# MSA101/MVE187 2020 Lecture 3.1 Mixtures

Petter Mostad

Chalmers University

September 15, 2020

- ▶ Basic concepts: Prior, posterior, ...
- Conjugate priors. Computations.
- ► The prior predictive and posterior predictive. Computations.
- ► The exponential family of distributions.

## Mixtures of conjugate priors

- A family of conjugate priors, with limited flexibility, can be greatly extended by also considering linear combinations of these prior densities.
- ► Example: The Poisson-Gamma conjugacy: Assume  $\pi(x \mid \theta) = e^{-\theta} \theta^x / x!$  and  $\pi(\theta) \propto_{\theta} \theta^{\alpha-1} \exp(-\beta\theta)$  so that  $\pi(\theta \mid x) \propto_{\theta} \theta^{\alpha+x-1} \exp(-(\beta+1)\theta)$ . Then a linear combination prior

$$\pi(\theta) = C_1 \theta^{\alpha_1 - 1} \exp(-\beta_1 \theta) + C_2 \theta^{\alpha_2 - 1} \exp(-\beta_2 \theta)$$

will result in a linear combination posterior

$$\pi(\theta \mid x) \propto_{\theta} C_1 \theta^{\alpha_1 + x - 1} \exp(-(\beta_1 + 1)\theta) + C_2 \theta^{\alpha_2 + x - 1} \exp(-(\beta_2 + 1)\theta).$$

- This works for any conjugate family, and any linear combination of priors from it.
- Note however that the weigts of the densities in the linear combination are updated!

#### Mixtures of conjugate priors: Formulas

Assume π(θ | λ) is a family of conjugate priors to π(x | θ). Given λ<sub>1</sub>,..., λ<sub>n</sub>, let g<sub>i</sub>(θ | x) and f<sub>i</sub>(x) denote the posterior and the prior predictive, respectively, when using the prior π(θ | λ<sub>i</sub>). Then

$$\pi(x \mid \theta)\pi(\theta \mid \lambda_i) = g_i(\theta \mid x)f_i(x).$$

Assume we use a linear combination prior

$$\pi( heta) = \sum_{i=1}^n w_i \pi( heta \mid \lambda_i) ext{ where } \sum_{i=1}^n w_i = 1.$$

For the prior predictive we get

$$\pi(x) = \int \pi(x \mid \theta) \sum_{i=1}^{n} w_i \pi(\theta \mid \lambda_i) d\theta = \sum_{i=1}^{n} w_i f_i(x).$$

for the posterior we get

$$\pi(\theta \mid x) = \frac{\pi(x \mid \theta)\pi(\theta)}{\pi(x)} = \frac{\pi(x \mid \theta)\sum_{j=1}^{n} w_j\pi(\theta \mid \lambda_j)}{\sum_{i=1}^{n} w_i f_i(x)}$$
$$= \frac{\sum_{j=1}^{n} w_j f_j(x)g_j(\theta \mid x)}{\sum_{i=1}^{n} w_i f_i(x)} = \sum_{j=1}^{n} w_j'g_j(\theta \mid x) \text{ where } w_j = \frac{w_j f_j(x)}{\sum_{i=1}^{n} w_i f_i(x)}.$$

### Example of mixtures

- ► We use a Binomial likelihood Binomial(4, p), with 3 successes observed in 4 trials.
- We use a mixture prior

$$\pi(p) = 0.5 \cdot \text{Beta}(p; 2.5, 2.5) + 0.5 \cdot \text{Beta}(p; 11, 31)$$

Recall that if x | p ~ Binomial(n, p) and p ~ Beta(α, β) then the prior predictive becomes

$$\pi(x) = \binom{n}{x} \frac{\mathsf{B}(\alpha + x, \beta + n - x)}{\mathsf{B}(\alpha, \beta)}$$

Thus the first updated weight becomes

$$w_1' = \frac{0.5 \cdot \binom{4}{3} \frac{B(2.5+3,2.5+1)}{B(2.5,2.5)}}{0.5 \cdot \binom{4}{3} \frac{B(2.5+3,2.5+1)}{B(2.5,2.5)} + 0.5 \cdot \binom{4}{3} \frac{B(11+3,31+1)}{B(11,31)}}{B(11,31)}$$

and similar for the second updated weight  $w'_2$ .

The posterior becomes

$$\pi(p \mid x) = w'_1 \operatorname{Beta}(x; 2.5 + 3, 2.5 + 1) + w'_2 \operatorname{Beta}(x; 11 + 3, 31 + 1).$$

# MSA101/MVE187 2020 Lecture 3.2 Some multivariate conjugacies

Petter Mostad

Chalmers University

September 3, 2020

# Multivariate conjugacy example: The normal likelihood, no parameters known

 Assume X ~ Normal(μ, 1/τ), with both μ and τ uncertain. The likelihood becomes

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

► Then the Normal-Gamma family is conjugate: The pair (μ, τ) has a Normal-Gamma distribution with parameters μ<sub>0</sub>, λ > 0, α > 0, β > 0 if the density has the form

$$\pi(\mu,\tau \mid \mu_{0},\lambda,\alpha,\beta) = \frac{\beta^{\alpha}\sqrt{\lambda}}{\Gamma(\alpha)\sqrt{2\pi}}\tau^{\alpha-1/2}\exp\left(-\beta\tau - \frac{\lambda\tau}{2}(\mu-\mu_{0})^{2}\right)$$

▶ Note: If  $(\mu, \tau)$  has the Normal-Gamma distribution above, we have  $\tau \sim \text{Gamma}(\alpha, \beta)$  and  $\mu \mid \tau \sim \text{Normal}(\mu_0, 1/(\lambda \tau))$ .

