

# **Introduction to Machine Learning: Homework II**

Due on Oct 26th, 2022 at 11:59pm

*Professor Ziping Zhao*

**Bingnan Li**  
2020533092

## 1. [Bayesian Decision Theory]

- (a) Specify decision rule via likelihood and posterior respectively, and give example with  $x = \text{light}$  to use the rules.

**Solution:**

Likelihood decision rule:

$$\text{choose} \begin{cases} C = C_1, & \text{if } P(x|C_1) > P(x|C_2) \\ C = C_2, & \text{elsewhere} \end{cases}$$

Posterior decision rule:

$$\text{choose} \begin{cases} C = C_1, & \text{if } P(C_1|x) > P(C_2|x) \\ C = C_2, & \text{elsewhere} \end{cases}$$

If  $x = \text{light}$ , then likelihood is

$$P(x = \text{light}|C_1) = \frac{2125}{2125 + 6375} = 0.25$$

$$P(x = \text{light}|C_2) = \frac{1000}{1000 + 500} = 0.667$$

Thus, by likelihood decision rule, we choose  $C_2$ . However, the posterior probability is

$$P(C_1|x = \text{light}) = \frac{2125}{2125 + 1000} = 0.68$$

$$P(C_2|x = \text{light}) = \frac{1000}{2125 + 1000} = 0.32$$

Thus, by posterior decision rule, we choose  $C_1$ .

- (b) (i) Find the optimal decision rule that will give the minimum expected loss.

**Solution:**

Let action  $\alpha_i$  be classifying input  $\mathbf{x}$  into class  $i$  and  $R(\alpha_i|\mathbf{x})$  be the expected loss of taking action  $\alpha_i$ , then we have

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^K \lambda_{i,j} P(C_j|\mathbf{x})$$

Then the optimal decision rule is

$$\begin{cases} \text{choose class } i, & \text{if } R(\alpha_i|\mathbf{x}) = \min\{R(\alpha_j|\mathbf{x})\}_{j=1}^K < \lambda \\ \text{reject } \mathbf{x}, & \text{if } R(\alpha_i|\mathbf{x}) = \min\{R(\alpha_j|\mathbf{x})\}_{j=1}^K \geq \lambda \end{cases}$$

- (ii) State the relationship between  $\lambda$  and the rejection threshold  $\theta$ .

**Solution:**

Given that  $\lambda_{i,k} = \begin{cases} 1, & \text{if } i \neq k \\ 0, & \text{if } i = k \end{cases}$  Then the corresponding expected loss is as following

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^K \lambda_{i,j} P(C_j|\mathbf{x}) = 1 - P(C_i|\mathbf{x})$$

In order to get minimum expected loss, we hope:

- When  $P(C_i|\mathbf{x}) > \theta$ , we hope that  $R(\alpha_i|\mathbf{x}) = \min\{R(\alpha_j|\mathbf{x})\}_{j=1}^K = 1 - P(C_i|\mathbf{x}) \leq \lambda$  so that we won't get expected loss any lower by manually choosing to reject  $\mathbf{x}$ .

In order to achieve that, we need:

$$\lambda \geq \max(1 - P(C_i|\mathbf{x})) = 1 - \theta$$

- When  $P(C_i|\mathbf{x}) \leq \theta$ , we hope that  $\lambda \leq R(\alpha_i|\mathbf{x}) = \min\{R(\alpha_j|\mathbf{x})\}_{j=1}^K = 1 - P(C_i|\mathbf{x})$  so that we won't get expected loss any lower by selecting to classifying  $\mathbf{x}$  into class  $i$ .

In order to achieve that, we need

$$\lambda \leq \min(1 - P(C_i|\mathbf{x})) = 1 - \theta$$

Based on that, we have

$$\lambda \geq 1 - \theta$$

$$\lambda \leq 1 - \theta$$

Therefore, we have

$$\lambda = 1 - \theta \Rightarrow \lambda + \theta = 1$$

## 2. [Parameter Estimation]

- (a) Let  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$  and  $\mathbf{y} = [y_1, \dots, y_n]^T$ . Assume that  $e_i \sim \mathcal{N}(0, \sigma^2)$  and  $\mathbf{X}^T \mathbf{X}$  is a full-rank matrix. Derive the maximum likelihood estimate  $\mathbf{w}_{ML} = (\mathbf{X}^T \mathbf{X}^{-1} \mathbf{X}^T \mathbf{y})$ .

