MSA101/MVE187 2021 Lecture 3 Low-dimensional Bayesian inference Mixtures Some multivariate conjugacies

Petter Mostad

Chalmers University

September 5, 2022

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

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

- Last time: Both likelihood and prior are from a list of elementary distributions, and are conjugate.
- Extension: Use discretization and computers: Works in low dimensions.
- Small extension: Use *mixtures* of priors.
- Small extension: Use multivariate conjugacies.
- ▶ Next time: Huge extension: Use simulation.

Bayesian inference with a discrete parameter θ

Assume θ has possible values $\theta_1, \ldots, \theta_n$.

• The prior $\pi(\theta)$ is represented as a vector $\mathbf{v} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$:

$$\mathbf{v}_i = \pi(\theta_i).$$

• The likelihood $\pi(y \mid \theta)$ is represented as a vector $w = (w_1, \dots, w_n)$:

$$w_i = \pi(y \mid \theta_i).$$

• The posterior is represented as a vector $z = (z_1, \ldots, z_n)$:

$$z_i = \frac{v_i \cdot w_i}{\sum_{j=1}^n v_j \cdot w_j}$$

The posterior predictive distribution can be computed for all values of Y_{pred} as a sum:

$$\pi(Y_{\mathsf{pred}} \mid Y_{\mathsf{data}}) = \sum_{i=1}^n \pi(Y_{\mathsf{pred}} \mid heta_i) z_i.$$

An experimental production process for an electronic component produces faulty components at a rate θ ; 17 tests have produced 2 faulty components; you want to predict probability of at most 1 faulty component in the next batch of 10.

- Prior (constructed based on earlier experience)
- Prior Likelihood





Prediction $\sum_{\theta} (\text{Binomial}(0; 10, \theta) + \text{Binomial}(1; 10, \theta)) \pi(\theta \mid \text{data}) = 0.4642503$

Posterio

Example: Braking distance for a bike, depending on speed



Braking distance for a bike has been measured at 5 different speeds: Data is $(x_1, y_1), \ldots, (x_5, y_5)$. At speed 30, what is ³ the probability that breaking distance will be more than 5?



Model: We assume y_i | a, b ~ Normal(ax_i + bx_i², 0.8²), and use a discrete prior on a grid for parameters (a, b).



Example: Braking distance for a bike

• Likelihood: $\pi(\text{data} \mid (a, b)) = \prod_{i=1}^{5} \text{Normal}(y_i; ax_i + bx_i^2, 0.8^2)$

Likelihood and posterior computed:



Prediction:

$$Pr(y > 5 \mid x = 30, data)$$

$$= \sum_{a,b} \left(\int_{5}^{\infty} Normal(y; a30 + b30^{2}, 0.8^{2}) \, dy \right) \pi(a, b \mid data)$$

$$= 0.9396133$$

- Two main strategies: Aiming for non-informative or informative priors.
- Non-informative examples:
 - With conjugacies, using improper distributions like Gamma(0,0) or Beta(0,0)
 - "Flat" densities...(but depends on scale!)
- Informative:
 - Use posteriors based on previous data, or
 - Check out the prior predictive: Does it "look reasonable" compared to what you expect for such data?

What happens with the discretization if θ is a high-dimensional variable?

In practice, we have to find other methods than discretization.

Numerical integration

The integrals of Bayesian inference

$$\pi(heta \mid Y_{\mathsf{data}}) = rac{\pi(Y_{\mathsf{data}} \mid heta)\pi(heta)}{\int_{ heta} \pi(Y_{\mathsf{data}} \mid heta)\pi(heta) \, d heta}$$

and

$$\pi(Y_{\text{pred}} \mid Y_{\text{data}}) = \int_{\theta} \pi(Y_{\text{pred}} \mid \theta) \pi(\theta \mid Y_{\text{data}}) d\theta$$
$$= \frac{\int_{\theta} \pi(Y_{\text{pred}} \mid \theta) \pi(Y_{\text{data}} \mid \theta) \pi(\theta) d\theta}{\int_{\theta} \pi(Y_{\text{data}} \mid \theta) \pi(\theta) d\theta}$$

can be computed with numerical integration.

- Can work slightly better than discretization (after all discretization is a primitive form of numerical integration).
- Suffers from the same curse of dimensionality as discretization.

