## MSA101/MVE187 2021 Lecture 2

#### Petter Mostad

Chalmers University

September 1, 2021

## Required knowledge

#### in basic probability theory:

- Basic knowledge of distributions, densities, conditional distributions, expectations ...
- Some familiarity with standard distributions such as Binomial, Poisson, Gamma (but no need to memorize).
- Consult your previous statistics/probability textbooks!

#### in clasical statistics:

...not much, you have mostly seen this in the first lecture.

#### in computation:

- We use R. Learn R now!
- ...in fact, no advanced programming is needed to get through this course.

- ▶ Definition and examples of conjugacy. How to compute in practice.
- Predictive distributions when using conjugate families.
- The exponential family of distributions.

- Prediction variable  $Y_{pred}$ , data  $Y_{data}$ , parameter (vector)  $\theta$ .
- Specify a complete model by specifying prior  $\pi(\theta)$ , likelihood  $\pi(Y_{data} \mid \theta)$ , and prediction distribution  $\pi(Y_{pred} \mid \theta)$ .
- Derive the posterior  $\pi(\theta \mid Y_{data})$ .
- Make predictions using

$$\pi(Y_{\textit{pred}} \mid Y_{\textit{data}}) = \int \pi(Y_{\textit{pred}} \mid \theta) \pi(\theta \mid Y_{\textit{data}}) \, d\theta$$

#### Another note on notation

- For standard distributions, we use similar but different notation for a random variable itself, and its density (or probability mass function).
- Example: We write

 $Y \sim \text{Binomial}(N, p)$  and  $\pi(y) = \text{Binomial}(y; N, p)$ 

so we have

$$\mathsf{Binomial}(y; N, p) = \binom{N}{y} p^{y} (1-p)^{N-y}.$$

Example: We write

 $Y \sim \text{Normal}(\mu, \sigma^2)$  and  $\pi(y) = \text{Normal}(y; \mu, \sigma^2)$ 

so we have

Normal
$$(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{1}{2\sigma^2}(y-\mu)^2\right)$$

Sometimes we write, e.g., π(y | μ, σ<sup>2</sup>) = Normal(y; μ, σ<sup>2</sup>) as μ and σ<sup>2</sup> could also be random variables.

- Y<sub>pred</sub> = 1 or 0 (heads or tails). Y<sub>data</sub>: Number of heads in N previous throws. θ: prob. of heads.
- We use  $Y_{data} = y \sim \text{Binomial}(N, \theta)$  and  $Y_{pred} \sim \text{Binomial}(1, \theta)$ .
- We first used a prior with two possible values for θ: 0.7 and 0.3, with equal probabilities.
- We now compute the posterior when the prior is  $\theta \sim \text{Uniform}(0,1)$ .

#### The Beta distribution

 $\theta$  has a Beta distribution on [0, 1], with parameters  $\alpha$  and  $\beta$ , if its density has the form

$$\pi( heta \mid lpha, eta) = rac{1}{\mathsf{B}(lpha, eta)} heta^{lpha - 1} (1 - heta)^{eta - 1}$$

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

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

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,  $\Gamma(n) = (n-1)! = 0 \cdot 1 \cdots (n-1)$ . 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  $\theta \sim \text{Beta}(\alpha, \beta)$ .

- We get from the definition of Beta density that  $\int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta = B(\alpha,\beta).$
- Thus the posterior in our case becomes

$$\pi(\theta \mid y) = \frac{\theta^{y}(1-\theta)^{N-y}}{B(y+1,N-y+1)}$$

We see that

$$\theta \mid y \sim \mathsf{Beta}(y+1, N-y+1)$$

NOTE: Computations can be made simpler, by not keeping track of factors not containing y! We define

expression 1  $\propto_{\theta}$  expression 2

to mean that the second expression is equal to the first expression except for a factor that does not contain the variable  $\theta.$ 

- We say that expression 2 is proportional to expression 1 as a function of θ.
- For example

$$\binom{N}{y} heta^y (1- heta)^{N-y} \propto_{ heta} heta^y (1- heta)^{N-y}$$

Another example:

$$rac{1}{\sqrt{2\pi\sigma^2}}\exp\left(-rac{1}{2\sigma^2}(y-\mu)^2
ight)\propto_\mu \exp\left(-rac{1}{2\sigma^2}(y-\mu)^2
ight).$$

### Using a Beta distribution as prior

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

$$\theta \mid y \sim \mathsf{Beta}(\alpha + y, \beta + N - y)$$

- ▶ DEFINITION: Given a likelihood model  $\pi(y \mid \theta)$ . A conjugate family of priors to this likelihood is a parametric family of distributions so that if the prior for  $\theta$  is in this family, the posterior  $\theta \mid y$  is also in the family.
- So the Beta family is conjugate to the Binomial likelihood: The Beta-Binomial conjugacy.
- NOTE: Uniform(0, 1) = Beta(1, 1), so our previous example is part of this example.

