

1. Consider the following linear regression model:

$$\sigma^2 \sim \text{InvGam}\left(\frac{\alpha_0}{2}, \frac{\delta_0}{2}\right) \quad (1)$$

$$\beta \sim \mathcal{N}(\beta_0, B_0) \quad (2)$$

$$Y | \beta, \sigma^2 \sim \mathcal{N}(X\beta, \sigma^2 I_T) \quad (3)$$

with $Y = T \times 1$ vector and $X = T \times k$ matrix. Then, derive the full conditional distribution of β and σ^2 .

Solution: It is well-known that the priors in this model are conditionally conjugate, which means the full conditionals are of the same parametric families as their prior distributions. Therefore, we can compute them as follows.

$$L(\beta, \sigma^2 | Y) \pi(\beta) \pi(\sigma^2) \propto \mathcal{N}(Y | X\beta, \sigma^2 I_T) \cdot \mathcal{N}(\beta | \beta_0, B_0) \cdot \text{InvGam}\left(\frac{\alpha_0}{2}, \frac{\delta_0}{2}\right) \quad (4)$$

$$\propto (\sigma^2)^T \exp\left(-\frac{1}{2\sigma^2} (y - X\beta)' (y - X\beta)\right) \times \exp\left(-\frac{1}{2} (\beta - \beta_0)' B_0^{-1} (\beta - \beta_0)\right) \quad (5)$$

$$\times (\sigma^2)^{-(\alpha_0/2+1)} \exp\left(-\frac{\delta_0}{2} \frac{1}{\sigma^2}\right) \quad (6)$$

Therefore, ignoring all the terms that don't have β , we get

$$\pi(\beta | Y, \sigma^2) \propto \exp\left(-\frac{1}{2} \left(\beta' \left(\frac{1}{\sigma^2} X'X + B_0^{-1}\right) \beta - 2 \left(\frac{1}{\sigma^2} X'y + B_0^{-1} \beta_0\right)' \beta\right)\right) \quad (7)$$

This is exactly the kernel of a normal distribution. Thus, $\beta | Y, \sigma^2 \sim \mathcal{N}(\beta_T, \Sigma_T)$ where

$$\Sigma_T = \left(\frac{1}{\sigma^2} X'X + B_0^{-1}\right)^{-1} \quad (8)$$

$$\beta_T = \Sigma_T \left(\frac{1}{\sigma^2} X'y + B_0^{-1} \beta_0\right) \quad (9)$$

Likewise,

$$\pi(\sigma^2 | Y, \beta) \propto (\sigma^2)^{-((T+\alpha_0)/2+1)} \exp\left(-\frac{1}{2\sigma^2} (\delta_0 + (y - X\beta)' (y - X\beta))\right) \quad (10)$$

Therefore,

$$\sigma^2 | Y, \beta \sim \text{InvGam}\left(\frac{\alpha_0 + T}{2}, \frac{1}{2} (\delta_0 + (y - X\beta)' (y - X\beta))\right) \quad (11)$$

2. (Greenberg, chapter 2) Consider the uniform distribution with density function $f(y_i | \theta) = \theta^{-1}, 0 \leq y_i \leq \theta$, and θ is unknown.

- (a) Show that the Pareto distribution, $\text{Pareto}(a, k)$

$$\pi(\theta) = \begin{cases} ak^a \theta^{-(a+1)}, & \text{if } \theta \geq k \text{ and } a > 0 \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

is a conjugate prior distribution for the uniform distribution assuming that $a > 1$. Hint: If $a > 1$, then $\mathbf{E}(\theta) = ak/(a - 1)$.

- (b) Show that $\hat{\theta} = \max(y_1, y_2, \dots, y_n)$ is the MLE of θ , where the y_i is the random variable from $f(y_i | \theta)$.
- (c) Find the posterior distribution of θ and the expected value.

Solution:

- (a) Conjugacy implies that the prior and posterior distributions are in the same parametric distributional family. Therefore, we simply need to show that the posterior distribution is also a Pareto distribution.

$$L(\theta | y_1, \dots, y_n) \propto \pi(\theta) \prod_{i=1}^n f(y_i | \theta) \mathbf{I}_{(y_{(n)}, \infty)}(\theta) \cdot \mathbf{I}_{(k, \infty)}(\theta) \quad (13)$$

$$\propto \theta^{-(\alpha+1)} \cdot \theta^{-1} \cdot \mathbf{I}_{(y_{(n)} \vee k, \infty)}(\theta) \quad (14)$$

$$\propto \theta^{-(\alpha+2)} \mathbf{I}_{(y_{(n)} \vee k, \infty)}(\theta) \quad (15)$$

Therefore,

$$\theta | y_1, \dots, y_n \sim \text{Pareto}(\alpha + 1, y_{(n)} \vee k) \quad (16)$$

where \vee denotes the maximum operator and $y_{(n)}$ is the maximum order statistic.

