

STAT 4010 – Bayesian Learning

TUTORIAL 4

Spring 2022

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

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

1 Decision Theory

Definition 1. (*Decision Theory*) Define the following settings,

$$\mathcal{X} = \text{Observation space}, \quad \Theta = \text{Parameter space}, \quad \mathcal{D} = \text{Decision space}.$$

Statistical inference is performed by using the observation $x \in \mathcal{X}$ to construct a decision $d(x) = d \in \mathcal{D}$ to understand the parameter $\theta \in \Theta$. The decision is made under the following systematic procedure.

- (Everyone) Fix a user-defined function, **utility function** for evaluating the goodness of a decision, defined as

$$U : \Theta \times \mathcal{D} \rightarrow (-\infty, 0].$$

- Frequentist procedure

- Compute the **loss function** to evaluate the loss of the decision, defined as

$$L(\theta, d) = -U(\theta, d).$$

- Compute the **(frequentist) risk** to average out x and the effect of sampling error, defined as

$$R(\theta, d) = \mathbb{E}_{x|\theta}[L(\theta, d) | \theta] = \int_{\mathcal{X}} L(\theta, d(x))f(x | \theta)dx.$$

The best decision is chosen by minimizing the risk.

- Bayesian procedure (Keep in mind θ is no longer a constant):

- Compute the **posterior loss function** to evaluate the loss of the decision, defined as

$$L(\pi, d) = \mathbb{E}[L(\theta, d) | x] = \int_{\Theta} L(\theta, d)\pi(\theta | x)d\theta.$$

- Compute the **Bayesian/average risk** to average out x as well as θ , defined as

$$R(\pi, d) = \mathbb{E}_{x,\theta}[L(\theta, d)] = \int_{\Theta} \int_{\mathcal{X}} L(\theta, d(x))f(x | \theta)\pi(\theta)dx d\theta.$$

The best decision is chosen by minimizing the bayesian risk.

Remark 1.1. Recall by Bayes' formula, we have

$$f(x, \theta) = f(x | \theta)f(\theta) = f(\theta | x)f(x).$$

And by tower property, we have

$$\mathbb{E}(L(\theta, d)) = \mathbb{E}_\theta \mathbb{E}_{x|\theta}(L(\theta, d)|\theta) = \mathbb{E}_x \mathbb{E}_{\theta|x}(L(\theta, d)|x).$$

The dependency on the data and the parameter is summarized below.

	Depend on x	Not depend on x
Depend on θ	$L(\theta, d)$	$R(\theta, d)$
Not depend on θ	$L(\pi, d)$	$R(\pi, d)$

Example 1.1. (Example of Frequentist and Bayesian procedure) Consider in a coin flipping experiment, we want to know the the probability of getting a head or θ . We have two observations from the following model,

$$x_{1:2} \stackrel{\text{iid}}{\sim} \text{Bern}(\theta).$$

We agree to pick the best estimator according to the absolute difference. There are two candidates of estimator to choose,

$$d_1(x) = \frac{1}{2}, \quad \text{and} \quad d_2(x) = \bar{X} = n^{-1} \sum_{i=1}^n x_i.$$

1. Describe the observation, parameter and decision space.
2. Describe the utility function.
3. Describe the frequentist procedure when θ is 0.3 and 0.5.
4. Describe the bayesian procedure, choose a conjugate prior for θ .

