

# Prior and Posterior

PUBH 8442: Bayes Decision Theory and Data Analysis

Eric F. Lock  
UMN Division of Biostatistics, SPH  
elock@umn.edu

01/27/2025

- ▶ Assume a *sampling model* for data  $\mathbf{y} = (y_1, \dots, y_n)$
- ▶ Specified by parameters  $\theta$ , which may be unknown
- ▶ Often represented as a probability density  $p(\mathbf{y} | \theta)$
- ▶ E.g., Gaussian model specified by  $\theta = (\mu, \sigma^2)$ :

$$p(y_1 | \theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{(y_1 - \mu)^2}{2\sigma^2} \right\}$$

And under independent  $y_i$ 's:

$$p(\mathbf{y} | \theta) = \left( \frac{1}{\sigma\sqrt{2\pi}} \right)^n \exp \left\{ -\frac{\sum_{i=1}^n (y_i - \mu)^2}{2\sigma^2} \right\}$$

- ▶ The probability density  $p(\mathbf{y} | \theta)$  is often called the *likelihood*
  - ▶ Sometimes with the notation  $L(\theta; \mathbf{y})$
- ▶ Use this notation when making inferences about  $\theta$
- ▶ Can choose  $\theta$  to maximize likelihood:

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta; \mathbf{y})$$

- ▶ i.e., estimate  $\theta$  to maximize density of observed data
- ▶ Called *maximum likelihood estimation* (MLE)
- ▶ Has been criticized for overfitting

- ▶ MLE approach assumes  $\theta$  is fixed (though unknown)
- ▶ Alternatively, treat  $\theta$  as a random variable
- ▶ Give  $\theta$  a probability density  $p(\theta)$ 
  - ▶ Potentially specified by *hyperparameters*  $\eta$ :  $p(\theta | \eta)$ .
- ▶ The *marginal* density of  $\mathbf{y}$  is

$$p(\mathbf{y}) = \int p(\mathbf{y} | \theta)p(\theta) d\theta.$$

- ▶ “averaging” over  $\theta$

- ▶ Bayes' rule for continuous random variables:

$$\begin{aligned} p(\theta | \mathbf{y}) &= \frac{p(\mathbf{y}, \theta)}{p(\mathbf{y})} \\ &= \frac{p(\mathbf{y} | \theta)p(\theta)}{\int p(\mathbf{y} | \theta)p(\theta) d\theta}. \end{aligned}$$

- ▶  $p(\theta)$  is the prior,  $p(\theta | \mathbf{y})$  the posterior distribution for  $\theta$
- ▶  $\int p(\mathbf{y} | \theta)p(\theta) d\theta$  is the normalizing constant
  - ▶ Assures the posterior integrates to 1

# Note on notation

- ▶ In class, we will use  $p(\cdot)$  for \*any\* pdf
  - ▶  $p(\mathbf{y} \mid \theta)$ ,  $p(\theta)$ ,  $p(\theta, \mathbf{y})$ ,  $p(\theta \mid \mathbf{y})$ ,  $p(\mathbf{y})$ , etc.
- ▶ For discrete variables, use  $P(\cdot)$  for \*any\* pmf.
- ▶ Common alternative notations:
  - ▶  $\pi$  for prior:  $\pi(\theta)$
  - ▶  $f$  for sampling model / likelihood:  $f(\mathbf{y} \mid \theta)$
  - ▶  $m$  for marginal distribution:  $m(\mathbf{y})$
  - ▶  $p$  for anything else. E.g.,  $p(\theta \mid \mathbf{y})$ ,  $p(\theta, \mathbf{y})$ .

- ▶ If  $\mathbf{y}$  is discrete, may write

$$p(\theta | \mathbf{y} = \mathbf{k}) = \frac{P(\mathbf{y} = \mathbf{k} | \theta)p(\theta)}{\int P(\mathbf{y} = \mathbf{k} | \theta)p(\theta) d\theta}.$$

- ▶ If  $\theta$  is discrete with possible values  $\Theta$ , may write

$$P(\theta = k | \mathbf{y}) = \frac{p(\mathbf{y} | \theta = k)P(\theta = k)}{\sum_{k \in \Theta} p(\mathbf{y} | \theta = k)P(\theta = k)}.$$

## Example: Sex proneness

- ▶ Research suggests the chance of a male or female birth depends on the family.
- ▶ The Gupta's (of CNN fame) have three daughters.
- ▶ What is the probability their next child will be a daughter?
- ▶ A (poor) maximum likelihood solution:
  - ▶ Let  $\theta$  be probability of a daughter for Gupta's
  - ▶ Binomial likelihood for first 3 births, with  $y=\#$  daughters, is

$$P(Y = 3 | \theta) = \theta^3$$

- ▶ Maximized at  $\hat{\theta} = 1$



## Example: Sex proneness

- Alternatively, use prior for  $\theta$
- Naive approach:  $\theta \sim \text{Uniform}(0, 1)$ .
  - So  $p(\theta) = 1$  for  $0 \leq \theta \leq 1$ .
- What is  $p(\theta|y = 3)$ ?
  