- (b) We take the logarithm of the likelihood and use the first-order condition to find the maximum.

$$\ell(\theta | \{y_i\}_{i=1}^n) \propto -n \log \theta \cdot \mathbf{I}_{(y_{(n)}, \infty)}(\theta) \quad (17)$$

$$\frac{d}{d\theta} \ell(\theta) = -\frac{n}{\theta} \quad (18)$$

Because the first-derivative is never zero, we need only use the monotone property of the likelihood function. Since it is monotonically decreasing, the maximum value is obtained when

θ is smallest, i.e., $\theta \in (y_{(n)}, \infty)$. Therefore,

$$\hat{\theta}^{\text{MLE}} = \inf \{ \theta \mid y_{(n)} < \theta \} = y_{(n)} \quad (19)$$

(c) We have already found the posterior of θ : $\text{Pareto}(\alpha + 1, y_{(n)} \vee k)$. If we let $w = y_{(n)} \vee k$, by definition, the expected value is

$$\mathbf{E}(\theta \mid \{y_i\}_{i=1}^n) = \int_w^\infty (\alpha + 1) w^{\alpha+1} \theta^{-(\alpha+1)} d\theta \quad (20)$$

$$= \frac{\alpha + 1}{\alpha} w, \quad \text{where } w > 0 \quad (21)$$

3. (Greenberg, chapter 2) The density function of the exponential distribution is

$$f(y_i \mid \theta) = \theta e^{-\theta y_i}, \quad \theta > 0, y_i > 0, \quad (22)$$

and let y_1, y_2, \dots, y_n be a random sample from the distribution.

(a) Show that the gamma distribution $\text{Ga}(\alpha, \beta)$ is a conjugate prior distribution for the exponential distribution. Hint: The density of $\text{Ga}(\alpha, \beta)$ is

$$\pi(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta} \quad (23)$$

(b) Show that $1/\bar{y}$ is the MLE for θ where \bar{y} is the sample mean of the observations.

(c) Write the mean of the posterior distribution as a weighted average of the mean of the prior distribution and the MLE.

(d) What happens to the weight on the prior mean as n becomes large?

Solution:

(a) Again, we need to show that the posterior is again a gamma distribution.

$$\pi(\theta \mid y_1, \dots, y_n) \propto \pi(\theta) \prod_{i=1}^n f(y_i \mid \theta) \quad (24)$$

$$\propto \theta^{n+\alpha-1} \exp\left(-\left(\beta + \sum_{i=1}^n y_i\right)\theta\right) \quad (25)$$

$$\sim \text{Ga}\left(\alpha + n, \beta + \sum_{i=1}^n y_i\right) \quad (26)$$

It is easily recognized that the kernel of the posterior is again a gamma distribution.

(b) Using the logarithm and the first-order condition,

$$\ell(\theta | y_1, \dots, y_n) \propto n \log \theta - \theta \sum_{i=1}^n y_i \quad (27)$$

$$\frac{d}{d\theta} \ell(\theta | y_1, \dots, y_n) = \frac{n}{\theta} - \sum_{i=1}^n y_i = 0 \quad (28)$$

$$\hat{\theta}^{\text{MLE}} = n / \left(\sum_{i=1}^n y_i \right) \quad (29)$$

$$= \frac{1}{\bar{y}} \quad (30)$$

(c) For a gamma distribution with parameters α and β , the expected value is α/β . Thus, the posterior expected value is

$$\mathbf{E}(\theta | y_1, \dots, y_n) = \frac{\alpha + n}{\beta + \sum_{i=1}^n y_i} \quad (31)$$

The posterior mean can be decomposed into a weighted average between the prior mean and the MLE as follows:

$$\frac{\alpha + n}{\beta + \sum_{i=1}^n y_i} = \frac{\beta}{\beta + \sum_{i=1}^n y_i} \times \frac{\alpha}{\beta} + \frac{\sum_{i=1}^n y_i}{\beta + \sum_{i=1}^n y_i} \times \frac{n}{\sum_{i=1}^n y_i} \quad (32)$$

(d) Taking the limit, the weight on the prior mean tends to infinity,

$$\frac{\beta}{\beta + \sum_{i=1}^n y_i} \xrightarrow{n \rightarrow \infty} 0 \quad (33)$$

because $\sum_{i=1}^n y_i \rightarrow \infty$. Furthermore, the weight on the MLE converges to 1 with which we can conclude that the more we data, the closer the posterior mean gets to MLE. In short, with a large amount of data, the posterior mean is not very different from MLE.