**Proof:**

Given that  $e_i \sim \mathcal{N}(0, \sigma^2)$  and  $y_i = \mathbf{w}^T \mathbf{x}_i + e_i$ , we have that

$$y_i \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}_i, \sigma^2)$$

Thus, the likelihood function is

$$\begin{aligned} \mathcal{L}(\mathbf{w}) &= \log(P(\mathbf{y}|\mathbf{X}, \mathbf{w})) = - \sum_{i=1}^n \left( \log \sqrt{2\pi\sigma} + \frac{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2} \right) \\ &= -n \log \sqrt{2\pi\sigma} - \frac{\sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2} \\ &= -n \log \sqrt{2\pi\sigma} - \frac{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2}{2\sigma^2} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= -\frac{1}{2\sigma^2} \frac{\partial \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}}{\partial \mathbf{w}} \\ &= -\frac{1}{2\sigma^2} (2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y}) = 0 \\ &\Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \\ &\Rightarrow \mathbf{w}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

- (b) Assume  $\mathbf{w} \sim \mathcal{N}(0, \nu^2 \mathbf{I})$ . Derive the maximum a posteriori estimate

$$\mathbf{w}_{MAP} = \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\nu^2} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

**Proof:**

Given that  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \nu^2 \mathbf{I})$ , we know that

$$P(\mathbf{w}) = \frac{1}{(2\pi)^{n/2} \nu} \exp\left[-\frac{1}{2}(\mathbf{w}^T \frac{1}{\nu^2} \mathbf{I} \mathbf{w})\right] = \frac{1}{(2\pi)^{n/2} \nu} \exp\left(-\frac{\mathbf{w}^T \mathbf{w}}{2\nu^2}\right)$$

Then by definition, we know that

$$\begin{aligned}
 \mathbf{w}_{MAP} &= \arg \max_{\mathbf{w}} P(\mathbf{w}|\mathcal{D}) \\
 &= \arg \max_{\mathbf{w}} \frac{P(\mathcal{D}|\mathbf{w})P(\mathbf{w})}{P(\mathcal{D})} \\
 &= \arg \max_{\mathbf{w}} P(\mathbf{y}, \mathbf{X}|\mathbf{w})P(\mathbf{w}) \\
 &= \arg \max_{\mathbf{w}} P(\mathbf{y}|\mathbf{X}, \mathbf{w})P(\mathbf{X}|\mathbf{w})P(\mathbf{w})
 \end{aligned}$$

Since the value of  $\mathbf{X}$  does not rely on  $\mathbf{w}$ , we have  $P(\mathbf{X}|\mathbf{w}) = P(\mathbf{X})$ . Sequentially, we have

$$\mathbf{w}_{MAP} = \arg \max_{\mathbf{w}} P(\mathbf{y}|\mathbf{X}, \mathbf{w})P(\mathbf{X})P(\mathbf{w}) = \arg \max_{\mathbf{w}} P(\mathbf{y}|\mathbf{X}, \mathbf{w})P(\mathbf{w})$$

Next, let  $\mathcal{L}(\mathbf{w}) = \log P(\mathbf{y}|\mathbf{X}, \mathbf{w})P(\mathbf{w}) = \log P(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \log P(\mathbf{w})$ , we have

$$\begin{aligned}
 \mathcal{L}(\mathbf{w}) &= -n \log \sqrt{2\pi}\sigma - \frac{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2}{2\sigma^2} + \log (2\pi)^{n/2}\nu - \frac{\mathbf{w}^T \mathbf{w}}{2\nu^2} \\
 \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= -\frac{\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}}{\sigma^2} - \frac{\mathbf{w}}{\nu^2} = 0 \\
 &\Rightarrow \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\nu^2} \mathbf{I} \right) \mathbf{w} = \mathbf{X}^T \mathbf{y}
 \end{aligned}$$

Given that  $\mathbf{X}^T \mathbf{X}$  is full-rank and  $\frac{\sigma^2}{\nu^2} \mathbf{I}$  is positive-definite, we know that  $\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\nu^2} \mathbf{I}$  is invertible.

Therefore, we know that

$$\mathbf{w}_{MAP} = \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\nu^2} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

### 3. [Parameter Estimation]

**Solution:**

Let  $\mathbf{z}_l = \text{sigmoid}(\mathbf{W}_1 \mathbf{x}_l + \mathbf{b}_1)$  and  $\mathbf{f}_l = \text{softmax}(\mathbf{W}_2 \mathbf{z}_l + \mathbf{b}_2)$ , then

$$\begin{aligned}
 \ell(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2) &= -\sum_{l=1}^L \mathbf{y}_l^T \log \mathbf{f}_l \\
 &= -\sum_{l=1}^L \sum_{k=1}^M y_{l,k} \log f_{l,k}
 \end{aligned}$$

where  $f_{l,k} = \text{softmax}(\mathbf{w}_{2,k}^T \mathbf{z}_l + b_{2,k})$ . Then, we have

$$\begin{aligned}
 \frac{\partial \ell}{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}} &= \frac{\partial \ell}{\partial f_{l,k}} \frac{\partial f_{l,k}}{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}} \\
 &= -\sum_{l=1}^L \sum_{k=1}^M \frac{y_{l,k}}{f_{l,k}} f_{l,k} (\delta_{ki} - f_{l,i}) \\
 &= -\sum_{l=1}^L \sum_{k=1}^M y_{l,k} (\delta_{ki} - f_{l,i}) \\
 &= -\sum_{l=1}^L (y_{l,i} - f_{l,i}) \\
 \text{where } \delta_{ki} &= \begin{cases} 1, & k = i \\ 0, & k \neq i \end{cases}
 \end{aligned}$$