  
  
  
  
  
  
  
  
  
- Estimate  $\theta$  using this posterior.

# Beta-Binomial model

- Assume  $\theta \sim \text{Beta}(a, b)$ :

$$p(\theta) = \frac{1}{B(a, b)} \theta^{a-1} (1 - \theta)^{b-1}$$

where  $B(\cdot, \cdot)$  is the *beta function*.

- If  $Y \sim \text{Binomial}(n, \theta)$ , then

$$p(\theta|y = k) = \text{Beta}(a + k, b + n - k)$$

## Example: Sex proneness (cont.)

- ▶ Parental sex proneness is thought to have a very small effect, if any.
- ▶ Based on data from many families, a more realistic prior for  $\theta$  is

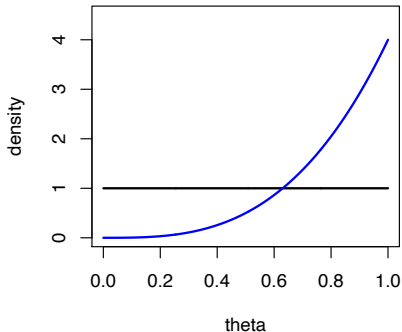
$$p(\theta) = \text{Beta}(39, 40).$$

- ▶ For the Gupta family,  $p(\theta | y = 3) = \text{Beta}(42, 40)$
- ▶ The expected value of  $\text{Beta}(a, b)$  is  $a/(a + b)$ :
  - ▶ Prior estimate is  $E_{\theta} = 39/79 = 0.494$
  - ▶ Posterior estimate is  $E_{\theta | y=3} = 42/82 = 0.512$

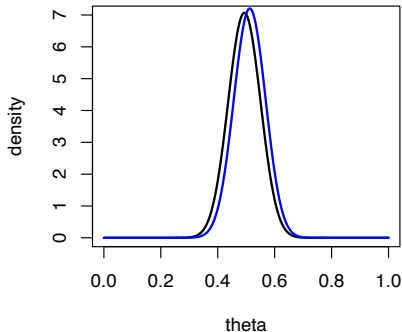
## Example: Sex proneness (cont.)

- **Prior** and **posterior** densities, given  $y = 3$ :

**Beta(1,1) prior**



**Beta(39,40) prior**



Code: [http:](http://www.ericfrazerlock.com/Prior_and_posterior_Rcode1.r)

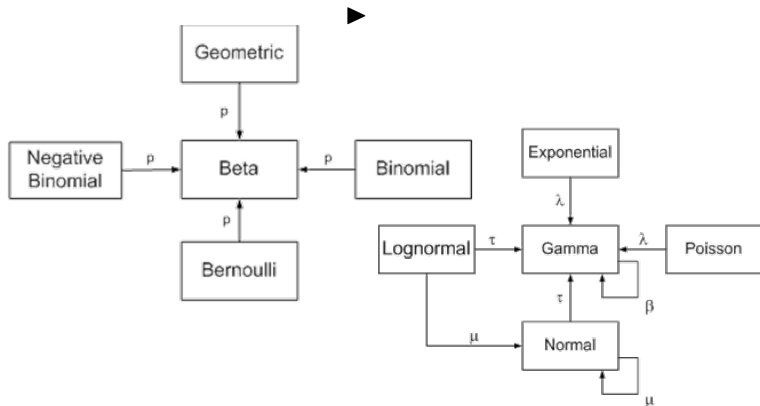
[//www.ericfrazerlock.com/Prior\\_and\\_posterior\\_Rcode1.r](http://www.ericfrazerlock.com/Prior_and_posterior_Rcode1.r)

# Conjugate priors

- ▶ A prior is *conjugate* for a given likelihood if its posterior belongs to the same distributional family.
- ▶ A beta prior is conjugate for binomial data
  - ▶ Both  $p(\theta)$  and  $p(\theta | y)$  give beta distributions.
- ▶ Conjugate priors facilitate computation of posteriors
- ▶ Particularly useful when updating the posterior adaptively
  - ▶ E.g., after one girl, posterior is Beta(40, 40)
  - ▶ After another girl, posterior is Beta (41, 40), etc.
- ▶ Otherwise, no profound theoretical justification.

# Conjugate priors

- ▶ Common conjugate families (credit: John D. Cook)



See [http://en.wikipedia.org/wiki/Conjugate\\_prior](http://en.wikipedia.org/wiki/Conjugate_prior)

# Normal-normal model

- Assume  $y \sim \text{Normal}(\mu, \sigma^2)$  with  $\sigma^2$  known
- If  $p(\mu) = \text{Normal}(\mu_0, \tau^2)$ , then

$$p(\mu | y) = \text{Normal} \left( \frac{\sigma^2 \mu_0 + \tau^2 y}{\sigma^2 + \tau^2}, \frac{\sigma^2 \tau^2}{\sigma^2 + \tau^2} \right)$$

# Normal-normal model

- ▶ Assume  $\mathbf{y} = (y_1, \dots, y_n)$  are iid with  $y_i \sim \text{Normal}(\mu, \sigma^2)$ , and  $\sigma^2$  known.
- ▶ If  $p(\mu) = \text{Normal}(\mu_0, \tau^2)$ , then

$$p(\mu | \mathbf{y}) = \text{Normal} \left( \frac{\sigma^2 \mu_0 + n\tau^2 \bar{y}}{\sigma^2 + n\tau^2}, \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2} \right)$$

where  $\bar{y} = \frac{\sum y_i}{n}$ .

