

# Analytically-Tractable Bayesian Models

# Conjugate Prior Distributions

For most Bayesian inference problems, the integrals needed to do inference and prediction are not analytically tractable — hence the need for numerical quadrature, Monte Carlo methods, or various approximations.

Most of the exceptions involve *conjugate priors*, which combine nicely with the likelihood to give a posterior distribution of the same form. Examples:

- 1) Independent observations from a finite set, with Beta / Dirichlet priors.
- 2) Independent observations of Gaussian variables with Gaussian prior for the mean, and either known variance or inverse-Gamma prior for the variance.
- 3) Linear regression with Gaussian prior for the regression coefficients, and Gaussian noise, with known variance or inverse-Gamma prior for the variance.

It's nice when a tractable model and prior are appropriate for the problem.

Unfortunately, people are tempted to use such models and priors even when they aren't appropriate.

# Independent Binary Observations with Beta Prior

We observe binary (0/1) variables  $Y_1, Y_2, \dots, Y_n$ .

We model these as being *independent*, and *identically distributed*, with

$$P(Y_i = y | \mu) = \begin{cases} \mu & \text{if } y = 1 \\ 1 - \mu & \text{if } y = 0 \end{cases} = \mu^y (1 - \mu)^{1-y}$$

Let's suppose that our prior distribution for  $\mu$  is Beta( $a, b$ ), with  $a$  and  $b$  being known positive reals. With this prior, the probability density over  $(0, 1)$  of  $\mu$  is:

$$P(\mu) = \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1 - \mu)^{b-1}$$

Here, the Gamma function,  $\Gamma(c)$ , is defined to be  $\int_0^\infty x^{c-1} \exp(-x) dx$ .

For integer  $c$ ,  $\Gamma(c) = (c - 1)!$ .

Note that when  $a = b = 1$  the prior is uniform over  $(0, 1)$ .

The prior mean of  $\mu$  is  $a/(a + b)$ . Big  $a$  and  $b$  give smaller prior variance.

# Posterior Distribution with Beta Prior

With this Beta prior, the posterior distribution is also Beta:

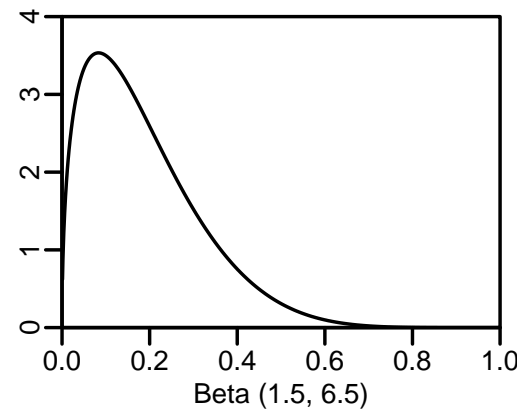
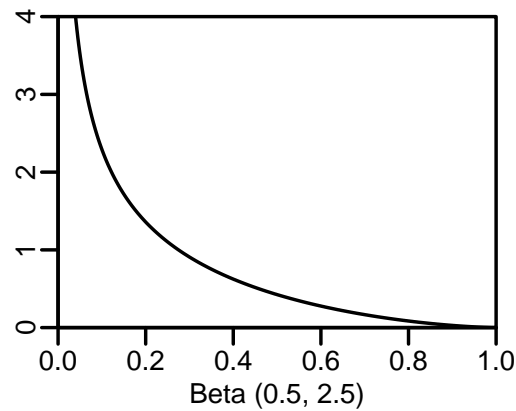
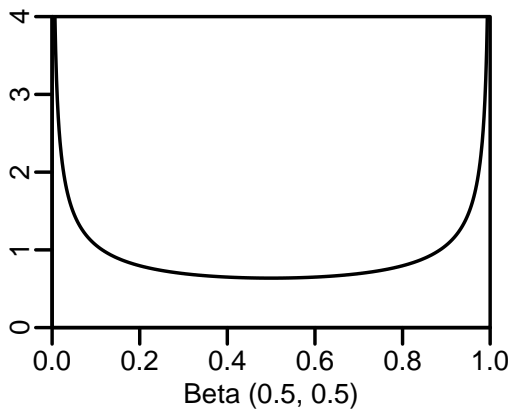
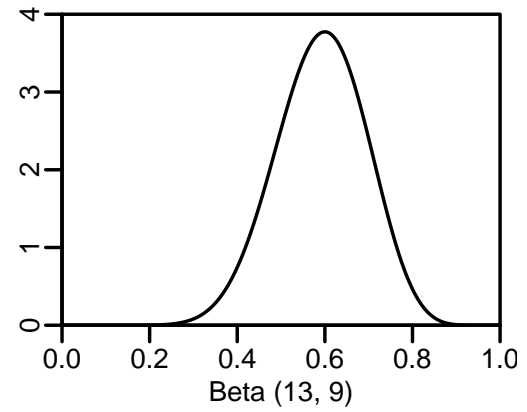
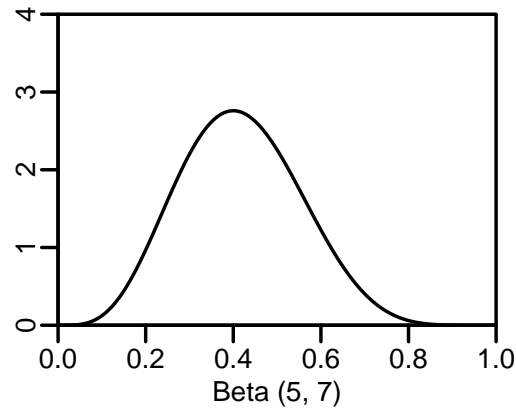
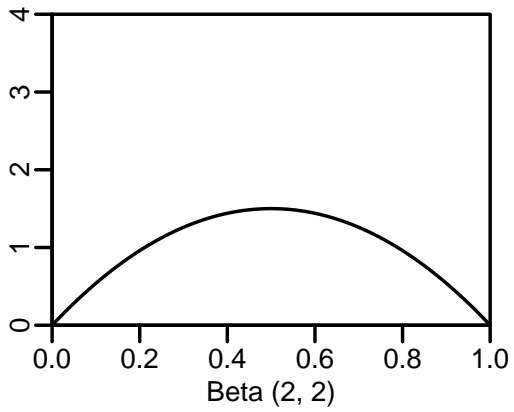
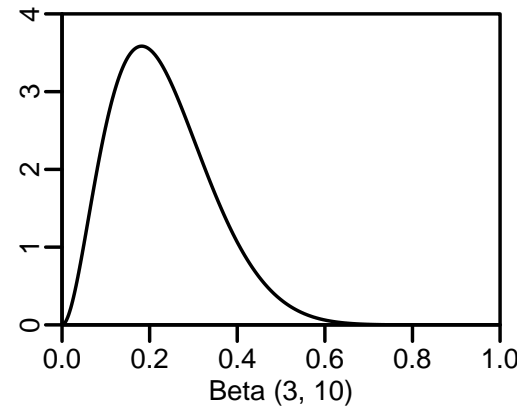
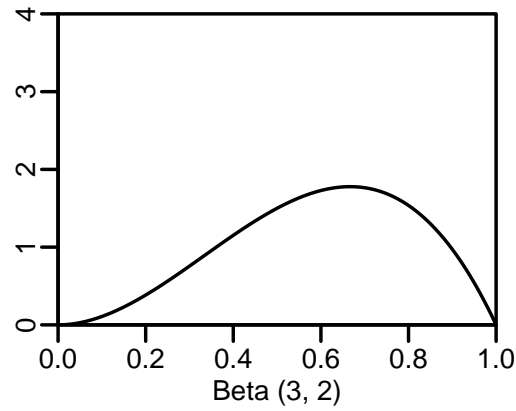
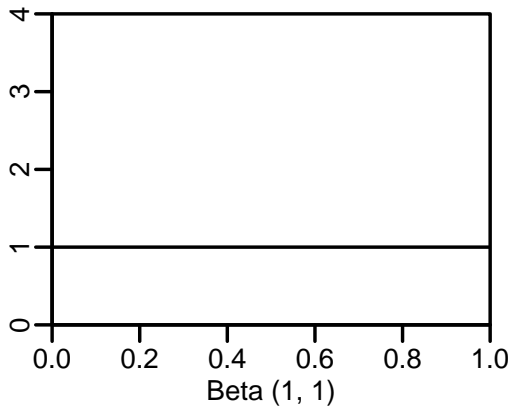
$$\begin{aligned} P(\mu | Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n) \\ &\propto P(\mu) \prod_{i=1}^n P(Y_i = y_i | \mu) \\ &\propto \mu^{a-1} (1-\mu)^{b-1} \prod_{i=1}^n \mu^{y_i} (1-\mu)^{1-y_i} \\ &\propto \mu^{\sum y_i + a - 1} (1-\mu)^{n - \sum y_i + b - 1} \end{aligned}$$

So the posterior distribution is Beta ( $\sum y_i + a, n - \sum y_i + b$ ).