Mixtures

A density written as a linear combination of other densities is called a mixture (where $\sum_{i=1}^{n} \nu_i = 1$):

$$\pi(\theta) = \sum_{i=1}^n \nu_i \pi_i(\theta).$$

Using a mixture prior gives a mixture prior predictive distribution:

$$\pi(\mathbf{y}) = \int \pi(\mathbf{y} \mid \theta) \sum_{i=1}^{n} \nu_i \pi_i(\theta) \, d\theta = \sum_{i=1}^{n} \nu_i \int \pi(\mathbf{y} \mid \theta) \pi_i(\theta) \, d\theta = \sum_{i=1}^{n} \nu_i \pi_i(\mathbf{y}).$$

• Defining
$$\pi_i(\theta \mid y) = \frac{\pi(y|\theta)\pi_i(\theta)}{\pi_i(y)}$$
, we also get a mixture posterior:

$$\pi(\theta \mid y) = \frac{\pi(y \mid \theta)\pi(\theta)}{\pi(y)} = \frac{\sum_{i=1}^n \nu_i \pi(y \mid \theta)\pi_i(\theta)}{\sum_{j=1}^n \nu_j \pi_j(y)} = \frac{\sum_{i=1}^n \nu_i \pi_i(y)\pi_i(\theta \mid y)}{\sum_{j=1}^n \nu_j \pi_j(y)}$$

$$= \sum_{i=1}^n \left(\frac{\nu_i \pi_i(y)}{\sum_{j=1}^n \mu_j \pi_j(y)}\right) \pi_i(\theta \mid y).$$

Finally, a mixture posterior predictive distribution:

$$\pi(y_{\mathsf{pred}} \mid y) = \int \pi(y_{\mathsf{pred}} \mid \theta) \pi(\theta \mid y) \, d\theta = \sum_{i=1}^n \left(\frac{\nu_i \pi_i(y)}{\sum_{j=1}^n \mu_j \pi_j(y)} \right) \pi_i(y_{\mathsf{pred}} \mid y).$$

Example: More on braking bikes



We now use the following pixture of priors:

We get the updated weights

 $\frac{0.95 \cdot \pi_1(\mathsf{data})}{0.95 \cdot \pi_1(\mathsf{data}) + 0.05 \cdot \pi_2(\mathsf{data})} = 0.8530929$

for Prior 1 and 1 - 0.8530929 = 0.1469071 for Prior 2.

 Using combined prior, the prediction hardly changes: 0.9399131 (before it was 0.9396133).

- When all the priors π₁(θ),..., π_n(θ) are conjugate to the likelihood π(data | θ), the mixture is also conjugate!
- Is a very powerful way to make families of conjugate priors more flexible!
- Note that we have formulas for all the weights and the posteriors occurring, no integration necessary.

Example: Using mixtures

If $y \sim \text{Geometric}(p)$ with $0 then <math>\pi(y \mid p) = p(1-p)^y$, and $p \sim \text{Beta}(\alpha, \beta)$ is a conjugate family.

If in some applied context our prior information is represented by, e.g., a histogram, we can model it as a Beta mixture:



In this simple case, you could alternatively use discretization.

Multivariate conjugacy example: The normal likelihood, no parameters known

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

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

Then the Normal-Gamma family is conjugate: The pair (μ, τ) has a Normal-Gamma distribution with parameters μ₀, 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 (μ, τ) has the Normal-Gamma distribution above, we have $\tau \sim \text{Gamma}(\alpha, \beta)$ and $\mu \mid \tau \sim \text{Normal}(\mu_0, 1/(\lambda \tau))$.

Computing the posterior

Assume x = (x₁, x₂,..., x_n) sampled from Normal(μ, 1/τ).
 Assume prior

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

Computing the posterior density using our proportionality method, the result is a Normal-Gamma density which can be expressed as

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

- Computations like these can get hairy; if you are lazy like me, consult, e.g., Wikipedia.
- ▶ Using improper prior $\pi(\mu, \tau) \propto_{\mu, \tau} 1/\tau$ gives posterior $\tau \mid x \sim \text{Gamma}(\frac{n-1}{2}, \frac{1}{2}\sum_{i=1}^{n}(x_i \overline{x})^2)$ and $\mu \mid \tau, x \sim \text{Normal}(\overline{x}, \frac{1}{n\tau})$.
- NOTE: The expectation of the posterior for τ then becomes 1 divided by the classical variance estimator, and the expectation for μ becomes x̄.