Since

$$\begin{aligned}
\frac{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}}{\partial w_{2,i,j}} &= z_{l,j} \\
\frac{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}}{\partial b_{2,i}} &= 1 \\
\frac{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}}{\partial w_{1,p,q}} &= \sum_{j=1}^H w_{2,i,j} \frac{\partial z_{l,j}}{\partial w_{1,p,q}} \\
\frac{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}}{\partial b_{1,p}} &= \sum_{j=1}^H w_{2,i,j} \frac{\partial z_{l,j}}{\partial b_{1,p}} \\
\frac{\partial z_{l,j}}{\partial w_{1,p,q}} &= \begin{cases} z_{l,q}(1 - z_{l,q})x_{l,q}, & j = p \\ 0, & j \neq p \end{cases} \\
\frac{\partial z_{l,j}}{\partial b_{1,p}} &= \begin{cases} z_{l,q}(1 - z_{l,q}), & j = p \\ 0, & j \neq p \end{cases}
\end{aligned}$$

Where  $z_{l,q} = \text{sigmoid}(\mathbf{w}_{1,q}^T \mathbf{x}_l + b_{1,q})$ .

Then, we have:

$$\begin{aligned}
\Delta w_{2,i,j} &= -\eta \frac{\partial \ell}{\partial w_{2,i,j}} = -\eta \frac{\partial \ell}{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}} \frac{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}}{\partial w_{2,i,j}} \\
&= \eta \sum_{l=1}^L (y_{l,i} - f_{l,i}) z_{l,j} \\
\Delta b_{2,i} &= -\eta \frac{\partial \ell}{\partial b_{2,i}} = -\eta \frac{\partial \ell}{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}} \frac{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}}{\partial b_{2,i}} \\
&= \eta \sum_{l=1}^L (y_{l,i} - f_{l,i}) \\
\Delta w_{1,p,q} &= -\eta \frac{\partial \ell}{\partial w_{1,p,q}} = -\eta \sum_{i=1}^M \frac{\partial \ell}{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}} \frac{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}}{\partial w_{1,p,q}} \\
&= \eta \sum_{l=1}^L \left[ \sum_{i=1}^M (y_{l,i} - f_{l,i}) w_{2,i,p} \right] z_{l,p}(1 - z_{l,p})x_{l,q} \\
\Delta b_{1,p} &= -\eta \frac{\partial \ell}{\partial b_{1,p}} = -\eta \sum_{i=1}^M \frac{\partial \ell}{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}} \frac{\partial \mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}}{\partial b_{1,p}} \\
&= \eta \sum_{l=1}^L \left[ \sum_{i=1}^M (y_{l,i} - f_{l,i}) w_{2,i,p} \right] z_{l,p}(1 - z_{l,p})
\end{aligned}$$

Where

$$\begin{aligned}
f_{l,i} &= \text{softmax}(\mathbf{w}_{2,i}^T \mathbf{z}_l + b_{2,i}) \\
z_{l,p} &= \text{sigmoid}(\mathbf{w}_{1,q}^T \mathbf{x}_l + b_{1,q})
\end{aligned}$$

## 4. [Bayesian Decision Theory, Linear Discrimination]

- (a) minimize the expected loss v.s. minimize the misclassification error.

**Solution:**

Let  $\lambda_{ij}$  be the loss for misclassifying  $\mathbf{x}$  into class  $i$  while it actually belongs to class  $j$ ,  $\alpha_i$  be the action of classifying  $\mathbf{x}$  into class  $i$  and  $R(\alpha_i|\mathbf{x})$  be the expected loss of classifying  $\mathbf{x}$  into class  $i$ .

Therefore,

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^2 \lambda_{ij} P(C_j|\mathbf{x})$$

Let  $g(\mathbf{x}) = R(\alpha_1|\mathbf{x}) - R(\alpha_2|\mathbf{x})$ , in order to minimize the expected loss, we have the following decision rule:

$$\begin{cases} \text{choose } C_1, \text{ if } g(\mathbf{x}) < 0 \\ \text{choose } C_2, \text{ elsewhere} \end{cases}$$

Then the decision boundary is the solution of  $g(\mathbf{x}) = 0$ .

