

## Problem class 2: Bayesian point estimation, and Credible sets

Lecturer: Georgios Karagiannis

georgios.karagiannis@durham.ac.uk

### 1 Bayesian point estimation

**Exercise 1.** (★★) Consider observables  $x = (x_1, \dots, x_n)$ . Consider the Bayesian model

$$\begin{cases} x_i | \theta & \stackrel{\text{iid}}{\sim} \text{N}(\theta, 1), \quad i = 1, \dots, n \\ \theta & \sim \Pi(\theta) \end{cases}$$

where  $\pi(\theta) \propto 1$  and that we have only one observable. Consider the LINEX loss function

$$\ell(\theta, \delta) = \exp(c(\theta - \delta)) - c(\theta - \delta) - 1$$

1. Show that  $\ell(\theta, \delta) \geq 0$
2. Find the Bayes estimator  $\hat{\delta}$  under LINEX loss function and under the given Bayesian model.

**Hint-1:** Random variable  $B$  follows a log-normal distribution  $B \sim \text{LN}(\mu_A, \sigma_A^2)$  with parameters  $\mu_A, \sigma_A^2$  if  $B = \exp(A)$  where  $A \sim \text{N}(\mu_A, \sigma_A^2)$ .

**Hint-2:** If  $B \sim \text{LN}(\mu_A, \sigma_A^2)$  then  $E_{\text{LN}(\mu_A, \sigma_A^2)}(B) = \exp(\mu_A + \frac{\sigma_A^2}{2})$ .

**Hint-3:** It is

$$-\frac{1}{2} \frac{(\mu - \mu_1)^2}{v_1^2} - \frac{1}{2} \frac{(\mu - \mu_2)^2}{v_2^2} \dots - \frac{1}{2} \frac{(\mu - \mu_n)^2}{v_n^2} = -\frac{1}{2} \frac{(\mu - \hat{\mu})^2}{\hat{v}^2} + C$$

where

$$\hat{v}^2 = \left( \sum_{i=1}^n \frac{1}{v_i^2} \right)^{-1}; \quad \hat{\mu} = \hat{v}^2 \left( \sum_{i=1}^n \frac{\mu_i}{v_i^2} \right); \quad C = \frac{1}{2} \frac{\hat{\mu}^2}{\hat{v}^2} - \frac{1}{2} \sum_{i=1}^n \frac{\mu_i^2}{v_i^2}$$

**Solution.** So

1. Let  $g(x) = \exp(cx) - cx - 1$  with  $x = \theta - \delta$ . I observe that  $g$  is differential with  $g'(x) = c(\exp(cx) - 1)$ . Also  $g(\cdot)$  has a minimum at  $x = 0$ , as  $g'(0) = 0$ ,  $g''(0) > 0$ , and in general

$$g'(x) : \begin{cases} < 0, & \text{for } x < 0 \\ = 0, & \text{for } x = 0 \\ > 0, & \text{for } x > 0 \end{cases}$$

Moreover, it is  $\lim_{x \rightarrow -\infty} g(x) = \lim_{x \rightarrow +\infty} g(x) = +\infty$  and  $g(0) = 0$ . Hence  $g(x) > 0, \forall x > 0$ . In other words,  $\ell(\theta, \delta) \geq 0, \forall \theta > 0$ .

- 2.

- It is

$$\begin{aligned} \rho(\pi, \delta | x) &= E_{\Pi}(\ell(\theta, \delta) | x) = E_{\Pi}(\exp(c(\theta - \delta)) - c(\theta - \delta) - 1 | x) \\ &= \exp(-cd) E_{\Pi}(\exp(c\theta) | x) - c(E_{\Pi}(\theta | x) - \delta) - 1 \end{aligned}$$