#### Biased coin example, continued





#### Biased coin example, continued



Figure: The probability of heads at each point in a sequence of observations, or the probability of "success", conditioning on the previous observations. The priors used are  $\theta \sim \text{Uniform}(0, 1)$  (left) and  $\theta \sim \text{Beta}(33.4, 33.4)$  (right).

### Example: The Poisson-Gamma conjugacy

• Assume the likelihood is  $\pi(y \mid \theta) = \text{Poisson}(y; \theta)$ , i.e., that

$$\pi(y \mid heta) = e^{- heta} rac{ heta^y}{y!}$$

Then π(θ | α, β) = Gamma(θ; α, β) where α, β are positive parameters, is a conjugate family. Recall that

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

- Compute the posterior!
- Specifically, we get

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

See Albert Section 3.3 for a computational example.

### Example: The Normal-Gamma conjugacy

Assume the likelihood is π(y | τ) = Normal(y; μ, 1/τ), so that y is normally distributed with known mean μ and unknown precision τ. The likelihood becomes

$$\pi(y \mid au) = rac{1}{\sqrt{2\pi 1/ au}} \exp\left(-rac{1}{2/ au} \left(y-\mu
ight)^2
ight) \propto_{ au} au^{1/2} \exp\left(-rac{1}{2} (y-\mu)^2 au
ight)$$

▶ Prove:  $\pi(\tau \mid \alpha, \beta) = \mathsf{Gamma}(\tau; \alpha, \beta)$  is a conjugate family, where

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

Specifically, we get the posterior below:

$$\pi(\tau \mid y) = \mathsf{Gamma}\left( au; \alpha + rac{1}{2}, eta + rac{1}{2}(y-\mu)^2
ight).$$

We can also describe this conjugacy using the variance σ<sup>2</sup> and an inverse Gamma (or inverse Chi-squared) distribution.

- Assume the likelihood is π(y | θ) = Normal(y; θ, 1/τ₀), where τ₀ is a known and fixed precision.
- Then π(θ | μ, τ) = Normal(θ; μ, 1/τ), where τ is positive and μ has any real value, is a conjugate family.
- Specifically, we have the posterior

$$\pi(\theta \mid y) = \mathsf{Normal}\left(\theta; \frac{\tau_0 y + \tau \mu}{\tau_0 + \tau}, \frac{1}{\tau_0 + \tau}\right)$$

PROOF: Use completion of squares.

$$\pi(\theta \mid y) \propto_{\theta} \pi(y \mid \theta)\pi(\theta)$$

$$\propto_{\theta} \exp\left(-\frac{\tau_{0}}{2}(y-\theta)^{2}\right)\exp\left(-\frac{\tau}{2}(\theta-\mu)^{2}\right)$$

$$= \exp\left(-\frac{1}{2}\left[\tau_{0}y^{2}-2\tau_{0}y\theta+\tau_{0}\theta^{2}+\tau\theta^{2}-2\tau\theta\mu+\tau\mu^{2}\right]\right)$$

$$\propto_{\theta} \exp\left(-\frac{1}{2}\left[(\tau_{0}+\tau)\theta^{2}-2(\tau_{0}y+\tau\mu)\theta\right]\right)$$

$$\propto_{\theta} \exp\left(-\frac{1}{2}(\tau_{0}+\tau)\left(\theta-\frac{\tau_{0}y+\tau\mu}{\tau_{0}+\tau}\right)^{2}\right)$$

$$\propto_{\theta} \operatorname{Normal}\left(\theta;\frac{\tau_{0}y+\tau\mu}{\tau_{0}+\tau},\frac{1}{\tau_{0}+\tau}\right)$$

## Conditionally independent data

Assume Y<sub>data</sub> = (y<sub>1</sub>, y<sub>2</sub>), where y<sub>1</sub> and y<sub>2</sub> are conditionally independent given θ, i.e.,

$$\pi(y_1 \mid \theta, y_2) = \pi(y_1 \mid \theta).$$

Then

$$\pi(\theta \mid y_1, y_2) \propto_{\theta} \pi(y_1, y_2 \mid \theta) \pi(\theta) = \pi(y_1 \mid \theta) \pi(y_2 \mid \theta) \pi(\theta)$$

- ► NOTE: We may first find the posterior given y<sub>2</sub>, then use this posterior as the prior when finding the posterior given y<sub>1</sub>: The result will be the posterior given y<sub>1</sub> and y<sub>2</sub>.
- NOTE: We may update the prior on θ sequentially with data y<sub>1</sub>, y<sub>2</sub>,..., y<sub>n</sub>, as long as all the y<sub>i</sub> are conditionally independent given θ.

