

STAT 4010 – Bayesian Learning

TUTORIAL 7

Spring 2022

Cheuk Hin (Andy) CHENG ([Email](#) | [Homepage](#))

Di SU ([Email](#) | [Homepage](#))

1 Region Estimation

Definition 1. Given a prior distribution $\theta \sim \pi(\theta)$, and $\alpha \in (0, 1)$. Then

1. A set \hat{I} is said to an $100(1 - \alpha)\%$ **credible set** if $P(\theta \in \hat{I} | x) \geq 1 - \alpha$.
2. The $100(1 - \alpha)\%$ **HPD credible set** (where HPD stands for highest posterior density) is

$$\hat{I}_{\text{HPD}} = \{\theta \in \Theta : \pi(\theta | x) \geq k\}$$

where k is the largest constant such that $P(\theta \in \hat{I}_{\text{HPD}} | x) \geq 1 - \alpha$.

Assume the posterior is continuous. Then the HPD credible set is the Bayes estimator under the loss

$$L(\theta, \hat{I}) = k|\hat{I}| + \mathbb{1}(\theta \notin \hat{I}), \quad (1.1)$$

which measures both the coverage and the width of \hat{I} .

3. The $100(1 - \alpha)\%$ **equal tailed credible set** is

$$\hat{I}_{\text{ET}} = [Q_{\alpha/2}(\theta | x), Q_{1-\alpha/2}(\theta | x)]$$

where $Q_p(\theta | x)$ is the p th quantile of $[\theta | x]$.

Lemma 1.1. Suppose the PDF $\theta \mapsto \pi(\theta | x)$ is (i) symmetric about $\theta = \theta_0$, (ii) continuous, and (iii) strictly increasing in $(-\infty, \theta_0]$ and is strictly decreasing in $[\theta_0, +\infty)$. Then, given the same credible level $1 - \alpha$,

$$\hat{I}_{\text{HPD}} = \hat{I}_{\text{ET}}$$

Definition 2. Let $g(\cdot)$ be a strictly increasing and continuous function, and $\phi = g(\theta)$. A principle for constructing region estimator is invariant under the transformation g if

$$[\hat{L}_\phi, \hat{U}_\phi] = [g(\hat{L}_\theta), g(\hat{U}_\theta)]$$

where $[\hat{L}_\phi, \hat{U}_\phi]$ and $[\hat{L}_\theta, \hat{U}_\theta]$ are region estimators for ϕ and θ , respectively.

Lemma 1.2. Suppose that $[\theta | x]$ is a continuous RV.

1. The "principle of HPD" for constructing credible sets is NOT invariant to continuous re-parametrization.
2. The "principle of equal tails" for constructing credible sets is invariant to continuous re-parametrization. (Because of the invariance of the quantile function)

Example 1.1. (The Bayes estimator of the confidence index) Under the quadratic loss

$$L(\theta, \hat{\alpha}) = \{\hat{\alpha} - \mathbb{1}(\theta \in I)\}^2,$$

where $\hat{\alpha} \in [0, 1]$, the Bayes estimator is given by $\hat{\alpha}_\pi = P(\theta \in I | x)$.

Remark 1.1. (Riemann Sum) Let $f : [a, b] \rightarrow \mathbb{R}$ be a function defined on a closed interval $I = [a, b]$ of the real numbers, \mathbb{R} , and

$$P = \{[x_0, x_1], [x_1, x_2], \dots, [x_{n-1}, x_n]\},$$

be a partition of I , where

$$a = x_0 < x_1 < x_2 < \dots < x_n = b.$$

A Riemann sum S of f over I with partition P is defined as

$$S = \sum_{i=1}^n f(x_i^*) \Delta x_i$$

where $\Delta x_i = x_i - x_{i-1}$ and $x_i^* \in [x_{i-1}, x_i]$. This limiting value, if it exists, is defined as the definite Riemann integral of the function over the domain,

$$\int_a^b f(x) dx = \lim_{\|\Delta x\| \rightarrow 0} \sum_{i=1}^n f(x_i^*) \Delta x_i.$$

Thus, to find $\int_a^b f(x) dx$, we can divide I as $I = [a, a+h] \cup [a+h, a+2h] \cup \dots \cup [a+mh, b]$ where h is a step size that is close to zero, and $m = \lfloor (a-b)/h \rfloor - 1$. Then we have the following numerical approximation

$$\int_a^b f(x) dx \approx \sum_{i=0}^m f(a+ih)h \approx (a-b)\bar{f}_i, \quad (1.2)$$

where \bar{f}_i is the average of $\{f(a+ih)\}_{i=0, \dots, m}$. Riemann sum is named after nineteenth century German mathematician Bernhard Riemann.