- I will minimize the posterior expected risk to find the Bayes estimator (rule). So, it is

$$\begin{aligned}\frac{d}{d\delta}\rho(\pi, \delta|x) &= -c \exp(-c\delta) E_{\Pi}(\exp(c\theta)|x) + c \\ 0 &= \left. \frac{d}{d\delta}\rho(\pi, \delta|x) \right|_{\delta=\delta^{\pi}} \\ 0 &= -c \exp(-c\delta^{\pi}) E_{\Pi}(\exp(c\theta)|x) + c \\ \delta^{\pi} &= \frac{1}{c} \log(E^{\pi}(\exp(c\theta)|x))\end{aligned}$$

- By using the Bayes theorem,

$$\begin{aligned}\pi(\theta|x) &\propto \prod_{i=1}^n N(x_i|\theta, 1) \pi(\theta) \propto \prod_{i=1}^n N(x_i|\theta, 1) \\ &\propto \exp\left(-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2\right) \propto \exp\left(-\frac{1}{2} \frac{(\theta - \hat{\theta})^2}{\hat{v}^2} + \text{const...}\right)\end{aligned}$$

with

$$\hat{v}^2 = \left(\sum_{i=1}^n \frac{1}{v_i^2}\right)^{-1} = \left(\sum_{i=1}^n \frac{1}{1}\right)^{-1} = \frac{1}{n}; \quad \hat{\theta} = \hat{v}^2 \left(\sum_{i=1}^n \frac{\mu_i}{v_i^2}\right) = \frac{1}{n} \left(\sum_{i=1}^n x_i\right) = \bar{x}$$

So the posterior distribution is

$$\theta|x \sim N(\bar{x}, 1/n)$$

- Now let's assume that  $E_{\Pi}(\exp(c\theta)|x) < \infty$ .
- Then  $E_{\Pi}(\exp(c\theta)|x) = E_{N(\bar{x}, 1/n)}(\underbrace{\exp(c\theta)}_{=\tilde{\theta}}|x) = E_{LN(c\bar{x}, c^2/2n)}(\tilde{\theta}|x) = \exp(c\bar{x} + c^2/2n)$  (as an expected value of LN distribution). Hence,

$$\delta^{\pi}(x) = c\bar{x} + c^2/2n$$

.

**Exercise 2.** (\*\*) Suppose we wish to estimate the values of a collection of discrete random variables  $\vec{X} = X_1, \dots, X_n$ . We have a posterior joint probability mass function for these variables,  $p(\vec{x}|y) = p(x_1, \dots, x_n|y)$  based on some data  $y$ . We decide to use the following loss function:

$$\ell(\hat{\vec{x}}, \vec{x}) = \sum_{i=1}^n (1 - \delta(\hat{x}_i, x_i)) \quad (1)$$

where  $\delta(a, b) = 1$  if  $a = b$  and zero otherwise.

1. Derive an expression for the estimated values, found by minimizing the expectation of the loss function. [Hint: use linearity of expectation.]
2. When the probability distribution is a posterior distribution in some problem, this type of estimate is sometimes called 'maximum posterior marginal' (MPM) estimate. Explain why this name is appropriate.
3. Explain in words what the loss function is measuring. Compare with the loss function for MAP estimation.

**Solution.**

1. We have that

$$E \left( \ell(\hat{\vec{x}}, \vec{X}) | y \right) = E \left( 1 - \sum_{i=1}^n \delta(\hat{x}_i, X_i) | y \right) \quad (2)$$

$$= n - \sum_{i=1}^n E(\delta(\hat{x}_i, X_i) | y) = n - \sum_{i=1}^n \sum_{x_i \in \mathcal{X}_i} \delta(\hat{x}_i, x_i) p(x_i | y) = n - \sum_{i=1}^n p(\hat{x}_i | y) \quad (3)$$

To minimize this, it suffices to minimize each term of the sum separately, and so we have for each  $i \in \{1, \dots, n\}$  separately:

$$\hat{x}_i^* = \arg \max_{x_i \in \mathcal{X}_i} p(x_i | y) \quad (4)$$

[The  $x_i$  in (4) corresponds to the  $\hat{x}_i$  in ; we simply drop the hat to keep notation as simple as possible.]

