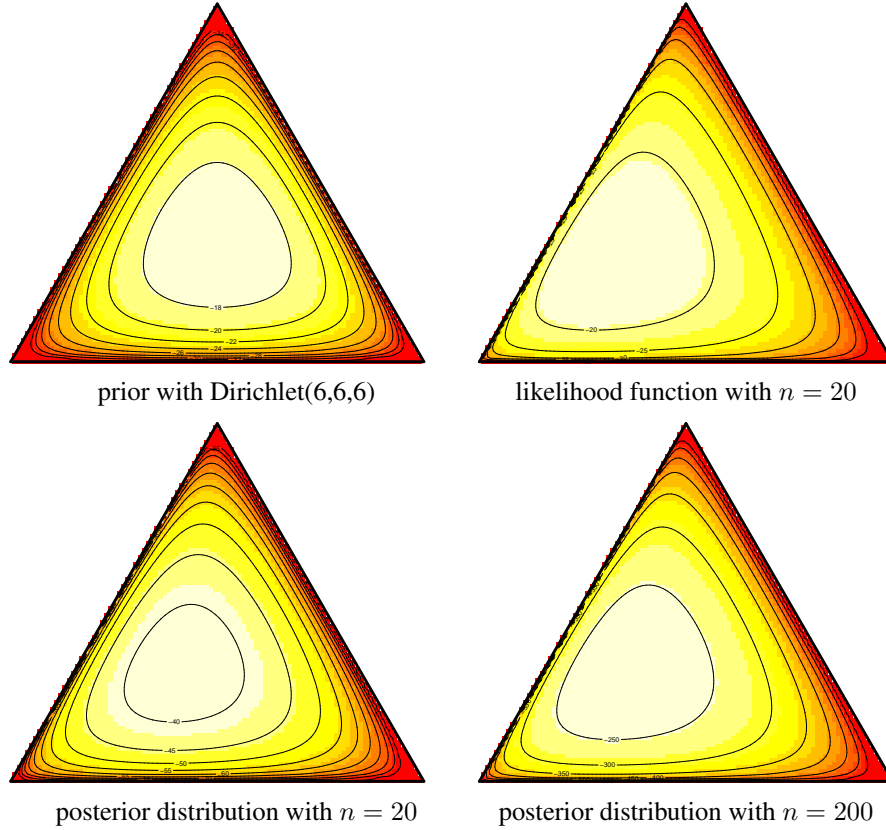


Figure 2: Illustration of Bayesian inference on multinomial data with the prior $\text{Dirichlet}(6, 6, 6)$. The contours of the prior, likelihood, and posteriors are plotted on a two-dimensional probability simplex (Starting from the bottom left vertex of each triangle, clock-wisely the three vertices correspond to $\theta_1, \theta_2, \theta_3$). We see that when the number of observed data is small, the posterior is affected by both the prior and the likelihood; when the number of observed data is large, the posterior is mainly dominated by the likelihood.

This again can be viewed as smoothing out the maximum likelihood estimate by allocating some additional probability mass to low frequency observations. The parameters $\alpha_1, \dots, \alpha_K$ act as “virtual counts” that don’t actually appear in the observed data.

An illustration of this example is shown in Figure 2. We use the multinomial model to generate $n = 20$ data points with parameter $\theta = (0.2, 0.3, 0.5)$. We adopt a prior $\text{Dirichlet}(6, 6, 6)$. The contours of the prior, likelihood, and posterior with $n = 20$ observed data are shown in the first three plots in Figure 2. As a comparison, we also provide the contour of the posterior with $n = 200$ observed data in the last plot. From this experiment, we see that when the number of observed data is small, the posterior is affected by both the prior and the likelihood; when the number of observed data is large, the posterior is mainly dominated by the likelihood.

In the previous two examples, the prior was a Dirichlet distribution and the posterior was also a

Dirichlet. When the prior and the posterior are in the same family, we say that the prior is *conjugate* with respect to the model; this will be discussed further below.