SOLUTION:

1. We have

$$\mathcal{X} = \{0, 1\}^2 = \{00, 10, 01, 11\}, \quad \Theta = [0, 1], \quad \mathcal{D} = [0, 1].$$

2. The utility function is $U(d, \theta) = -|d - \theta|$, the closer d to θ , the larger the utility, which means the better d is.
3. The loss function is $L(\theta, d) = -U(d, \theta) = |d - \theta|$. The risk is

$$R(\theta, d) = \mathbb{E}[L(\theta, d) | \theta] = \sum_{x_{1:2} \in \mathcal{X}} |d - \theta| f(x_{1:2} | \theta) = \sum_{x_{1:2} \in \mathcal{X}} |d - \theta| \theta^{S_2} (1 - \theta)^{n - S_2}, \quad (1.1)$$

where $n = 2$ and $S_n = \sum_{i=1}^n x_i$. For d_1 , we have

$$R(\theta = 0.5, d_1) = 0, \quad \text{and} \quad R(\theta = 0.3, d_1) = 0.2.$$

For d_2 , we have

$$R(\theta = 0.5, d_2) = 0.25, \quad \text{and} \quad R(\theta = 0.3, d_2) = 0.294.$$

For both values of θ , we choose d_1 .

```

1 ##compute frequentist risk
2 get_joint_density <- function(theta=0.5){ ##c(f(00),f(10),f(01),f(11))
3   cx = c(0,1,1,2)
4   den = (theta^sx)*(1-theta)^(2-sx)
5   return(den)
6 }
7 get_risk = function(d,theta=0.5){ ##d = c(d(00),d(10),d(01),d(11))
8   den = get_joint_density(theta)
9   risk = sum(abs(theta-d)*den)
10  return(risk)
11 }
12 ##compute risk of d1
13 d1_risk_05 = get_risk(c(0.5,0.5,0.5,0.5),theta = 0.5)
14 d1_risk_03 = get_risk(c(0.5,0.5,0.5,0.5),theta = 0.3)
15 d1_risk_05;d1_risk_03
16 ##compute risk of d2
17 d2_risk_05 = get_risk(c(0,0.5,0.5,1),theta = 0.5)
18 d2_risk_03 = get_risk(c(0,0.5,0.5,1),theta = 0.3)
19 d2_risk_05;d2_risk_03

```

4. The conjugate prior for θ is $\text{Beta}(\alpha, \beta)$. Choose $\alpha = 1/2$ and $\beta = 1/2$. Then, the posterior is $\text{Beta}(\alpha + n\bar{x}, \beta + n(1 - \bar{x}))$. We can simulate the posterior loss and bayesian risk.

x	00	10	01	11
$L(\pi, d_1)$	0.353	0.207	0.210	0.354
$L(\pi, d_2)$	0.161	0.207	0.210	0.170

Also, $R(\pi, d_1) = 0.317$ and $R(\pi, d_2) = 0.182$. We choose d_2 accordingly.

```

1 ##compute posterior loss
2 get_posterior_loss <- function(x,seed = 4010,a=0.5,b=0.5,d1 = 0.5,nRep =
3   2^10){
4   bar_x = mean(x) ##d2 is bar_x
5   set.seed(seed)
6   theta = rbeta(nRep,a+n*bar_x,b+n*(1-bar_x))
7   p_loss = array(NA,2)
8   p_loss[1] = mean(abs(theta-d1)) ##posterior loss for d1
9   p_loss[2] = mean(abs(theta-bar_x)) ##posterior loss for d2
10  return(p_loss)
11 }
12 get_posterior_loss(x = c(0,0),seed=4010) ##posterior loss when x = 00
13 get_posterior_loss(x = c(1,0),seed=4011) ##posterior loss when x = 10
14 get_posterior_loss(x = c(0,1),seed=4012) ##posterior loss when x = 01
15 get_posterior_loss(x = c(1,1),seed=4013) ##posterior loss when x = 11
16
17 ##compute bayesian risk by MC

```

```

18 a=0.5
19 b=0.5
20 d1 = 0.5
21 nRep = 2^10
22
23 set.seed(4010)
24 theta = rbeta(nRep, a, b)
25 x = array(NA, c(nRep, 2))
26 for (i in 1:nRep) {
27   x[i,] = rbinom(n=2, size=1, theta[i])
28 }
29 d2 = apply(x, 1, mean)
30 mean(abs(d1-theta)) ##bayesian risk for d1
31 mean(abs(d2-theta)) ##bayesian risk for d2