Contrast this with the MAP estimate, which requires optimisation over the full joint pmf:

$$\hat{\vec{x}}^* = \arg \max_{\vec{x} \in \vec{\mathcal{X}}} p(\hat{\vec{x}} | y) \quad (5)$$

2. The individual terms of the sum in (3) are the posterior marginal distributions for each  $x_i$  found from the joint distribution  $p(\hat{x}_i | y)$ . The estimate for each  $x_i$  is found by maximizing its own posterior marginal distribution, hence the name.

3. MPM estimation has attractive properties. Like MAP estimation, it can be defined on any set (albeit with the same caveats as for MAP with continuous variables). On the other hand, it does not insist that all variables ‘match’ to get the minimum loss. Rather, it counts the number that match by summing over the individual zero-one losses. The loss function for MAP estimation, on the other hand, imposes a different penalization since it is

$$\ell(\hat{\vec{x}}, \vec{x}) = 1 - \delta(\hat{\vec{x}}, \vec{x}) = 1 - \prod_{i=1}^n \delta(\hat{x}_i, x_i) \quad (6)$$

where by taking the product, the loss will be one unless *all* variables match.

## 2 Credible sets

**Exercise 3.** (★★) (Example from the Lecture’s handout) Consider a Bayesian model

$$\begin{cases} y_i | \mu & \stackrel{\text{iid}}{\sim} N_d(\mu, \Sigma), & i = 1, \dots, n \\ \mu & \sim N_d(\mu_0, \Sigma_0) \end{cases}$$

where uncertain  $\mu \in \mathbb{R}^d$ ,  $d \geq 1$ , and known  $\Sigma > 0$ ,  $\mu_0, \Sigma_0 > 0$ . Find the  $C_a$  parametric HPD credible set for  $\mu$ .

**Hint-1:** If  $z = (z_1, \dots, z_d)^\top$  such as  $z_j \stackrel{\text{iid}}{\sim} N(0, 1)$  for  $j = 1, \dots, d$ , and  $\xi = z^\top z = \sum_{j=1}^d z_j^2$ , then  $\xi \sim \chi_d^2$

**Hint-2:** It is

$$\begin{aligned}
-\frac{1}{2} \sum_{i=1}^n (x - \mu_i)^\top \Sigma_i^{-1} (x - \mu_i) &= -\frac{1}{2} (x - \hat{\mu})^\top \hat{\Sigma}^{-1} (x - \hat{\mu}) + C(\hat{\mu}, \hat{\Sigma}) \quad ; \\
\hat{\Sigma} &= \left( \sum_{i=1}^n \Sigma_i^{-1} \right)^{-1}; \quad \hat{\mu} = \hat{\Sigma} \left( \sum_{i=1}^n \Sigma_i^{-1} \mu_i \right); \\
C(\hat{\mu}, \hat{\Sigma}) &= \underbrace{\frac{1}{2} \left( \sum_{i=1}^n \Sigma_i^{-1} \mu_i \right)^\top \left( \sum_{i=1}^n \Sigma_i^{-1} \right)^{-1} \left( \sum_{i=1}^n \Sigma_i^{-1} \mu_i \right) - \frac{1}{2} \sum_{i=1}^n \mu_i^\top \Sigma_i^{-1} \mu_i}_{=\text{independent of } x}
\end{aligned}$$

**Solution.**

I will use the Definition of HPD credible interval.

- First, I compute the posterior of  $\mu$ . It is

$$\begin{aligned}
\pi(\mu|y) &\propto f(y|\mu)\pi(\mu) = \prod_{i=1}^n \mathcal{N}_d(y_i|\mu, \Sigma) \mathcal{N}_d(\mu|\mu_0, \Sigma_0) \\
&\propto \exp \left( -\frac{1}{2} \sum_{i=1}^n (y_i - \mu)^\top \Sigma^{-1} (y_i - \mu) - \frac{1}{2} (\mu - \mu_0)^\top \Sigma_0^{-1} (\mu - \mu_0) \right) \\
&\propto \exp \left( -\frac{1}{2} (\mu - \hat{\mu}_n)^\top \hat{\Sigma}_n^{-1} (\mu - \hat{\mu}_n) \right)
\end{aligned}$$

where

$$\hat{\Sigma}_n = (n\Sigma^{-1} + \Sigma_0^{-1})^{-1}; \quad \hat{\mu}_n = \hat{\Sigma}_n (n\Sigma^{-1}\bar{y} + \Sigma_0^{-1}\mu_0)$$