Example 1.2. (Wild guess - A4 Fall 2019). After grading the mid-term exam, Keith would like to analyze how well the students did in the multiple-choice questions (MCQs). Let $x_{ij} = \mathbb{1}$ (the i th student answer the j MCQ correctly), $i = 1, \dots, n$, $j = 1, \dots, m$, where $n = 42$ and $m = 8$. Consider the model

$$[x_{ij} | \theta] \stackrel{\text{iid}}{\sim} \text{Bern}(\theta), \quad i = 1, \dots, n, \quad j = 1, \dots, m,$$

and the prior $\theta \sim \pi(\theta)$ defined by

$$\theta \sim \text{Beta}(\alpha, \beta),$$

where $\alpha = \beta = 1/2$. Denote the entire dataset by $X = \{x_{ij} : 1 \leq i \leq n, 1 \leq j \leq m\}$. The dataset MCQ.txt can be downloaded from the course website. Because of the privacy concern, not less than 10% of the data have been replaced by imputed values such that the structure of the privacy-protected dataset is close to that of the actual one.

Consider $\hat{I} \in \{[a, b] : 0 \leq a \leq b \leq 1\}$, and the loss

$$L(\theta, \hat{I}) = k|\hat{I}| + \mathbb{1}(\theta \notin \hat{I})$$

where $|\hat{I}|$ denotes the width of the interval \hat{I} , and $k = 1$ is fixed.

1. Compute the Bayes estimator \hat{I}_π . What is the credible level $P(\theta \in \hat{I}_\pi | X)$ of \hat{I}_π ?
2. Compute the 95% HPD credible interval \hat{I}_{HPD} of θ .
3. Compute the 95% equal-tailed credible interval \hat{I}_{ET} of θ .

(Hints: (a) Example 5.5. (b): Example 5.7. (c) Example 5.1.)

SOLUTION:

1. By Example 2.13,

$$[\theta | x_{ij}] \sim \text{Beta}(\alpha + S_n, \beta + n - S_n),$$

where $S_n = \sum_{i=1}^n \sum_{j=1}^m x_{ij}$. After some calculation, we get $\alpha + S_n = 140.5$, $\beta + n - S_n = 196.5$. By Example 5.5, we have that

$$\hat{I}_\pi = \hat{I}_{\text{HPD}},$$

with $k = 1$. By R, we have $\hat{I}_{\text{HPD}} = [0.355, 0.480]$ is the Bayes estimator up to 3 decimal places. Its credible level is $P(\theta \in \hat{I}_\pi | X) = P(\theta \in \hat{I}_{\text{HPD}} | X) = 0.980$.

```

1 ### Step 1: Initialization
2 data = read.table("MCQ.txt", header = T);
3 X = as.matrix(data);
4 a = 0.5; b=0.5;
5 n = dim(X)[1];
6 m = dim(X)[2];
7 alpha.b = sum(X) + a;
8 beta.b = n*m - sum(X) + b
9 # alpha.b = 140.5
10 # beta.b = 196.5
11
12 ### Step 2: Find Bayes estimator
13 theta = seq(0, 1, length = 5001)
14 # notice that the posterior has a closed form
15 d = dbeta(theta, alpha.b, beta.b)
16 L = theta[which.min(abs(d[1:which.max(d)]-1))]
17 U = theta[which.min(abs(d[which.max(d)+1:length(d)]-1))+which.max(d)]
18 L
19 [1] 0.355
20 U
21 [1] 0.4796
22
23 ### Step 3: Find Credible level
24 pbeta(U, alpha.b, beta.b) - pbeta(L, alpha.b, beta.b)
25 [1] 0.9800575