```

2 Bayes estimator

Definition 2. (*Bayes Estimator*)

1. The Bayes estimator under the prior $\theta \sim \pi(\theta)$ and the loss $L(\theta, \hat{\theta})$ is defined as

$$\hat{\theta}_\pi = \arg \min_{\hat{\theta}} R(\pi, \hat{\theta}).$$

2. The value $R(\pi, \hat{\theta}_\pi)$ is called the Bayes risk.

Theorem 2.1. (*A simpler way of finding Bayes estimators*). Assume there is an estimator $\hat{\theta}_0$ such that $R(\theta, \hat{\theta}_0) < \infty$. Then, for each $x \in X$, the Bayes estimator is

$$\hat{\theta}_\pi(x) = \arg \min_{\hat{\theta}} L(\pi, \hat{\theta})$$

Theorem 2.2. (*Bayes estimators under special losses*). Suppose the conditions in Theorem 2.1 hold.

1. Quadratic loss. Let $L(\theta, \hat{\theta}) = (\hat{\theta} - \theta)^2$. (See Example 3.13 for a generalization.) The Bayes estimator is

$$\hat{\theta}_\pi = E(\theta | x).$$

2. Absolute loss. Let $L(\theta, \hat{\theta}) = |\hat{\theta} - \theta|$. (See Example 3.14 for a generalization.) The Bayes estimator is

$$\hat{\theta}_\pi = Q_{1/2}(\theta | x).$$

3. Zero-one loss. Let $L(\theta, \hat{\theta}) = \mathbb{1}(|\hat{\theta} - \theta| > \varepsilon)$ for some $\varepsilon > 0$. The Bayes estimator is

$$\hat{\theta}_\pi = \arg \max_{\hat{\theta}} P(\theta \in [\hat{\theta} - \varepsilon, \hat{\theta} + \varepsilon] | x).$$

Example 2.1. (Q4 Midterm 2019 - Bayes Estimator (20%)). Let Θ be the parameter space. Suppose that the decision space is $\mathcal{D} = \Theta$. Consider the generalized squared-error loss function $L : \Theta \times \mathcal{D} \rightarrow \mathbb{R}$ defined as

$$L(\theta, \hat{\theta}) = \{g(\hat{\theta}) - g(\theta)\}^2 \quad (2.1)$$

where $g : \Theta \rightarrow \mathbb{R}$ is a continuous and strictly increasing function whose inverse function g^{-1} exists. (Note: g may NOT be differentiable.) Let $\theta \sim \pi(\theta)$ such that the Bayesian risk exists.

1. Find the Bayes estimator $\hat{\theta}_\pi$ of θ .
2. Let $W \sim N(\mu, \sigma^2)$. Prove that

$$E(e^W) = \exp\left\{\mu + \frac{\sigma^2}{2}\right\}$$

3. For parts (3) and (4), consider $g(t) = e^t$, $[x_1, \dots, x_n | \theta] \stackrel{\text{ID}}{\sim} N(\theta, \sigma^2)$ and $\theta \sim N(\theta_0, \tau_0^2)$. Prove that the Bayes estimator of θ under the loss (2.1) is $\hat{\theta}_\pi = \theta_n + \frac{\tau_n^2}{2}$, where

$$\theta_n = \frac{\tau_0^2 \bar{x}_n + \sigma^2 \theta_0 / n}{\tau_0^2 + \sigma^2 / n}, \quad \tau_n^2 = \frac{\tau_0^2 \sigma^2 / n}{\tau_0^2 + \sigma^2 / n}, \quad \bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i. \quad (2.2)$$

4. Comment the following statement: “Under the squared-error loss, the Bayes estimator of θ is $\hat{\theta}'_\pi = \theta_n$, which is different from the $\hat{\theta}_\pi$ in (2.2). In practice, we should choose the estimator which has a smaller mean squared error.”

