

MSA101/MVE187 2018 Lecture 2

Petter Mostad

Chalmers University

September 20, 2018

Example: Learning about a proportion

- ▶ An experiment is performed n times. We assume there is a probability p for "success" each time, and that the outcomes are independent. Let X be the observed number of successes. We get $X \sim \text{Binomial}(n, p)$. Given $X = x$, what do we know about p ?
- ▶ For a Bayesian analysis, we need a joint probability density (or mass function) $\pi(X, p)$. We have defined $\pi(X | p)$ (the *likelihood*). We need to define $\pi(p)$ (the *prior*).
- ▶ Let us first try with the prior $p \sim \text{Uniform}[0, 1]$.
- ▶ The conditional model $\pi(p | X = x)$ (the *posterior* for p) can be computed with Bayes formula. We get

$$\pi(p | X = x) \propto_p p^x (1 - p)^{n-x}.$$

- ▶ We can recognize this as a Beta distribution:
 $p | X = x \sim \text{Beta}(x + 1, n - x + 1)$

Review of definition: The Beta distribution

θ has a Beta distribution on $[0, 1]$, with parameters α and β , if its density has the form

$$\pi(\theta \mid \alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

where $B(\alpha, \beta)$ is the Beta *function* defined by

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

where $\Gamma(t)$ is the *Gamma function* defined by

$$\Gamma(t) = \int_0^{\infty} x^{t-1} e^{-x} dx$$

Recall that for positive integers, $\Gamma(n) = (n - 1)! = 0 \cdot 1 \cdot \dots \cdot (n - 1)$. See for example Wikipedia for more properties of the Beta distribution, and the Beta and Gamma functions. We write $\pi(\theta \mid \alpha, \beta) = \text{Beta}(\theta; \alpha, \beta)$ for the Beta density; we then also write $\theta \sim \text{Beta}(\alpha, \beta)$.

Using a Beta distribution as prior

- ▶ Assume the prior is $p \sim \text{Beta}(\alpha, \beta)$.
- ▶ The posterior becomes

$$p \mid (X = x) \sim \text{Beta}(\alpha + x, \beta + n - x)$$

- ▶ DEFINITION: Given a likelihood model $\pi(x \mid \theta)$. A *conjugate family of priors* to this likelihood is a parametric family of distributions so that if the prior for θ is in this family, the posterior $\theta \mid x$ is also in the family.

Using a discrete prior

- ▶ What if the prior for p is a discrete distribution, i.e., $\pi(p) = \sum_{i=1}^k I(p = p_i)q_i$ where p_1, \dots, p_k are points in the interval $[0, 1]$ and q_1, \dots, q_k are their probabilities?
- ▶ The conditional model is obtained with Bayes theorem:

$$P(p = p_i | x) = \frac{\pi(x | p = p_i)q_i}{\sum_{i=1}^k \pi(x | p = p_i)q_i} = \frac{p_i^x (1 - p_i)^{n-x} q_i}{\sum_{j=1}^k p_j^x (1 - p_j)^{n-x} q_j}.$$

- ▶ Computationally, you compute the vector of likelihoods, multiply termwise with the vector (q_1, \dots, q_k) of prior probabilities, and normalize to 1.

Using discretization

- ▶ Assume you have ANY prior, with density $\pi(p)$ on $[0, 1]$. This density can be approximated, generally with reasonable accuracy, with a discrete distribution, a *discretization*.
- ▶ The corresponding posterior produced by discretization can be easily produced by computer: Compute the likelihood on a grid over p , compute the prior on the same grid, multiply, and normalize.
- ▶ NOTE: This works for ANY likelihood, as long as the parameter p has a prior distribution on a bounded set.

Example: The Poisson-Gamma conjugacy

- ▶ Assume $\pi(x | \theta) = \text{Poisson}(x; \theta)$, i.e., that

$$\pi(x | \theta) = e^{-\theta} \frac{\theta^x}{x!}$$

- ▶ Then $\pi(\theta | \alpha, \beta) = \text{Gamma}(\theta; \alpha, \beta)$ where α, β are positive parameters, is a conjugate family. Recall that

$$\text{Gamma}(\theta; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} \exp(-\beta\theta).$$

- ▶ Specifically, we have the posterior

$$\pi(\theta | x) = \text{Gamma}(\theta; \alpha + x, \beta + 1).$$

- ▶ See Albert Section 3.3 for a computational example.