```

2. Using R, $\hat{I}_{\text{HPD}} = [0.4164, 0.470]$ with credible level: 0.950.

```

1 # Method 1
2 {
3   alpha = 0.05
4   ### step 1: values of posterior at different values of theta in [a,b]
5   theta = seq(from=0.3, to=0.5, length=501)
6   d = dbeta(theta, alpha.b, beta.b)
7
8   ### step 2: find the theta that satisfy the credibility requirement
9   # O stores indices
10  O = order(d, decreasing=TRUE)
11  # confidence w.r.t. all interval candidates
12
13  ##---Learning Moment: use Riemann sum to approximate integral---##
14  step = theta[2]-theta[1]
15  conf = cumsum(d[O]*step)
16
17  # check confidence condition
18  N = sum(conf<(1-alpha))+1
19  selected.index = O[1:N]
20  selected = theta[selected.index]
21  # results
22  k = d[O[N]]
23  selected[1]
24  selected[N]
25  alpha.hat = conf[N]
26  alpha.hat
27
28  ### step 3: plot
29  plot(theta,d, type="l", lwd=2, col="red4",
30        xlab=expression(theta),
31        ylab=expression(pi(theta~"|"~italic(x[1:n]))))
32  abline(v=selected,col="pink")
33  abline(h=c(0,k),v=c(a,b),lty=3, lwd=.75)
34
35  ### exploration
36  abline(v = theta[O[1:10]],col='yellow')
37  abline(v = theta[O[1:50]],col='yellow')
38  abline(v = theta[O[1:100]],col='yellow')
39  abline(v = theta[O[1:150]],col='yellow')
40
41  abline(h = d[O[10]],lty = 2,col='orange')
42  abline(h = d[O[50]],lty = 2,col='orange')
43  abline(h = d[O[100]],lty = 2,col='orange')
44  abline(h = d[O[150]],lty = 2,col='orange')
45
46  conf[1:50]
47
48 }
49
50 # Method 2
51 {
52   ### search for all possible HPD set (with different credible level)
53   i.max = which.max(d)
54   out = array(NA, dim=c(i.max, 3))
55   colnames(out) = c("Lower-bound", "Upper-bound", "Credible-level")
56   for(i.left in 1:i.max){
57     delta = abs(d[i.left]-d[-(1:i.max)])
58     out[i.left,1] = theta[i.left]

```

```

59     i.right = i.max+which.min(delta)
60     out[i.left,2] = theta[i.right]
61     out[i.left,3] = pbeta(out[i.left,2], alpha.b, beta.b) - pbeta(out[i.
        left,1], alpha.b, beta.b)
62 }
63
64 ### select the HPD set with a desired creible level
65 (result = out[which.min(abs(out[,3]-0.95)),])
66
67 ### plot
68 abline(v=c(result[1:2]), col="blue4", lwd=3)
69
70 # exploration
71 plot(theta,d, type="l", lwd=2, col="red4",
72       xlab=expression(theta),
73       ylab=expression(pi(theta~"|"~italic(x[1:n]))))
74 abline(h=c(0,k), v=c(a,b), lty=3, lwd=.75)
75 abline(v=c(result[1:2]), col="blue4", lwd=3)
76 abline(v = c(out[50,1], out[50,2]), col='blue', lty=4)
77 abline(v = c(out[100,1], out[100,2]), col='blue', lty=4)
78 abline(v = c(out[150,1], out[150,2]), col='blue', lty=4)
79
80 }
81
82 ##---Learning Moment: use Riemann sum to approximate integral---##
83 riemann = function(f, left, right, step){
84   theta = seq(from = left, to = right, by = step)
85   f.all = f(theta)
86   sum1 = sum(f.all*step)
87   sum2 = mean(f.all)*(right-left)
88   sum3 = pbeta(right, alpha.b, beta.b) - pbeta(left, alpha.b, beta.b)
89   out = array(c(sum1, sum2, sum3), dimnames=list(c('Riemann Sum', 'Sample
        average', 'pbeta')))
90   out
91 }
92 post = function(x) dbeta(x, alpha.b, beta.b)
93 l=0.3
94 r=0.5
95 riemann(f = post, left = l, right = r, step = 0.01)

```

Remark 1.2. Method 1 is better than Method 2 in the sense that Method 1 starts with the shortest intervals containing the MAP, and needs to consider less candidates. In assignments and projects, write a HPD function with suitable inputs and outputs.