4. (Greenberg, chapter 3) Compute the predictive distribution for y_{n+1} if the y_i have independent normal distributions $\mathcal{N}(\mu, 1)$, where the prior distribution for μ is $\mathcal{N}(\mu_0, \sigma_0^2)$.

Solution: First, the posterior distribution of μ is

$$\pi(\mu | y_1, \dots, y_n) \propto \exp\left(-\frac{1}{2} \sum_{i=1}^n (y_i - \mu)^2\right) \times \exp\left(-\frac{1}{2\sigma_0^2} (\mu - \mu_0)^2\right) \quad (34)$$

$$\propto \exp\left(-\frac{1}{2} \left(\sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \mu)^2\right) - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2\right) \quad (35)$$

$$\propto \exp\left(-\frac{n}{2} (\mu^2 - 2\bar{y}\mu) - \frac{1}{2\sigma_0^2} (\mu^2 - 2\mu_0\mu)\right) \quad (36)$$

$$\propto \exp\left(-\frac{1}{2} \left(\left(n + \frac{1}{\sigma_0^2}\right) \mu^2 - 2\left(n\bar{y} + \frac{1}{\sigma_0^2} \mu_0\right) \mu\right)\right) \quad (37)$$

Therefore, $\mu | y_1, \dots, y_n \sim \mathcal{N}(\mu_n, \sigma_n^2)$ where

$$\sigma_n^2 = \left(n + \frac{1}{\sigma_0^2}\right)^{-1} \quad (38)$$

$$\mu_n = \sigma_n^2 \left(n\bar{y} + \frac{\mu_0}{\sigma_0^2}\right) \quad (39)$$

By definition of the predictive distribution,

$$p(y_{n+1} | y_n, \dots, y_1) = \int_{-\infty}^{\infty} p(y_{n+1} | \mu) \cdot \pi(\mu | y_1, \dots, y_n) d\mu \quad (40)$$

$$= \frac{e^{y_{n+1}^2 + \mu_n^2 / \sigma_n^2}}{2\pi\sigma_n} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2} \left(\left(1 + \frac{1}{\sigma_n^2}\right) \mu^2 - 2(y_{n+1} + \mu_n) \mu\right)\right) d\mu \quad (41)$$

$$= \frac{\exp(y_{n+1}^2 + \mu_n^2 / \sigma_n^2)}{\sqrt{2\pi\sigma_n^2}} \cdot \left(1 + \frac{1}{\sigma_n^2}\right)^{-1} \quad (42)$$

5. Show that the median of the posterior distribution minimizes the absolute loss function.

Solution: Let X be a random variable with a distribution function F which is a legitimate measure and let \mathbf{E} denote the expectation with respect to the measure F . Then, if we let m denote the median, we can assume that F is zero to the left of an arbitrary constant A and one to the right of B .

$$\phi = \mathbf{E} |X - m| \quad (43)$$

$$= \int_A^m (m - x) dF + \int_m^B (x - m) dF \quad (44)$$

By interchanging the differentiation and integral (since we need to differentiate so as to get the

minimum—first-order condition),

$$\frac{d\phi}{dm} = \int_A^m 1 dF - \int_m^B 1 dF = 0 \quad (45)$$

We should solve

$$\int_A^m 1 dF = \int_m^B 1 dF \quad (46)$$

which is essentially comparing two probability measures, $\Pr(X < m) = \Pr(X > m)$. Since the probability below a constant m and the probability above it are identical, by definition, m is the median that minimizes the absolute loss function. By the same logic, the posterior density must be a Radon-Nikodym derivative of some absolutely continuous measure P with respect to either the Lebesgue measure or the counting measure. Then, the posterior median must be the solution to

$$\int_A^m 1 dP = \int_m^B 1 dP \quad (47)$$

which translates to $\Pr(\theta < m | X) = \Pr(\theta > m | X)$. Therefore, the posterior median minimizes the absolute loss function $\mathbf{E}(|\theta - m| | X)$.

6. A vector of observations at time t , $x_t = (x_{1,t}, x_{2,t}, \dots, x_{k,t})'$ has a multinomial distribution, a priori, with $\sum_{i=1}^k x_{i,t} = N$ and probabilities $p = (p_1, p_2, \dots, p_k)$. p_i is the probability of $x_{i,t}$ and $\sum_{i=1}^k p_i = 1$. Then, its joint probability of x_t is given by

$$f(x_t | p) = \frac{N!}{\prod_{i=1}^k x_{i,t}!} \prod_{i=1}^k p_i^{x_{i,t}}, \quad 0 \leq x_{i,t} \leq N. \quad (48)$$

Note that x_t is identically and independently distributed for $t = 1, 2, \dots, T$. Then, we assume that the joint prior distribution for p is a Dirichlet distribution with parameter $(\alpha_1, \alpha_2, \dots, \alpha_k)$ and its density is given by

$$\frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{i=1}^k p_i^{\alpha_i-1}, \quad 0 < p_i < 1, \quad \alpha_i > 0, \quad i = 1, 2, \dots, k \quad (49)$$

Derive the posterior distribution of p .