### Example

- Assume  $x = (x_1, x_2, ..., x_n)$  observed from Normal $(\mu, 1/\tau)$ .
- Assume prior

 $au \sim \mathsf{Gamma}(lpha,eta) \;\; \; ext{ and } \;\; \mu \mid au \sim \mathsf{Normal}(\mu_0,1/(\lambda au))$ 

The posterior becomes (compute or use Wikipedia...)

$$\tau \mid x \sim \operatorname{Gamma}\left(\alpha + \frac{n}{2}, \beta + \frac{1}{2}\sum_{i=1}^{n} (x_{i} - \overline{x})^{2} + \frac{n\lambda}{\lambda + n} \frac{(\overline{x} - \mu_{0})^{2}}{2}\right)$$
$$\mu \mid \tau, x \sim \operatorname{Normal}\left(\frac{\lambda\mu_{0} + n\overline{x}}{\lambda + n}, \frac{1}{(\lambda + n)\tau}\right)$$

▶ Numerical example: *x* = (3.1, 4.2, 2.9, 3.7, 3.9). Plot:

- Likelihood
- Prior

$$au \sim \mathsf{Gamma}(1,1/3)$$
 and  $\mu \mid au \sim \mathsf{Normal}(3,3/ au)$ 

Posterior

 $au \sim \mathsf{Gamma}(3.5, 0.9783)$  and  $\mu \mid au \sim \mathsf{Normal}(3.525, 1/(5.33 \cdot au))$ 

# MSA101/MVE187 2020 Lecture 3.3 Some multivariate conjugacies

Petter Mostad

Chalmers University

September 3, 2020

- Where does the prior come from?
- From "other information" than the current data; possibly older data. The context.
- What if the context indicates "no prior information"?
- ► For many likelihood models there exist "non-informative priors".
- ► For the Normal-Gamma prior to the Normal distribution, try the prior  $\alpha = -1/2, \beta = \lambda = \mu_0 = 0.$
- ▶ NOTE: The posterior now becomes  $\tau \sim \text{Gamma}(\frac{n-1}{2}, \frac{1}{2}\sum_{i=1}^{n}(x_i \overline{x})^2)$  and  $\mu \mid \tau \sim \text{Normal}(\overline{x}, \frac{1}{n\tau})$ , corresponding to standard frequentist analysis.
- ▶ NOTE: Plugging in the suggested parameters into the prior density gives  $\pi(\mu, \tau) \propto_{\mu, \tau} \tau^{-1}$ . The prior is an *improper prior*. It does not integrate to 1!

#### Stepwise Bayesian updating

• Assume  $x = (x_1, \ldots, x_m)$  is a random sample, so that

$$\pi(x \mid heta) = \prod_{i=1}^m \pi(x_i \mid heta)$$

• Using a prior  $\pi(\theta)$  the posterior becomes

$$\pi( heta \mid x) \propto_{ heta} \prod_{i=1}^m \pi(x_i \mid heta) \pi( heta)$$

▶ If we first update only with the observations *x*<sub>1</sub>,...,*x*<sub>k</sub>, we get the posterior

$$\pi( heta \mid x_1, \dots, x_k) \propto_{ heta} \prod_{i=1}^k \pi(x_i \mid heta) \pi( heta)$$

- We see that if we use this as the prior and update with the remaining data (x<sub>k+1</sub>,..., x<sub>m</sub>), we get the same result as before.
- In Bayesian statistics, we may subdivide the data into data subsets and update the model stepwise with the data, as long as all the data sets are mutually independent given the model parameter θ.

## Multinomial-Dirichlet conjugacy

Assume x = (x<sub>1</sub>,...,x<sub>n</sub>) ~ Multinomial(m, θ<sub>1</sub>, θ<sub>2</sub>,..., θ<sub>n</sub>), with θ<sub>1</sub> + ··· + θ<sub>n</sub> = 1, so that x<sub>i</sub> counts the number of results of type i in m independent trials, if results of type i have probability θ<sub>i</sub>. The probability mass function is

$$\pi(x \mid \theta_1, \ldots, \theta_n) = \frac{m!}{x_1! \ldots x_k!} \theta_1^{x_1} \ldots \theta_n^{x_n}$$

► (θ<sub>1</sub>,...,θ<sub>n</sub>) has a Dirichlet distribution with parameters α<sub>1</sub>,...,α<sub>n</sub> if the density can be written as

$$\pi(\theta_1,\ldots,\theta_n \mid \alpha_1,\ldots,\alpha_n) = \frac{\Gamma(\alpha_1+\cdots+\alpha_n)}{\Gamma(\alpha_1)\ldots\Gamma(\alpha_n)} \theta_1^{\alpha_1-1}\ldots\theta_n^{\alpha_n-1}$$

Prove that the Dirichlet family is a conjugate family to the Multinomial likelhiood!