SOLUTION:

1. There are two methods.

- (Method 1) The posterior loss is

$$\begin{aligned} L(\pi, \hat{\theta}) &= E\{L(\theta, \hat{\theta}) | x\} \\ &= E\left[\{g(\hat{\theta}) - E\{g(\theta) | x\} + E\{g(\theta) | x\} - g(\theta)\}^2 | x\right] \\ &= \left[g(\hat{\theta}) - E\{g(\theta) | x\}\right]^2 + E\left(\left[E\{g(\theta) | x\} - g(\theta)\right]^2 | x\right) \\ &\quad + 2 \left[g(\hat{\theta}) - E\{g(\theta) | x\}\right] \underbrace{E\left[E\{g(\theta) | x\} - g(\theta) | x\right]}_{=0} \\ &= \underbrace{\left[g(\hat{\theta}) - E\{g(\theta) | x\}\right]^2}_{\star} + \underbrace{E\left(\left[E\{g(\theta) | x\} - g(\theta)\right]^2 | x\right)}_{\text{does not depends on } \hat{\theta}} \end{aligned}$$

Hence, $L(\pi, \hat{\theta})$ is minimized (with respect to $\hat{\theta}$) iff \star is minimized. Clearly, the minimum is achieved iff $g(\hat{\theta}) - E\{g(\theta) | x\} = 0$. Solving for $\hat{\theta}$, we know that the Bayes estimator is

$$\hat{\theta}_\pi = g^{-1}(E\{g(\theta) | x\})$$

- (Method 2) We first find the minimizer with respect to $\hat{\phi} = g(\hat{\theta})$. Note that the posterior loss and its derivatives are

$$\begin{aligned} L(\pi, \hat{\theta}) &= \hat{\phi}^2 - 2\hat{\phi}\mathbb{E}\{g(\theta) \mid x\} + \mathbb{E}\{g^2(\theta) \mid x\} \\ \frac{d}{d\hat{\phi}}L(\pi, \hat{\theta}) &= 2\hat{\phi} - 2\mathbb{E}\{g(\theta) \mid x\} \\ \frac{d^2}{d\hat{\phi}^2}L(\pi, \hat{\theta}) &= 2 > 0 \end{aligned}$$

Clearly, $\hat{\phi} = \mathbb{E}\{g(\theta) \mid x\}$ minimizes $L(\pi, \hat{\theta})$. Since g is a strictly increasing function, we know that $\hat{\theta} = g^{-1}(\mathbb{E}\{g(\theta) \mid x\})$ minimizes $L(\pi, \hat{\theta})$. Hence, the Bayes estimator is

$$\hat{\theta}_\pi = g^{-1}(\mathbb{E}\{g(\theta) \mid x\}).$$

2. Represent $W = \mu + \sigma Z$, where $Z \sim N(0, 1)$. So, by matching the Normal density, we have

$$\begin{aligned} \mathbb{E}(e^W) &= \mathbb{E}(e^{\mu + \sigma Z}) = e^\mu \mathbb{E}(e^{\sigma Z}) \\ &= e^\mu \int_{-\infty}^{\infty} e^{\sigma z} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz \\ &= e^\mu e^{\sigma^2/2} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(z - \sigma)^2\right\} dz \\ &= \exp\left\{\mu + \frac{\sigma^2}{2}\right\}. \end{aligned}$$

3. Recall that $[\theta \mid x_{1:n}] \sim N(\theta_n, \tau_n^2)$, where θ_n and τ_n^2 are given by (2.2). By the result of parts (1) and (2), we have

$$\hat{\theta}_\pi = \log \mathbb{E}(e^\theta \mid x_{1:n}) = \log \exp\left\{\theta_n + \frac{\tau_n^2}{2}\right\} = \theta_n + \frac{\tau_n^2}{2}$$

4. It is true that these two Bayes estimators ($\hat{\theta}_\pi$ and $\hat{\theta}'_\pi$) are different because they are evaluated under different loss functions. However, there is no absolute rule for comparing these two estimators.