# Example from Lecture 1: Normal distribution with fixed variance 1

• Assume  $Y_{data} = (y_1, y_2, \dots, y_n)$  where, independently given  $\theta$ ,

 $y_1, y_2, \ldots, y_n \sim \mathsf{Normal}(\theta, 1)$ 

• If the prior is  $heta \sim \mathsf{Normal}(\mu, 1/ au)$ , we get

$$\theta \mid y_1 \sim \mathsf{Normal}\left(rac{y_1 + au\mu}{1 + au}, rac{1}{1 + au}
ight)$$

• Repeated updates give (writing  $\overline{y} = (y_1 + \cdots + y_n)/n$ )

$$\theta \mid y_1, \dots, y_n \sim \mathsf{Normal}\left(\frac{n\overline{y} + \tau\mu}{n + \tau}, \frac{1}{n + \tau}\right)$$

• Similar computations give, for any prior  $\pi(\theta)$ ,

$$\pi(\theta \mid y_1, \ldots, y_n) \propto_{\theta} \mathsf{Normal}\left(\theta; \overline{y}, 1/n\right) \pi(\theta)$$

- An improper distribution is represented by a function of the variable. However it integrates or sums to ∞, and is thus not an actual density or probability mass function.
- Two such functions are regarded as representing the same distribution if they are proportional.
- Extending theory so that such functions can be used as priors turns out to be extremely useful.
- There are no extra problems as long as you make sure the *posterior* density you use for predictions is *proper* (i.e., integrates or sums to 1).

#### Predictive distributions

Recall: We use for prediction

$$\pi(Y_{\textit{pred}} \mid Y_{\textit{data}}) = \int \pi(Y_{\textit{pred}} \mid \theta) \pi(\theta \mid Y_{\textit{data}}) \, d\theta$$

•  $Y_{pred} \mid Y_{data}$  is called the *posterior predictive distribution*.

 In Bayesian statistics, we may even compute the marginal distribution for Y<sub>pred</sub>, for example as

$$\pi(Y_{pred}) = \int \pi(Y_{pred} \mid \theta) \pi(\theta) \, d\theta$$

(as long as we use a proper prior).

- This is called the *prior predictive distribution*.
- It describes your beliefs about Y<sub>pred</sub> before you have looked at any data.

• If the prior is conjugate to  $Y_{pred} \mid \theta$  the prior predictive density is always easy to compute:

$$\pi(Y_{pred}) = rac{\pi(Y_{pred} \mid heta)\pi( heta)}{\pi( heta \mid Y_{pred})}$$

for any  $\theta$ . The densities on the right-hand side can all be computed!

- If the prior is conjugate also to the likelihood Y<sub>data</sub> | θ, the prior and the posterior π(θ | Y<sub>data</sub>) are in the same family, so the posterior will also be conjugate to Y<sub>pred</sub> | θ as above.
- ► Thus, to compute the posterior predictive density, use the same formula as above, just use as prior for  $\theta$  the posterior distribution  $\pi(\theta \mid Y_{data})$ .

# Example: Predictive distribution for the Beta-Binomial conjugacy

- Assume  $\pi(y \mid \theta) = \text{Binomial}(y; N, \theta) \text{ and } \pi(\theta) = \text{Beta}(\theta; \alpha, \beta).$
- We get for the prior predictive

$$\pi(y) = \frac{\pi(y \mid \theta)\pi(\theta)}{\pi(\theta \mid y)} = \frac{\text{Binomial}(y; N, \theta) \text{Beta}(\theta; \alpha, \beta)}{\text{Beta}(\theta; \alpha + y, \beta + N - y)}$$
$$= \frac{\binom{N}{y} \theta^{y} (1 - \theta)^{N - y} \theta^{\alpha - 1} (1 - \theta)^{\beta - 1} / B(\alpha, \beta)}{\theta^{\alpha + y - 1} (1 - \theta)^{\beta + N - y - 1} / B(\alpha + y, \beta + N - y)}$$
$$= \binom{N}{y} \frac{B(\alpha + y, \beta + N - y)}{B(\alpha, \beta)}$$

• This is the Beta-Binomial distribution with parameters N,  $\alpha$ , and  $\beta$ :