3. Using R, we have $\hat{I}_{ET} = [0.365, 0.470]$.

```

1   qbeta(c(0.025, 0.975), alpha.b, beta.b)
2   [1] 0.3648527 0.4699094

```

2 Theoretical Justification

This section shows that the Bayesian methods studied in previous chapters are theoretically sensible.

Definition 3. Given any DGP $f_*(x)$ and model $\mathcal{F} = \{f(x | \theta) : \theta \in \Theta\}$. Denote the expectation and variance under the DPG $f_*(x)$ by E_* and Var_* . Define

$$\theta_* = \arg \max_{\theta \in \Theta} E_* \{ \log f(x_1 | \theta) \},$$

and

$$I_* = \left[\text{Var}_* \left\{ \frac{d}{d\theta} \log f(x_1 | \theta) \right\} \right]_{\theta=\theta_*} \quad J_* = \left[-E_* \left\{ \frac{d^2}{d\theta^2} \log f(x_1 | \theta) \right\} \right]_{\theta=\theta_*},$$

provided that the expectations exist. The quantities I_* and J_* are called Fisher information. If \mathcal{F} well specifies f_* , then $\theta_* = \theta_0$ and $I_* = J_*$, where θ_0 is the true DGP parameter.

Theorem 2.1. (Consistency of posterior). Assume regularity conditions (RCs). If n is large enough, then

$$\hat{\theta}_{\text{MLE}} \approx \theta_* \quad \text{and} \quad [\theta | x_{1:n}] \approx \theta_*.$$

Theorem 2.2. (Asymptotic distributions of posterior). Assume RCs. If n is large enough, then

$$\hat{\theta}_{\text{MLE}} \approx N \left(\theta_*, \frac{J_*^{-1} I_* J_*^{-1}}{n} \right) \quad \text{and} \quad [\theta | x_{1:n}] \approx N \left(\hat{\theta}_{\text{MLE}}, \frac{J_*^{-1}}{n} \right).$$

If the model is well-specified, the precision of Bayesian framework and frequentist framework are consistent.

Theorem 2.3. (Asymptotic representation of posterior mean). Assume RCs. If n is large enough, then

$$E(\theta | x_{1:n}) \approx \hat{\theta}_{\text{MLE}}.$$

Remark 2.1. Intuitively, if the model is well-specified, we will have the following approximation

$$[\theta | x_{1:n}] \text{ “} \approx \text{” } \hat{\theta}_{\text{MLE}} \approx \theta_* = \theta_0.$$

Be careful that $\hat{\theta}_{\text{MLE}}$ is the maximizer $\log f(x | \theta)$, whereas θ_* is the maximizer of $\log f(x | \theta)$ after taken expectation w.r.t. $f_*(x)$. In practice, f_* may not be known.

Theorem 2.4. We have the following bi-directional relation

$$x_{1:n} \text{ are exchangeable with joint density } f(x_{1:n}) \\ \iff \exists \theta \in \Theta, f(x | \theta), \pi(\theta) \text{ s.t. } \begin{cases} [x_{1:n} | \theta] \stackrel{IID}{\sim} f(x_{1:n} | \theta) \\ \theta \sim \pi(\theta). \end{cases}$$

The direction “ \implies ” is stated in theorem 6.5. De Finetti Theorem, and the direction “ \impliedby ” is given in proposition 6.4.

Remark 2.2. Bayesian model enables us to work with IID data, by “extracting out” the dependence within exchangeable data.

Proof. (of Proposition 6.4.) Denote the PDF of $y_{1:n}$ by $f(x_1, \dots, x_n)$. For any permutation σ of $\{1, \dots, n\}$, we have

$$\begin{aligned}
 f(x_1, \dots, x_n) &= \int_{\Theta} \pi(\theta) f(x_1, \dots, x_n \mid \theta) \, d\theta && \because \text{By the definition of marginal PDF} \\
 &= \int_{\Theta} \pi(\theta) \prod_{i=1}^n f(x_i \mid \theta) \, d\theta && \because y_{1:n} \text{ are conditionally independent given } \theta \\
 &= \int_{\Theta} \pi(\theta) \prod_{i=1}^n f(x_{\sigma(i)} \mid \theta) \, d\theta && \because \text{Product is commutative} \\
 &= f(x_{\sigma(1)}, \dots, x_{\sigma(n)}) && \because \text{By the definition of marginal PDF}
 \end{aligned}$$

Thus, $x_{1:n}$ are exchangeable.