I recognize that  $\pi(\mu|y) = \mathcal{N}_d(\mu|\hat{\mu}_n, \hat{\Sigma}_n)$ , and hence  $\mu|y \sim \mathcal{N}_d(\hat{\mu}_n, \hat{\Sigma}_n)$

- Now let's implement Definition of HPD credible interval. So,

$$\begin{aligned}
C_a &= \{ \mu \in \mathbb{R}^d : \pi(\mu|y) \geq k_a \} \\
&= \{ \mu \in \mathbb{R}^d : \mathcal{N}_d(\mu|\hat{\mu}_n, \hat{\Sigma}_n) \geq k_a \} \\
&= \left\{ \mu \in \mathbb{R}^d : (\mu - \hat{\mu}_n)^\top \hat{\Sigma}_n^{-1} (\mu - \hat{\mu}_n) \leq \underbrace{-2 \log \left( (2\pi)^{\frac{d}{2}} \det(\hat{\Sigma}_n) k_a \right)}_{=\tilde{k}_a} \right\} \quad (7)
\end{aligned}$$

and I want the smallest constant  $\tilde{k}_a$  (aka the largest constant  $k_a$ ) such that

$$\begin{aligned}
&\Pr_{\Pi}(\mu \in C_a|y) \geq 1 - a \iff \\
&\Pr_{\Pi} \left( \underbrace{(\mu - \hat{\mu}_n)^\top \hat{\Sigma}_n^{-1} (\mu - \hat{\mu}_n)}_{=\xi} \leq \tilde{k}_a \right) \geq 1 - a \quad (8)
\end{aligned}$$

- I need to find quantile  $\tilde{k}_a$ . This requires to find the distribution of  $\xi$ . I know that

$$\xi = (\mu - \hat{\mu}_n)^\top \hat{\Sigma}_n^{-1} (\mu - \hat{\mu}_n) \sim \chi_d^2 \quad (9)$$

because  $\xi = z^\top z = \sum_{j=1}^n z_j$  with  $z = L^{-1}(\mu - \hat{\mu}_n) \sim N_d(0, I_d)$  where  $L$  is the lower matrix of the Cholesky decomposition of  $\hat{\Sigma}_n = L^\top L$ .

Hence Eq. 8, (due to Eqs. 7, 9) becomes

$$\Pr_{\chi_d^2}((\mu - \hat{\mu}_n)^\top \hat{\Sigma}_n^{-1}(\mu - \hat{\mu}_n) \leq \tilde{k}_a) = 1 - a \quad (10)$$

which means that,  $\tilde{k}_a$  is the  $1 - a$  quantile of the  $\chi_d^2$  distribution, aka  $\tilde{k}_a = \chi_{d,1-a}^2$

- Hence, the  $C_a$  parametric HPD credible set for  $\mu$  is

$$C_a = \{\mu \in \mathbb{R}^d : (\mu - \hat{\mu}_n)^\top \hat{\Sigma}_n^{-1}(\mu - \hat{\mu}_n) \leq \chi_{d,1-a}^2\}$$

**Example 4.** (★★) (Example from the Lecture's handout) Assume an 1- dimensional random quantity  $x \sim Q(x|y)$ . In the Lecture Handout (Handout 11: Bayesian point estimation), discussed the following Hint:

**Hint:** The Bayes estimate  $\hat{\delta}$  of  $x$  under the linear loss function

$$\ell(x, \delta; \varpi) = (1 - \varpi)(\delta - x)1_{x \leq \delta} + \varpi(x - \delta)1_{x > \delta},$$

where  $\varpi \in [0, 1]$ , is the  $\varpi$ -th quantile of distribution  $Q$ , let's denote it as  $x_\varpi$ .