- ▶ Homework.



## Example: Coke bottles

- ▶ Coca-Cola bottling machines fill with known variance 0.05 oz
- ▶ Each machine is calibrated to fill mean capacity  $\mu$
- ▶ Bottles are filled with Gaussian error:  $\text{Normal}(\mu, 0.05)$
- ▶ Historical data show machine calibrations are approximately  $\text{Normal}(12, 0.01)$ .

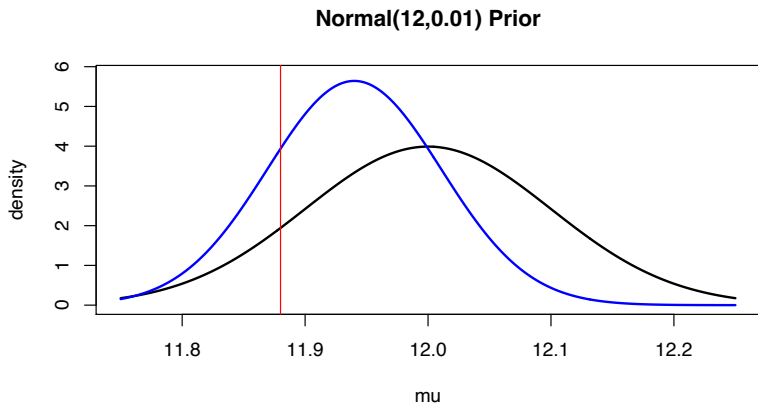
## Example: Coke bottles

- ▶ Five randomly selected bottles from a given machine have sample mean  $\bar{y} = 11.88$
- ▶ What is the posterior for the calibration of this machine?

- ▶  $p(\mu | \mathbf{y}) = \text{Normal}(11.94, 0.005)$

# Example: Coke bottles

- **Prior** and **posterior** densities

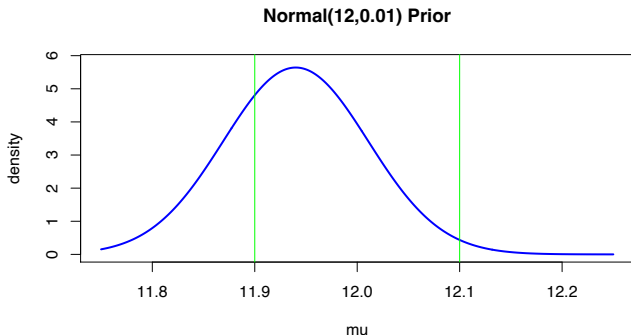


Code: [http:](http://www.ericfrazerlock.com/Prior_and_posterior_Rcode2.r)

[//www.ericfrazerlock.com/Prior\\_and\\_posterior\\_Rcode2.r](http://www.ericfrazerlock.com/Prior_and_posterior_Rcode2.r)

## Example: Coke bottles

- Re-calibrate machines if do not fill within 11.9 and 12.1 oz on average
- What is the probability this machine needs recalibration?



$$\Phi\left(\frac{11.9 - 11.94}{\sqrt{0.005}}\right) + 1 - \Phi\left(\frac{12.1 - 11.94}{\sqrt{0.005}}\right) \approx 0.30$$

# Prior and posterior predictive

- ▶ Often we are not interested in making inference on  $\theta$ , but rather on our best guess for the distribution of  $y_i$
- ▶ Can estimate density of the full model by integrating over  $\theta$
- ▶ The *prior predictive* is the marginal distribution of an observation given the prior:

$$p(y_1) = \int p(y_1 | \theta)p(\theta)d\theta$$

- ▶ The *posterior predictive* is the distribution for a future observation  $y_{n+1}$  given the data so far. If  $y_1, \dots, y_n, y_{n+1} \stackrel{iid}{\sim} p(y | \theta)$ ,

$$p(y_{n+1} | \mathbf{y}) = \int p(y_{n+1} | \theta)p(\theta | \mathbf{y})d\theta.$$

# Normal-normal predictives

- ▶ For the normal-normal model introduced earlier, the prior predictive is

$$p(y_i) = \text{Normal}(\mu_0, \tau^2 + \sigma^2)$$

- ▶ The posterior predictive is

$$p(y_{n+1} | \mathbf{y}) = \text{Normal}\left(\frac{\sigma^2 \mu_0 + n\tau^2 \bar{y}}{\sigma^2 + n\tau^2}, \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2} + \sigma^2\right)$$

- ▶ Homework

## Example: Coke bottles (cont)

- The **predictive for a single bottle from a given machine** is

$$\text{Normal}(12, 0.06)$$

- The **predictive for the sixth bottle after the five observed** is

$$\text{Normal}(11.94, 0.055)$$

Normal(12,0.01) Prior