Example: The Normal-Gamma conjugacy

- ▶ Assume $\pi(x | \tau) = \text{Normal}(x; \mu, 1/\tau)$, so that x is normally distributed with known mean μ and unknown precision τ . The likelihood becomes

$$\pi(x | \tau) = \frac{1}{\sqrt{2\pi 1/\tau}} \exp\left(-\frac{1}{2/\tau} (x - \mu)^2\right) \propto_{\tau} \tau^{1/2} \exp\left(-\frac{1}{2}(x - \mu)^2 \tau\right)$$

- ▶ Then $\pi(\tau | \alpha, \beta) = \text{Gamma}(\tau; \alpha, \beta)$ is a conjugate family, so that

$$\pi(\tau | \alpha, \beta) \propto_{\tau} \tau^{\alpha-1} \exp(-\beta\tau).$$

- ▶ Specifically, we get the posterior below.

$$\pi(\tau | x) = \text{Gamma}\left(\tau; \alpha + \frac{1}{2}, \beta + \frac{1}{2}(x - \mu)^2\right).$$

- ▶ We can also describe this conjugacy using the variance σ^2 and an inverse Gamma (or inverse Chi-squared) distribution.

Example: the Normal-Normal conjugacy

- ▶ Assume $\pi(x | \theta) = \text{Normal}(x; \theta, 1/\tau_0)$, where τ_0 is a known and fixed *precision*.
- ▶ Then $\pi(\theta | \mu, \tau) = \text{Normal}(\theta; \mu, 1/\tau)$, where τ is positive and μ has any real value, is a conjugate family.
- ▶ Specifically, we have the posterior

$$\pi(\theta | x) = \text{Normal} \left(\theta; \frac{\tau_0 x + \tau \mu}{\tau_0 + \tau}, \frac{1}{\tau_0 + \tau} \right)$$

- ▶ PROOF: Use completion of squares.

$$\begin{aligned}
\pi(\theta | x) &\propto_{\theta} \pi(x | \theta)\pi(\theta) \\
&\propto_{\theta} \exp\left(-\frac{\tau_0}{2}(x - \theta)^2\right) \exp\left(-\frac{\tau}{2}(\theta - \mu)^2\right) \\
&= \exp\left(-\frac{1}{2}[\tau_0 x^2 - 2\tau_0 x \theta + \tau_0 \theta^2 + \tau \theta^2 - 2\tau \theta \mu + \tau \mu^2]\right) \\
&\propto_{\theta} \exp\left(-\frac{1}{2}[(\tau_0 + \tau)\theta^2 - 2(\tau_0 x + \tau \mu)\theta]\right) \\
&\propto_{\theta} \exp\left(-\frac{1}{2}(\tau_0 + \tau) \left(\theta - \frac{\tau_0 x + \tau \mu}{\tau_0 + \tau}\right)^2\right) \\
&\propto_{\theta} \text{Normal}\left(\theta; \frac{\tau_0 x + \tau \mu}{\tau_0 + \tau}, \frac{1}{\tau_0 + \tau}\right)
\end{aligned}$$

The Bayesian paradigm implies:

- ▶ The usefulness of a model lies in its ability to predict.
- ▶ We create a joint probability model for the parameters θ , the observed data x , and data we would like to predict x_{new} . Often on the form $\pi(\theta, x, x_{new}) = \pi(\theta)\pi(x | \theta)\pi(x_{new} | \theta)$.
- ▶ The distribution for x_{new} is given by conditioning on the observed x and marginalizing out θ :

$$\begin{aligned}\pi(x_{new} | x) &= \int_{\theta} \pi(\theta, x_{new} | x) d\theta = \int_{\theta} \pi(x_{new} | \theta, x)\pi(\theta | x) d\theta \\ &= \int_{\theta} \pi(x_{new} | \theta)\pi(\theta | x) d\theta\end{aligned}$$

This is called the *posterior predictive distribution*.

- ▶ It is also possible to look at the predictive distribution for x before it has been observed. This is called the *prior predictive distribution*:

$$\pi(x) = \int_{\theta} \pi(x, \theta) d\theta = \int_{\theta} \pi(x | \theta)\pi(\theta) d\theta$$

Predictive distributions when using conjugate priors

- ▶ When using a conjugate prior, not only do we have an analytic expression for the posterior density for θ , we also have analytic expressions for the prior predictive density and the posterior predictive density.
- ▶ To see this for the prior predictive density, use this formula derived from Bayes formula:

$$\pi(x) = \frac{\pi(x | \theta)\pi(\theta)}{\pi(\theta | x)}$$

The prior predictive density is on the left and all expressions on the right have analytic formulas.