- On one side, using the mean-squared error (MSE) is possible. Recall that MSE is the frequentist risk under the squared-error loss. In this case, we compare the frequentist's properties (i.e., MSE) of two Bayesian estimators. The MSE is useful for understanding, e.g., consistency and the L^2 -convergence rate. So, it is not wrong to do so.
- But, on the other side, comparing the estimators by MSE may not be relevant or fair from the Bayesian point of view. For Bayesian, once we fix our loss function (e.g., (4.2) or the squared-error loss), the goal is to minimize the Bayesian risk but not other frequentist risk (e.g., the MSE). Hence, "evaluating an estimator that optimizes criteria A by another criteria B " may not be fair. It is possible in practice because the loss function $L(\theta, \hat{\theta})$ may be selected due to some problem-specific purposes, e.g., handling asymmetric loss.

Example 2.2. (Q4 Midterm 2021- Bayes Estimator). Let $y_i \in \mathbb{R}$ be the response value at time i , and $x_i \in \{0, 1\}$ be a binary covariate. (For example, y_i and x_i may denote the rate of return a stock and the market condition, respectively.) Suppose that $y_0 = 0$ and $x_{1:n}$ are non-random and always known. Consider a regime-switching AR (1) model:

$$[y_i | \theta_1, \theta_0, y_{i-1}] \sim N(\theta_1^{x_i} \theta_0^{1-x_i} y_{i-1}, 1), \quad i = 1, \dots, n,$$

$$\theta_1, \theta_0 \stackrel{\text{iid}}{\sim} \text{Unif}(-1, 1).$$

1. Derive the posterior distribution of $\theta = (\theta_0, \theta_1)^\top$.
2. Consider the loss function:

$$L(\theta, \hat{\theta}) = \sum_{i=1}^n \left(\theta_1^{x_i} \theta_0^{1-x_i} - \hat{\theta}_1^{x_i} \hat{\theta}_0^{1-x_i} \right)^2$$

where $\hat{\theta} = (\hat{\theta}_0, \hat{\theta}_1)^\top$.

- (a) Interpret the loss function.
- (b) Derive and interpret the Bayes estimator.

Use no more than about 50 words for the interpretations in (a) and (b).

SOLUTION:

1. Let $n_1 = \sum_{i=1}^n x_i$ be the number of i such that $x_i = 1$. Then $n_0 = n - n_1$ is the number of i such that $x_i = 0$. Denote all the indexes $i \in \{1, 2, \dots, n\}$ satisfying $x_i = 1$ to be j_1, j_2, \dots, j_{n_1} . Similarly, denote all the indexes $i \in \{1, 2, \dots, n\}$ satisfying $x_i = 0$ to be k_1, k_2, \dots, k_{n_0} . The sampling distribution can be written as

$$\begin{aligned} f(y_{1:n} | \theta_1, \theta_0) &= f(y_{1:n} | \theta_1, \theta_0, y_0) \\ &= \prod_{i=1}^n f(y_i | \theta_1, \theta_0, y_{i-1}) \\ &= (2\pi)^{-n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (y_i - \theta_1^{x_i} \theta_0^{1-x_i} y_{i-1})^2 \right\} \\ &= (2\pi)^{-n/2} \exp \left\{ -\frac{1}{2} \sum_{s=1}^{n_1} (y_{j_s} - \theta_1 y_{j_s-1})^2 \right\} \exp \left\{ -\frac{1}{2} \sum_{t=1}^{n_0} (y_{k_t} - \theta_0 y_{k_t-1})^2 \right\} \end{aligned}$$