## Applied example: Forensic DNA matches

- DNA matching between a trace and a person may be used as proof in criminal cases.
- Given that there is a match, how strong is the evidence? Could the match have been by chance?
- DNA traces collect information from (mostly independent) *loci* (i.e., locations) in the genome. At these loci, each person has two copies of an *allele*, various alleles may exist in the population.
- Popular forensic loci are STRs, where alleles differ by the number of repetitions of a short sequence (such as, e.g., CGTT).
- To answer question above, one needs to establish the probaility of observing in a random person each possible allele. For this databases are collected.
- Lets say a (toy) database of size 16 contans the counts: 14: 3; 15: 9, 16: 3, 17: 0, 18: 1. What is the population probabilities of each allele? How certain can you be?
- It is then common, and advantageous, to use the Multinomial-Dirichlet model, with *pseudocounts*.

MSA101/MVE187 2020 Lecture 3.4 Total expectation, total variance The multivariate normal distribution

Petter Mostad

Chalmers University

September 4, 2020

Assume you have a joint density  $\pi(x, \theta)$ . Note the formulas

 $\mathsf{E}(x) = \mathsf{E}(\mathsf{E}(x \mid \theta))$ 

and

$$Var(x) = E(Var(x \mid \theta)) + Var(E(x \mid \theta)).$$

You should be able to use and prove these formulas.

### The multivariate normal distribution

We say X has a multivariate (n-variate) normal distribution, if it is a vector of length n with density

$$\pi(X) = rac{1}{|2\pi\Sigma|^{1/2}} \exp\left(-rac{1}{2}(X-\mu)\Sigma^{-1}(X-\mu)^t
ight)$$

where the vector  $\mu$  is the expectation and the  $n \times n$  symmetric matrix  $\Sigma$  is the covariance matrix.  $|2\pi\Sigma|$  is the determinant of  $2\pi\Sigma$ .

• We write 
$$X \sim \text{Normal}(\mu, \Sigma)$$
.

- Just as in the 1-dimensional case: If Y | X ~ Normal(AX + B, Σ<sub>1</sub>) and X ~ Normal(μ, Σ<sub>0</sub>), and if we look at Y | X as a likelihood and π(X) as a prior, then this is a conjugate prior.
- We usually express this by using that
  - ▶ In the case above, the *joint* density for X and Y is multivariate normal.
  - For a multivariate normal vector, the *conditional* vector when fixing one or more components in the vector is also multivariate normal.

## The joint multivariate normal distribution

 Assume Y | X ~ Normal(AX + B, Σ<sub>1</sub>) and X ~ Normal(μ, Σ<sub>0</sub>). Then

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim \mathsf{Normal} \left( \begin{bmatrix} \mu \\ A\mu + B \end{bmatrix}, \begin{bmatrix} \Sigma_0 & \Sigma_0 A^t \\ A\Sigma_0 & A\Sigma_0 A^t + \Sigma_1 \end{bmatrix} \right)$$

One can prove this directly from the definitions, or use

- Prove first that the joint distribution must be multivariate normal.
- Then, compute the expectation and the covariance matrix of the joint vector, using, e.g., the formulas for total expectation and variation, or matrix algebra.

# The conditional and the marginal in a multivariate normal distribution

Assume the joint distribution for two vectors X and Y is multivariate normal. Then

- If we integrate out one of them, e.g. Y, the marginal for X, the remaining one, is multivariate normal. The parameters can be read off the expectation and the covariance matrix of the joint distribution.
- ► If we fix Y, then the conditional distribution X | Y is also multivariate normal. In fact, if

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim \text{Normal} \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix}^{-1} \right)$$

we have

$$X \mid Y \sim \text{Normal}(\mu_1 - P_{11}^{-1}P_{12}(Y - \mu_2), P_{11}^{-1})$$

Prove the algebraic matrix identity

$$\begin{pmatrix} \left[ \begin{array}{c} \theta_1 \\ \theta_2 \end{array} \right] - \left[ \begin{array}{c} \mu_1 \\ \mu_2 \end{array} \right] \end{pmatrix}^t \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix} \begin{pmatrix} \left[ \begin{array}{c} \theta_1 \\ \theta_2 \end{array} \right] - \left[ \begin{array}{c} \mu_1 \\ \mu_2 \end{bmatrix} \end{pmatrix} \\ = & \left( \theta_1 - \mu_1 + P_{11}^{-1} P_{12} (\theta_2 - \mu_2) \right)^t P_{11} \left( \theta_1 - \mu_1 + P_{11}^{-1} P_{12} (\theta_2 - \mu_2) \right) \\ + \left( \theta_2 - \mu_2 \right)^t (P_{22} - P_{21} P_{11}^{-1} P_{12}) (\theta_2 - \mu_2).$$

► Use the definition of the joint density for X and Y, and rewrite it as two factors, one depending only on Y.