- ▶ Note that, when using the right hand side for computing, θ will necessarily eventually disappear.
- ▶ As the posterior predictive distribution is on the same form as the prior predictive, we also get an analytic formula for it. Specifically, we can write

$$\pi(x_{new} | x) = \frac{\pi(x_{new} | \theta)\pi(\theta | x)}{\pi(\theta | x_{new}, x)}.$$

Example: Predictive distribution for the Beta-Binomial conjugacy

- ▶ Assume $\pi(x | \theta) = \text{Binomial}(x; n, \theta)$ and $\pi(\theta) = \text{Beta}(\theta; \alpha, \beta)$.
- ▶ We get for the prior predictive

$$\begin{aligned}\pi(x) &= \frac{\pi(x | \theta)\pi(\theta)}{\pi(\theta | x)} \\ &= \frac{\text{Binomial}(x; n, \theta) \text{Beta}(\theta; \alpha, \beta)}{\text{Beta}(\theta; \alpha + x, \beta + n - x)} \\ &= \frac{\binom{n}{x} \theta^x (1 - \theta)^{n-x} \theta^{\alpha-1} (1 - \theta)^{\beta-1} / \text{B}(\alpha, \beta)}{\theta^{\alpha+x-1} (1 - \theta)^{\beta+n-x-1} / \text{B}(\alpha + x, \beta + n - x)} \\ &= \binom{n}{x} \frac{\text{B}(\alpha + x, \beta + n - x)}{\text{B}(\alpha, \beta)}\end{aligned}$$

- ▶ This is the Beta-Binomial distribution with parameters n , α , and β .

Example: Predictive distribution for the Normal-Normal conjugacy

- ▶ Assume $\pi(x | \theta) = \text{Normal}(x; \theta, 1/\tau_0)$ and $\pi(\theta) = \text{Normal}(\mu, 1/\tau)$.
- ▶ Instead of using the type of computations above, the following is simpler:
 - ▶ We know from general theory of the normal distribution that $\pi(x)$ is normal.
 - ▶ $E(x) = E(E(x | \theta)) = E(\theta) = \mu$.
 - ▶ $\text{Var}(x) = \text{Var}(E(x | \theta)) + E(\text{Var}(x | \theta)) = \text{Var}(\theta) + E(1/\tau_0) = 1/\tau + 1/\tau_0$.
- ▶ So for the prior predictive we get

$$\pi(x) = \text{Normal}(x; \mu; 1/\tau + 1/\tau_0)$$

Mixtures of conjugate distributions

- ▶ Assume we have a model $\pi(x | \theta)$ and a conjugate family of priors with densities $g(\theta; \gamma)$, where $\gamma \in Q$. For a fixed integer $k > 1$ define a new family of prior densities as consisting of all sums

$$\sum_{i=1}^k \alpha_i g(\theta; \gamma_i)$$

where $\alpha_i > 0$, $\sum_{i=1}^k \alpha_i = 1$, and $\gamma_i \in Q$. Then, the new family is also a conjugate family.

- ▶ To assemble a proof: First, write $f_i(x)$ for the prior predictive density when using the prior $g(\theta; \gamma_i)$. We have shown above that it has an analytic form. Also, we know that, when using this prior, the posterior for θ has the form $g(\theta; \gamma'_i)$ for some $\gamma'_i \in Q$. So we can write $\pi(x | \theta)g(\theta; \gamma_i) = f_i(x)g(\theta; \gamma'_i)$.

Mixtures of conjugate distributions, cont.

We can compute the prior predictive as

$$\begin{aligned}\pi(x) &= \int \pi(x | \theta) \left[\sum_{i=1}^k \alpha_i g(\theta; \gamma_i) \right] d\theta \\ &= \sum_{i=1}^k \alpha_i \int \pi(x | \theta) g(\theta; \gamma_i) d\theta = \sum_{i=1}^k \alpha_i f_i(x)\end{aligned}$$

We get the posterior distribution

$$\pi(\theta | x) = \frac{\pi(x | \theta)\pi(\theta)}{\pi(x)} = \sum_{i=1}^k \frac{\alpha_i}{\pi(x)} \pi(x | \theta) g(\theta; \gamma_i) = \sum_{i=1}^k \frac{\alpha_i f_i(x)}{\pi(x)} g(\theta; \gamma_i)$$

Thus the posterior has the same form as the prior: We have conjugacy.

The exponential family of distributions

- ▶ The exponential family of distributions over x with parameters η have densities

$$\pi(x | \eta) = h(x)g(\eta) \exp(\eta \cdot u(x))$$

where η and $u(x)$ are vectors and $\eta \cdot u(x)$ is their dot product.

- ▶ All of the families of distributions we have seen so far, and many more, can be written in this way.
- ▶ Essentially all families that have conjugate prior families are of this type.
- ▶ Read more in Bishop Chapter 2.4