Solution: Multiplying the likelihood function and the prior,

$$\pi(p | \mathbf{x}) \propto \pi(p) \prod_{t=1}^T f(x_t | p) \quad (50)$$

$$\propto \prod_{i=1}^k p_k^{\alpha_i - 1} \prod_{t=1}^T \prod_{i=1}^k p_i^{x_{i,t}} \quad (51)$$

$$\propto \prod_{i=1}^k p_i^{\sum_{t=1}^T x_{i,t} + \alpha_i - 1} \quad (52)$$

This is the kernel of a Dirichlet distribution. Thus,

$$p | \mathbf{x} \sim \text{Dir} \left(\sum_{t=1}^T x_{1,t} + \alpha_1, \dots, \sum_{t=1}^T x_{k,t} + \alpha_k \right) \quad (53)$$

Since the posterior belongs to the same parametric family as the prior, it is conjugate.

7. Show that the mode of the posterior distribution minimizes the following loss function

$$L_3(\hat{\theta} | \theta) = \mathbf{I}(|\hat{\theta} - \theta| > b) \quad (54)$$

for any small $b > 0$.

Solution: The expected zero-one loss is computed as follows.

$$\mathbf{E}(\mathbf{I}(|\theta - \hat{\theta}| > b) | X) = \int_{-\infty}^{\infty} \mathbf{I}(|\theta - \hat{\theta}| > b) p(\theta | X) d\theta \quad (55)$$

$$= \int_{-\infty}^{\infty} (1 - \mathbf{I}(|\theta - \hat{\theta}| \leq b)) p(\theta | X) d\theta \quad (56)$$

$$= 1 - \int_{|\theta - \hat{\theta}| \leq b} p(\theta | X) d\theta \quad (57)$$

$$= 1 - \Pr(|\theta - \hat{\theta}| \leq b | X) \quad (58)$$

$$= 1 - \Pr(\hat{\theta} - b \leq \theta \leq \hat{\theta} + b | X) \quad (59)$$

With a small arbitrary constant b , the posterior probability measure on the neighborhood of $\hat{\theta}$ of distance b is maximized when $\hat{\theta}$ is the posterior mode. This is equivalent to saying 1 minus the posterior probability measure on the neighborhood of $\hat{\theta}$ of distance b is minimized when $\hat{\theta}$ is the posterior mode. Thus, the posterior mode minimizes the expected zero-one loss.

8. Explain the following terminologies.

- (1) Prior distribution
- (2) Model
- (3) Posterior distribution
- (4) Likelihood
- (5) Marginal likelihood
- (6) Posterior predictive distribution
- (7) Posterior predictive density
- (8) Predictive likelihood
- (9) Posterior probability of models
- (10) Full conditional distribution
- (11) Gibbs sampling

Solution:

- (1) A prior distribution is the distribution that encodes the belief that one has about the parameter before observing data.
- (2) A model is essentially a representation which describes the data generating process that the researcher believes in.
- (3) The posterior distribution is the updated distribution of the parameter after having incorporated all the observed data.
- (4) The likelihood is a function of the parameter given data. Although the likelihood is the same as the joint probability of the data, the distinction should be made between the two in that we only discuss *probability* before we observe the data whereas the likelihood is discussed after we have obtained the data.
- (5) A marginal likelihood is mathematically the expression in which the parameter has been integrated out. Intuitively, it is the distribution of the data that involves the uncertainty of the parameter. It also has an interpretation in terms of the predictive distribution.
- (6) The posterior predictive distribution is the distribution of the first observation, y_{T+1} in the future given the observations up to time T , y_1, \dots, y_T .

- (7) The posterior predictive density is the distribution of future H observations, y_{T+1}, \dots, y_{T+H} given the observations up to time T , y_1, \dots, y_T .
- (8) The predictive likelihood is
- (9) The posterior probability of models is the probability of postulated models, M_1, M_2, \dots after having observed data. Normally, the prior for each model is assumed to be equal across all models.
- (10) The full conditional distribution is the distribution of a chosen parameter given everything else.
- (11) Gibbs sampling is an MCMC algorithm that can be used if every parameter is conditionally conjugate, which is a special case of the Metropolis-Hastings algorithm.