Predictive distributions

Given parameters ν > 0, μ, and σ², a real variable x has a generalized t-distribution, x ~ t(ν, μ, σ²), when the density is

$$t(x;\nu,\mu,\sigma^{2}) = \frac{1}{\sqrt{\nu\sigma^{2}}B(\nu/2,1/2)} \left[1 + \frac{1}{\nu} \left(\frac{x-\mu}{\sigma}\right)^{2}\right]^{-\frac{\nu+1}{2}}$$

When x | τ ~ Normal(μ, 1/λτ) and τ ~ Gamma(α, β), the marginal (i.e. prior predictive) becomes

$$\pi(\mathbf{x}) = \mathsf{t}\left(\mathbf{x}; 2\alpha, \mu, \frac{\beta}{\alpha\lambda}\right)$$

• When $x \mid \mu, \tau \sim \text{Normal}(\mu, 1/\tau)$, $\mu \mid \tau \sim \text{Normal}(\mu_0, \frac{1}{\lambda\tau})$, and $\tau \sim \text{Gamma}(\alpha, \beta)$, then the marginal becomes

$$\pi(x) = t\left(x; 2\alpha, \mu_0, \frac{\beta(\lambda+1)}{\alpha\lambda}\right)$$

To derive this, marginalize first over the normal-normal conjugacy.

Multinomial-Dirichlet conjugacy

Assume x = (x₁,...,x_n) ~ Multinomial(m, θ₁, θ₂,..., θ_n), with θ₁ + ··· + θ_n = 1, so that x_i counts the number of results of type i in m independent trials, if results of type i have probability θ_i. The probability mass function is

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

• $\theta = (\theta_1, \dots, \theta_n)$ with $\theta_i > 0$ and $\sum_{i=1}^n \theta_i = 1$ has a Dirichlet distribution with parameters $\alpha_1, \dots, \alpha_n$ 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!
- With a Dirichlet(α₁,..., α_n) prior, one can show that the probability of observing a type *i* result in the next trial becomes

$$\frac{\alpha_i+x_i}{\sum_{j=1}^n(\alpha_j+x_j)}.$$

Applied example: Forensic DNA matches

- DNA matching between a trace and a person may be used as proof in criminal cases: For this, one needs to compute the strength of evidence when there is a match at some investigated *loci*.
- At an STR locus in a chromosome, a person has a particular allele (variant): Variants there differ by the number of repetitions of a short sequence (such as CAAT).
- The probability that a random person has a particular allele at this chromosome needs to be computed.
- To do so, population databases of alleles are collected. A small database might look like

10	11	12	13	14	15	16	17	18
1	0	5	89	143	9	3	0	2

- What is the probability that a random person has 17 repetitions as his allele?
- lt is common to use the Multinomial-Dirichlet model together with *pseudocounts*, i.e., values for α_i , for example $\alpha_i = 0.5$ or $\alpha_i = 1$.
- Probabilities get a reasonable value, instead of zero.

The multivariate normal distribution

We say X has a multivariate (n-variate) normal distribution, if it is a real 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 μ is the expectation and the $n \times n$ symmetric matrix Σ 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, Σ₁) and X ~ Normal(μ, Σ₀), 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.

Assume Y | X ~ Normal(AX + B, Σ₁) and X ~ Normal(μ, Σ₀). 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 θ_1 and θ_2 is multivariate normal. Then

- If we integrate out one of them, e.g. θ₂, the marginal for θ₁ is multivariate normal. The parameters can be read off the expectation and the covariance matrix of the joint distribution.
- If we fix θ₂, then the *conditional distribution* θ₁ | θ₂ is also multivariate normal. In fact, if

$$\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \sim \mathsf{Normal} \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix}^{-1} \right)$$

we have

$$\theta_1 \mid \theta_2 \sim \mathsf{Normal}(\mu_1 - P_{11}^{-1}P_{12}(Y - \mu_2), P_{11}^{-1})$$

Prove the algebraic matrix identity

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

Use the definition of the joint density for θ₁ and θ₂, and rewrite it as two factors, one depending only on θ₂.