Example 2.3. Let $X \sim N(\theta, \sigma^2)$ and $\mathcal{D}_n = \{x_1, \dots, x_n\}$ be the observed data. For simplicity, let us assume that σ is known and we want to estimate $\theta \in \mathbb{R}$. Suppose we take as a prior $\theta \sim N(a, b^2)$. Let $\bar{x} = \sum_{i=1}^n x_i/n$ be the sample mean. It can be shown that the posterior for θ is

$$\theta | \mathcal{D}_n \sim N(\bar{\theta}, \tau^2) \quad (13)$$

where

$$\begin{aligned} \bar{\theta} &= w\hat{\theta} + (1-w)a, \\ \hat{\theta} &= \bar{x}, \quad w = \frac{\frac{1}{se^2}}{\frac{1}{se^2} + \frac{1}{b^2}}, \quad \frac{1}{\tau^2} = \frac{1}{se^2} + \frac{1}{b^2}, \end{aligned}$$

and $se = \sigma/\sqrt{n}$ is the standard error of the maximum likelihood estimate $\hat{\theta}$. This is another example of a conjugate prior. Note that $w \rightarrow 1$ and $\tau/se \rightarrow 1$ as $n \rightarrow \infty$. So, for large n , the posterior is approximately $N(\hat{\theta}, se^2)$. The same is true if n is fixed but $b \rightarrow \infty$, which corresponds to letting the prior become very flat.

Continuing with this example, let us find $C = (c, d)$ such that $\mathbb{P}(\theta \in C | \mathcal{D}_n) = 0.95$. We can do this by choosing c and d such that $\mathbb{P}(\theta < c | \mathcal{D}_n) = 0.025$ and $\mathbb{P}(\theta > d | \mathcal{D}_n) = 0.025$. More specifically, we want to find c such that

$$\mathbb{P}(\theta < c | \mathcal{D}_n) = \mathbb{P}\left(\frac{\theta - \bar{\theta}}{\tau} < \frac{c - \bar{\theta}}{\tau} \mid \mathcal{D}_n\right) = \mathbb{P}\left(Z < \frac{c - \bar{\theta}}{\tau}\right) = 0.025$$

where $Z \sim N(0, 1)$ is a standard Gaussian random variable. We know that $\mathbb{P}(Z < -1.96) = 0.025$. So,

$$\frac{c - \bar{\theta}}{\tau} = -1.96$$

implying that $c = \bar{\theta} - 1.96\tau$. By similar arguments, $d = \bar{\theta} + 1.96\tau$. So a 95 percent Bayesian credible interval is $\bar{\theta} \pm 1.96\tau$. Since $\bar{\theta} \approx \hat{\theta}$ and $\tau \approx se$ when n is large, the 95 percent Bayesian credible interval is approximated by $\hat{\theta} \pm 1.96 se$ which is the frequentist confidence interval.

3. Bayesian Prediction

After the data $\mathcal{D}_n = \{x_1, \dots, x_n\}$ have been observed, the Bayesian framework allows us to predict the distribution of a future data point x conditioned on \mathcal{D}_n . To do this, we first obtain the posterior

$p(\theta | \mathcal{D}_n)$. Then

$$p(x | \mathcal{D}_n) = \int p(x, \theta | \mathcal{D}_n) d\theta \quad (14)$$

$$= \int p(x | \theta, \mathcal{D}_n) p(\theta | \mathcal{D}_n) d\theta \quad (15)$$

$$= \int p(x | \theta) p(\theta | \mathcal{D}_n) d\theta. \quad (16)$$

Where we use the fact that $p(x | \theta, \mathcal{D}_n) = p(x | \theta)$ since all the data are conditionally independent given θ . From the last line, the predictive distribution $p(x | \mathcal{D}_n)$ can be viewed as a weighted average of likelihood $p(x | \theta)$. The weights are determined by the posterior distribution of θ .

Example 3.1. Under a Bernoulli model $X \sim \text{Bernoulli}(\theta)$, let $\mathcal{D}_n = \{x_1, \dots, x_n\}$ be the observed data and $\pi(\theta) = 1$ so that $\theta | \mathcal{D}_n \sim \text{Beta}(s + 1, n - s + 1)$ with $s = \sum_{i=1}^n x_i$. We define $\psi = \log(\theta/(1 - \theta))$. Then