 $y \sim \text{Beta-Binomial}(N, \alpha, \beta)$ 

## Predictive distribution for the Poisson-Gamma conjugacy

- ▶ We have seen: If  $y | \theta \sim \text{Poisson}(\theta)$  and  $\theta \sim \text{Gamma}(\alpha, \beta)$  then  $\theta | y \sim \text{Gamma}(\alpha + y, \beta + 1)$ .
- When Y<sub>pred</sub> = y and y ~ Poisson(θ), direct computation gives the prior predictive distribution

$$\pi(y) = \frac{\pi(y \mid \theta)\pi(\theta)}{\pi(\theta \mid y)} = \frac{\beta^{\alpha}\Gamma(\alpha + y)}{(\beta + 1)^{\alpha + y}\Gamma(\alpha)y!}$$

Note that the positive integer x has a Negative-Binomial distribution with parameters r and p if its probability mass function is

$$\pi(x \mid r, p) = \binom{x+r-1}{x} \cdot (1-p)^{x} p^{r} = \frac{\Gamma(x+r)}{\Gamma(x+1)\Gamma(r)} (1-p)^{x} p^{r}$$

- We get that the prior predictive is Negative-Binomial( $\alpha, \beta/(1+\beta)$ ).
- Note that we can get the posterior predictive by simply replacing the α and β of the prior with the corresponding parameters after the update with data.

#### Poisson-Gamma example



Figure: Two different ways of predicting the values of  $k_4$ , given the observations  $k_1 = 20$ ,  $k_2 = 24$ ,  $k_3 = 23$  when  $k_i \mid \theta \sim \text{Poisson}(\theta)$  and an improper Gamma(0,0) prior. The pluses represent the Bayesian predictions using the posterior predictive; the circles represent the Frequentist predictions, using the Poisson distribution with parameter (20 + 24 + 23)/3 = 22.33.

# Example: Predictive distribution for the Normal-Normal conjugacy

- Assume  $\pi(y \mid \theta) = \text{Normal}(y; \theta, 1/\tau_0)$  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 π(x) is normal.
  - $E(y) = E(E(y \mid \theta)) = E(\theta) = \mu$ .
  - $Var(y) = Var(E(y \mid \theta)) + E(Var(y \mid \theta)) = Var(\theta) + E(1/\tau_0) = 1/\tau + 1/\tau_0.$
- So for the prior predictive we get

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

## The exponential family of distributions

Many parametric families of distributions can be written in a particular form:

$$\pi(x \mid \eta) = h(x)g(\eta)\exp\left(\eta \cdot u(x)\right)$$

where  $\eta$  and u(x) are vectors,  $\eta \cdot u(x)$  is their dot product, and  $\eta$  is called the "natural parameters" of the family.

- Some examples of exponential families of distributions, corresponding to particular choices of g, h, and u:
  - Normal distributions.
  - Beta distributions.
  - Poisson distributions.
  - Gamma distributions.
  - Bernoulli distributions and Binomial distributions for a fixed N.
  - Multinomial distributions for a fixed *N*.
  - ....and many more.
- Exponential families of distributions share many properties and can be studied together.

If π(x | η) = h(x)g(η) exp(η ⋅ u(x)), then a conjugate family of priors for η is given as

$$\pi(\eta \mid \nu, \beta) \propto_{\eta} g(\eta)^{\nu} \exp(\eta \cdot \beta).$$

The posterior becomes

$$\pi(\eta \mid x) \propto_{\eta} g(\eta)^{\nu+1} \exp\left(\eta \cdot (\beta + u(x))\right).$$

- Essentially all examples of conjugacy fit into the framework above, so the above describes conjugacy in general.
- Note that the conjugate family of priors is also an exponential family.

#### Some properties

Assume  $\pi(x \mid \eta) = h(x)g(\eta) \exp(\eta \cdot u(x))$ .

The expectation (and further moments) of u(x) can be expressed with a differentiation of g(η):

$$\mathsf{E}_{x|\eta}[u(x)] = -\nabla_{\eta} \log g(\eta).$$

• Given data  $x_1, x_2, \ldots, x_N$  and a prior  $\pi(\eta \mid \nu, \beta) \propto_{\eta} g(\eta)^{\nu} \exp(\eta \cdot \beta)$ the posterior becomes

$$\pi(\eta \mid x_1,\ldots,x_N) \propto_{\eta} g(\eta)^{\nu+N} \exp\left(\eta \cdot \left(\beta + \sum_{i=1}^N u(x_i)\right)\right).$$

- ▶ With for example a flat prior ( $\mu = 0, \beta = 0$ ), the posterior is  $\propto_{\eta} g(\eta)^N \exp\left(\eta \cdot \sum_{i=1}^N u(x_i)\right)$  and
  - The posterior (i.e., likelihood) depends only on  $\sum_i u(x_i)$ .
  - The maximum posterior (i.e., maximum likelihood) is the  $\hat{\eta}$  satisfying

$$-
abla_\eta \log g(\hat{\eta}) = rac{1}{N} \sum_{i=1}^N u(x_i).$$