The posterior is

$$\begin{aligned} f(\theta_1, \theta_0 | y_{1:n}) &\propto f(y_{1:n} | \theta_1, \theta_0) f(\theta_1, \theta_0) \\ &\propto \exp \left\{ -\frac{1}{2} \sum_{s=1}^{n_1} (y_{j_s} - \theta_1 y_{j_s-1})^2 \right\} \exp \left\{ -\frac{1}{2} \sum_{t=1}^{n_0} (y_{k_t} - \theta_0 y_{k_t-1})^2 \right\} \mathbb{1}(-1 < \theta_1, \theta_0 < 1) \\ &\propto \exp \left[-\frac{1}{2} \left\{ \left(\sum_{s=1}^{n_1} y_{j_s-1}^2 \right) \theta_1^2 - 2 \left(\sum_{s=1}^{n_1} y_{j_s} y_{j_s-1} \right) \theta_1 \right\} \right] \\ &\times \exp \left[-\frac{1}{2} \left\{ \left(\sum_{t=1}^{n_0} y_{k_t-1}^2 \right) \theta_0^2 - 2 \left(\sum_{t=1}^{n_0} y_{k_t} y_{k_t-1} \right) \theta_0 \right\} \right] \mathbb{1}(-1 < \theta_1, \theta_0 < 1). \end{aligned}$$

Denote

$$\mu_0 = \frac{\sum_{k=1}^{n_0} y_{k_m} y_{k_m-1}}{\sum_{k=1}^{n_0} y_{k_m-1}^2}, \quad \sigma_0 = \frac{1}{\sqrt{\sum_{k=1}^{n_0} y_{k_m-1}^2}}, \quad \mu_1 = \frac{\sum_{l=1}^{n_1} y_{j_l} y_{j_l-1}}{\sum_{l=1}^{n_1} y_{j_l-1}^2} \quad \text{and} \quad \sigma_1 = \frac{1}{\sqrt{\sum_{l=1}^{n_1} y_{j_l-1}^2}}.$$

Then

$$\begin{aligned} f(\theta_1, \theta_0 \mid y_{1:n}) &\propto \text{dnorm}(\theta_1, \mu_1, \sigma_1) \text{dnorm}(\theta_0, \mu_0, \sigma_0) \mathbb{1}(-1 < \theta_1, \theta_0 < 1) \\ &= \text{dTN}(\theta_1, \mu_1, \sigma_1, -1, 1) \text{dTN}(\theta_0, \mu_0, \sigma_0, -1, 1). \end{aligned}$$

2. Observe that the loss function can be written as

$$\begin{aligned} L(\theta, \hat{\theta}) &= \sum_{i=1}^n \left(\theta_1^{x_i} \theta_0^{1-x_i} - \hat{\theta}_1^{x_i} \hat{\theta}_0^{1-x_i} \right)^2 \\ &= n_1 \left(\theta_1 - \hat{\theta}_1 \right)^2 + n_0 \left(\theta_0 - \hat{\theta}_0 \right)^2 \end{aligned}$$

Since n_1 and n_0 are given, we can find the minimizer of the \mathcal{L}^2 loss of θ_0 and of θ_1 respectively. By Theorem 3.2 and Example 3.13 from the lecture notes, the Bayes estimator of θ is

$$\hat{\theta} = \left(\hat{\theta}_1^\pi, \hat{\theta}_0^\pi \right)^T = \left(\text{E}(\theta_1 \mid y_{1:n}), \text{E}(\theta_0 \mid y_{1:n}) \right)^T.$$

Specifically,

$$\hat{\theta}_j^\pi = \text{E}(\theta_j \mid y_{1:n}) = \mu_j - \sigma_j \frac{\text{dnorm}(u_j) - \text{dnorm}(l_j)}{\text{pnorm}(u_j) - \text{pnorm}(l_j)},$$

where $u_j = (1 - \mu_j) / \sigma_j$ and $l_j = (-1 - \mu_j) / \sigma_j$ for $j = 0, 1$. For the interpretations, any reasonable answers are accepted.