One way this is sometimes visualized is as the prior being equivalent to  $a$  fictitious observations with  $Y = 1$  and  $b$  fictitious observations with  $Y = 0$ .

Note that all that is used from the data is  $\sum y_i$ , which is a *minimal sufficient statistic*, whose values are in one-to-one correspondence with possible likelihood functions (ignoring constant factors).

# Examples of Beta Priors and Posteriors



## Predictive Distribution from Beta Posterior

From the Beta  $(\sum y_i + a, n - \sum y_i + b)$  posterior distribution, we can make a probabilistic prediction for the next observation:

$$\begin{aligned} &P(Y_{n+1} = 1 \mid Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n) \\ &= \int_0^1 P(Y_{n+1} = 1 \mid \mu) P(\mu \mid Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n) d\mu \\ &= \int_0^1 \mu P(\mu \mid Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n) d\mu \\ &= \int_0^1 \mu \frac{\Gamma(n + a + b)}{\Gamma(\sum y_i + a)\Gamma(n - \sum y_i + b)} \mu^{\sum y_i + a - 1} (1 - \mu)^{n - \sum y_i + b - 1} d\mu \\ &= \frac{\Gamma(n + a + b)}{\Gamma(\sum y_i + a)\Gamma(n - \sum y_i + b)} \frac{\Gamma(1 + \sum y_i + a)\Gamma(n - \sum y_i + b)}{\Gamma(1 + n + a + b)} \\ &= \frac{\sum y_i + a}{n + a + b} \end{aligned}$$

This uses the fact that  $c\Gamma(c) = \Gamma(1 + c)$ .

## Generalizing to More Than Two Values

For i.i.d. observations with a finite number,  $K$ , of possible values, with  $K > 2$ , the conjugate prior for the probabilities  $\mu_1, \dots, \mu_K$  is the Dirichlet distribution, with the following density on the simplex where all  $\mu_k > 0$  and  $\sum \mu_k = 1$ :

$$P(\mu_1, \dots, \mu_K) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

The parameters  $\alpha_1, \dots, \alpha_K$  can be any positive reals.

The posterior distribution after observing  $n$  items, with  $m_1$  having value 1,  $m_2$  having value 2, etc. is Dirichlet with parameters  $\alpha_1 + m_1, \dots, \alpha_K + m_K$ .

The predictive distribution for item  $n + 1$  is

$$P(Y_{n+1} = k | Y_1 = y_1, \dots, Y_K = y_k) = \frac{m_k + \alpha_k}{n + \sum \alpha_k}$$

# Independent Observations from a Gaussian Distribution

We observe real variables  $Y_1, Y_2, \dots, Y_n$ .

We model these as being independent, all from some Gaussian distribution with unknown mean,  $\mu$ , and known variance,  $\sigma^2$ .

The conjugate prior for  $\mu$  is Gaussian with some mean  $\mu_0$  and variance  $\sigma_0^2$ .

Rather than talk about the variance, it is more convenient to talk about the *precision*, equal to the reciprocal of the variance. A data point has precision  $\tau = 1/\sigma^2$  and the prior has precision  $\tau_0 = 1/\sigma_0^2$ .

The posterior distribution for  $\mu$  is also Gaussian, with precision  $\tau_n = \tau_0 + n\tau$ , and with mean

$$\mu_n = \frac{\tau_0\mu_0 + n\tau\bar{y}}{\tau_0 + n\tau}$$

where  $\bar{y}$  is the sample mean of the observations  $y_1, \dots, y_n$ .

The predictive distribution for  $Y_{n+1}$  is Gaussian with mean  $\mu_n$  and variance  $(1/\tau_n) + \sigma^2$ .



## Gaussian with Unknown Variance

What if both the mean and the variance (precision) of the Gaussian distribution for  $Y_1, \dots, Y_n$  are unknown?

There is still a conjugate prior, but in it,  $\mu$  and  $\tau$  are dependent:

$$\begin{aligned}\tau &\sim \text{Gamma}(a, b) \\ \mu | \tau &\sim N(\mu_0, c/\tau)\end{aligned}$$

for some constants  $a$ ,  $b$ , and  $c$ .

It's hard to imagine circumstances where our prior information about  $\mu$  and  $\tau$  would have a dependence of this sort. But unfortunately, people use this conjugate prior anyway, because it's convenient.