1. Derive the  $(1 - a)$ -credible interval  $C_a = [L, U]$  for  $x$  as a Bayesian rule  $C_a$  under the loss function

$$\ell(x, C_a; \varpi_L, \varpi_U) = \ell(x, L; \varpi_L) + \ell(x, U; \varpi_U) \quad (11)$$

by computing  $L$  and  $U$ .

2. Your client is worried the same both for under-estimation and over-estimation; derive a suitable  $(1 - a)$ -credible interval  $C_a = [L, U]$  based on (11) by computing  $L$ , and  $U$ .
3. Your client is worried only for over-estimation; derive a suitable  $(1 - a)$ -credible interval  $C_a = [L, U]$  based on (11) by computing  $L$  and  $U$ .

**Solution.** It is given that

$$\begin{aligned} 0 &= \frac{d}{d\delta} E_Q(\ell(x, \delta; \varpi)|y) \Big|_{\delta=\hat{\delta}} = \frac{d}{d\delta} \int \ell(x, \delta; \varpi) dQ(x|y) \Big|_{\delta=\hat{\delta}} \implies \hat{\delta} = x_\varpi \\ &= (1 - \varpi) \Pr_Q(x \leq \hat{\delta}|y) - \varpi \Pr_Q(x > \hat{\delta}|y) \implies \hat{\delta} = x_\varpi \end{aligned}$$

1. The decision space is  $\mathcal{D} = \{C_a = [L, U] : \Pr_Q(x \in C_a|y) = 1 - a\}$ . Therefore, to find the Bayes rule (or Bayes estimate) of  $C_a = [L, U]$  I need to minimize the expected posterior loss  $E_Q(\ell(x, C_a; \varpi_L, \varpi_U)|y)$  with respect to  $C_a$  or equivalently  $L, U$ , so

$$\begin{aligned} 0 &= \frac{d}{dL} E_Q(\ell(x, C_a; \varpi_L, \varpi_U)|y) \Big|_{C_a=[\hat{L}, \hat{U}]} = E_Q(\ell(x, L; \varpi_L)|y) \Big|_{L=\hat{L}} \implies \hat{L} = x_{\varpi_L} \\ 0 &= \frac{d}{dU} E_Q(\ell(x, C_a; \varpi_L, \varpi_U)|y) \Big|_{C_a=[\hat{L}, \hat{U}]} = E_Q(\ell(x, U; \varpi_U)|y) \Big|_{U=\hat{U}} \implies \hat{U} = x_{\varpi_U} \end{aligned}$$

So  $x \in [x_{\varpi_L}, x_{\varpi_U}]$  where  $\varpi_U + \varpi_L = 1 - a$ . It is the minimum because

$$\frac{d^2}{dU^2} E_Q(\ell(x, C_a; \varpi_L, \varpi_U)|y) \Big|_{C_a=[\hat{L}, \hat{U}]} = q(\hat{U}|y) > 0$$

$$\frac{d^2}{dL^2} E_Q (\ell(x, C_a; \varpi_L, \varpi_U) | y) \Big|_{C_a = [\hat{L}, \hat{U}]} = q(\hat{L} | y) > 0$$

$$\frac{d}{dU} \frac{d}{dL} E_Q (\ell(x, C_a; \varpi_L, \varpi_U) | y) \Big|_{C_a = [\hat{L}, \hat{U}]} = 0$$

and hence the determinant of the Hessian is positive.

- (a) Then I can use the equi-tail interval:  $x \in [x_{a/2}, x_{1-a/2}]$  with  $\varpi_L = a/2$  and  $\varpi_U = 1 - a/2$
- (b) Then I can use the lower-tail interval:  $x \in (-\infty, x_{1-a}]$  with  $\varpi_L = 0$  and  $\varpi_U = 1 - a$ .