$$\begin{aligned} H(t | \mathcal{D}_n) &= \mathbb{P}(\psi \leq t | \mathcal{D}_n) = \mathbb{P}\left(\log\left(\frac{\theta}{1 - \theta}\right) \leq t \mid \mathcal{D}_n\right) \\ &= \mathbb{P}\left(\theta \leq \frac{e^t}{1 + e^t} \mid \mathcal{D}_n\right) \\ &= \int_0^{e^t/(1+e^t)} p(\theta | \mathcal{D}_n) d\theta \\ &= \frac{\Gamma(n + 2)}{\Gamma(s + 1)\Gamma(n - s + 1)} \int_0^{e^t/(1+e^t)} \theta^s (1 - \theta)^{n-s} d\theta \end{aligned}$$

and

$$p(\psi | \mathcal{D}_n) = H'(\psi | \mathcal{D}_n) = \frac{\Gamma(n + 2)}{\Gamma(s + 1)\Gamma(n - s + 1)} \left(\frac{e^\psi}{1 + e^\psi}\right)^s \left(\frac{1}{1 + e^\psi}\right)^{n-s+2}$$

for $\psi \in \mathbb{R}$.

4 Multiparameter Problems

Let $\mathcal{D}_n = \{x_1, \dots, x_n\}$ be the observed data. Suppose that $\theta = (\theta_1, \dots, \theta_d)$ with some prior distribution $\pi(\theta)$. The posterior density is still given by

$$p(\theta | \mathcal{D}_n) \propto \mathcal{L}_n(\theta) \pi(\theta). \quad (17)$$

The question now arises of how to extract inferences about one single parameter. The key is to find the marginal posterior density for the parameter of interest. Suppose we want to make inferences

about θ_1 . The marginal posterior for θ_1 is

$$p(\theta_1 | \mathcal{D}_n) = \int \cdots \int p(\theta_1, \dots, \theta_d | \mathcal{D}_n) d\theta_2 \dots d\theta_d. \quad (18)$$

In practice, it might not be feasible to do this integral. Simulation can help: we draw randomly from the posterior:

$$\theta^1, \dots, \theta^B \sim p(\theta | \mathcal{D}_n)$$

where the superscripts index different draws. Each θ^j is a vector $\theta^j = (\theta_1^j, \dots, \theta_d^j)$. Now collect together the first component of each draw: $\theta_1^1, \dots, \theta_1^B$. These are a sample from $p(\theta_1 | \mathcal{D}_n)$ and we have avoided doing any integrals. One thing to note is, sampling B data from a multivariate distribution $p(\theta | \mathcal{D}_n)$ is challenging especially when the dimensionality d is large. We will discuss this topic further in the chapter on computing.

Example 4.2. (Comparing Two Binomials) Suppose we have n_1 control patients and n_2 treatment patients and that x_1 control patients survive while x_2 treatment patients survive. We assume the Binomial model:

$$X_1 \sim \text{Binomial}(n_1, \theta_1) \text{ and } X_2 \sim \text{Binomial}(n_2, \theta_2).$$

We want to estimate $\tau = g(\theta_1, \theta_2) = \theta_2 - \theta_1$.

If $\pi(\theta_1, \theta_2) = 1$, the posterior is

$$p(\theta_1, \theta_2 | x_1, x_2) \propto \theta_1^{x_1} (1 - \theta_1)^{n_1 - x_1} \theta_2^{x_2} (1 - \theta_2)^{n_2 - x_2}.$$

Notice that (θ_1, θ_2) live on a rectangle (a square, actually) and that

$$p(\theta_1, \theta_2 | x_1, x_2) = p(\theta_1 | x_1) p(\theta_2 | x_2)$$

where

$$p(\theta_1 | x_1) \propto \theta_1^{x_1} (1 - \theta_1)^{n_1 - x_1} \text{ and } p(\theta_2 | x_2) \propto \theta_2^{x_2} (1 - \theta_2)^{n_2 - x_2}$$

which implies that θ_1 and θ_2 are independent under the posterior. Also, $\theta_1 | x_1 \sim \text{Beta}(x_1 + 1, n_1 - x_1 + 1)$ and $\theta_2 | x_2 \sim \text{Beta}(x_2 + 1, n_2 - x_2 + 1)$. If we simulate $\Theta_1^1, \dots, \Theta_1^B \sim \text{Beta}(x_1 + 1, n_1 - x_1 + 1)$ and $\Theta_2^1, \dots, \Theta_2^B \sim \text{Beta}(x_2 + 1, n_2 - x_2 + 1)$, then $\tau_b = \Theta_2^b - \Theta_1^b$, $b = 1, \dots, B$, is a sample from $p(\tau | x_1, x_2)$.

5. Simulation

Since the posterior distribution $p(\theta | \mathcal{D}_n)$ generally involves high dimensional integrals, it is generally approximated by simulation. Suppose we draw $\theta^1, \dots, \theta^B \sim p(\theta | \mathcal{D}_n)$. Then a histogram of $\theta^1, \dots, \theta^B$ approximates the posterior density $p(\theta | \mathcal{D}_n)$. An approximation to the posterior mean $\bar{\theta}_n = \mathbb{E}(\theta | \mathcal{D}_n)$ is $B^{-1} \sum_{j=1}^B \theta^j$. The posterior $1 - \alpha$ interval can be approximated by $(\theta_{\alpha/2}, \theta_{1-\alpha/2})$ where $\theta_{\alpha/2}$ is the $\alpha/2$ sample quantile of $\theta^1, \dots, \theta^B$.

Once we have a sample $\theta^1, \dots, \theta^B$ from $p(\theta | \mathcal{D}_n)$, let $\tau^i = g(\theta^i)$. Then τ^1, \dots, τ^B is a sample from $p(\tau | \mathcal{D}_n)$. This avoids the need to do any analytical calculations. Simulation techniques are discussed in more detail in a later chapter on statistical computing.

Example 5.1. Consider again Example 4.2. We can approximate the posterior for ψ without doing any calculus. Here are the steps:

1. Draw $\theta^1, \dots, \theta^B \sim \text{Beta}(s+1, n-s+1)$.
2. Let $\psi^i = \log(\theta^i/(1-\theta^i))$ for $i = 1, \dots, B$.

Now ψ^1, \dots, ψ^B are i.i.d. draws from the posterior density $p(\psi | \mathcal{D}_n)$. A histogram of these values provides an estimate of $p(\psi | \mathcal{D}_n)$.

6. Bayesian Linear Models

Many frequentist methods can be viewed as the *maximum a posteriori* (MAP) estimator under a Bayesian framework. As an example, we consider Gaussian linear regression:

$$Y = \beta_0 + \sum_{j=1}^d \beta_j X_j + \epsilon, \quad \epsilon \sim N(0, \sigma^2). \quad (19)$$

Here we assume that σ is known. Let $\mathcal{D}_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ be the observed data points. The conditional likelihood of $\beta = (\beta_0, \beta_1, \dots, \beta_d)$ can be written as

$$\mathcal{L}(\beta) = \prod_{i=1}^n p(y_i | x_i, \beta) \propto \exp\left(-\frac{\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij}\right)^2}{2\sigma^2}\right). \quad (20)$$

Using a Gaussian prior $\pi_\lambda(\beta) \propto \exp(-\lambda \|\beta\|_2^2/2)$, the posterior of β can be written as

$$p(\beta | \mathcal{D}_n) \propto \mathcal{L}(\beta) \pi_\lambda(\beta). \quad (21)$$

The “maximum a posteriori estimator” or MAP estimator $\hat{\beta}^{\text{MAP}}$ takes the form

$$\hat{\beta}^{\text{MAP}} = \arg \max_{\beta} p(\beta | \mathcal{D}_n) = \arg \min_{\beta} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij})^2 + \lambda \sigma^2 \|\beta\|_2^2 \right\}. \quad (22)$$

This is exactly the ridge regression with the regularization parameter $\lambda' = \lambda \sigma^2$. If we adopt the Laplacian prior $\pi_\lambda(\beta) \propto \exp(-\lambda \|\beta\|_1/2)$, we get the Lasso estimator

$$\hat{\beta}^{\text{MAP}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij})^2 + \lambda \sigma^2 \|\beta\|_1 \right\}. \quad (23)$$

Instead of using the MAP point estimate, a full Bayesian inference aims at obtaining the whole posterior distribution $p(\beta \mid \mathcal{D}_n)$. In general, $p(\beta \mid \mathcal{D}_n)$ does not have an analytic form and we need to resort to simulation to approximate the posterior.