$$\begin{aligned} g(\mathbf{x}) &= 0 \\ \Downarrow \\ (\lambda_{11} - \lambda_{21}) \frac{P(\mathbf{x}|C_1)P(C_1)}{P(\mathbf{x})} \\ &= (\lambda_{22} - \lambda_{12}) \frac{P(\mathbf{x}|C_2)P(C_2)}{P(\mathbf{x})} \\ \Downarrow \\ \log(|\lambda_{11} - \lambda_{21}|) + \log P(\mathbf{x}|C_1) + \log P(C_1) \\ &= \log(|\lambda_{22} - \lambda_{12}|) + \log P(\mathbf{x}|C_2) + \log P(C_2) \\ \Downarrow \\ \log(|\lambda_{11} - \lambda_{21}|) + \frac{1}{2\sigma^2}(\mathbf{2}\boldsymbol{\mu}_1^T \mathbf{x} - \boldsymbol{\mu}_1^T \boldsymbol{\mu}_1) + \log P(C_1) \\ &= \log(|\lambda_{22} - \lambda_{12}|) + \frac{1}{2\sigma^2}(\mathbf{2}\boldsymbol{\mu}_2^T \mathbf{x} - \boldsymbol{\mu}_2^T \boldsymbol{\mu}_2) + \log P(C_2) \\ \Downarrow \\ \frac{1}{\sigma^2}(\boldsymbol{\mu}_1^T - \boldsymbol{\mu}_2^T)\mathbf{x} &= \log \frac{(\lambda_{22} - \lambda_{12})P(C_2)}{(\lambda_{11} - \lambda_{21})P(C_1)} + \frac{1}{2\sigma^2}(\boldsymbol{\mu}_1^T \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\mu}_2) \end{aligned}$$

Define  $\mathbf{w} = \frac{1}{\sigma^2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$  and  $w_0 = \log \frac{(\lambda_{11} - \lambda_{21})P(C_1)}{(\lambda_{22} - \lambda_{12})P(C_2)} - \frac{1}{2\sigma^2}(\boldsymbol{\mu}_1^T \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\mu}_2)$ , we have

$$\mathbf{w}^T \mathbf{x} + w_0 = \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0) = 0$$

where  $\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2) - \frac{1}{\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2} \sigma^2 \log \frac{(\lambda_{11} - \lambda_{21})P(C_1)}{(\lambda_{22} - \lambda_{12})P(C_2)} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$

However, in order to minimize the probability of error, we define

$$\begin{aligned}
 g'(\mathbf{x}) &= \log P(C_1|\mathbf{x}) - \log P(C_2|\mathbf{x}) \\
 &= \log P(\mathbf{x}|C_1) - \log P(\mathbf{x}|C_2) + \log \frac{P(C_1)}{P(C_2)} \\
 &= \frac{1}{\sigma^2}(\boldsymbol{\mu}_1^T - \boldsymbol{\mu}_2^T)\mathbf{x} - \frac{1}{2\sigma^2}(\boldsymbol{\mu}_1^T \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\mu}_2) + \log \frac{P(C_1)}{P(C_2)} \\
 &= \mathbf{w}'^T \mathbf{x} + w_0' = \mathbf{w}'^T (\mathbf{x} - \mathbf{x}_0') = 0 \\
 \Rightarrow \mathbf{x}_0' &= \frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2) - \frac{1}{\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2} \sigma^2 \log \frac{P(C_1)}{P(C_2)} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)
 \end{aligned}$$

The corresponding decision rule is:

$$\begin{cases} \text{choose } C_1, \text{ if } g'(\mathbf{x}) > 0 \\ \text{choose } C_2, \text{ elsewhere} \end{cases}$$

Based on that, we know

$$\begin{aligned}
 \mathbf{w} &= \mathbf{w}' \\
 x_0 &= x_0' - \frac{1}{\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2} \sigma^2 \log \frac{\lambda_{11} - \lambda_{21}}{\lambda_{22} - \lambda_{12}} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)
 \end{aligned}$$

Therefore, two decision boundaries have the same normal vector but different bias.

- (b) By minimizing the misclassification error, obtain and draw the decision boundary when  $\mu_{11} = 1, \mu_{12} = 1, \mu_{21} = 3, \mu_{22} = 5, \sigma = 1$  and  $P(C_1) = P(C_2)$ .

**Solution:**

By Bayesian rule, we have

$$P(C_i|\mathbf{x}) = \frac{P(\mathbf{x}|C_i)P(C_i)}{P(\mathbf{x})}$$

Given that  $x_1, x_2$  are independent following the Laplace distribution, we have

$$P(C_i|\mathbf{x}) = \frac{P(x_1|C_i)P(x_2|C_i)P(C_i)}{P(\mathbf{x})}$$

Then define discrimination function

$$\begin{aligned}
 g_i(\mathbf{x}) &= \log P(C_i|\mathbf{x}) \\
 &= \log P(x_1|C_i) + \log P(x_2|C_i) + \log P(C_i) - \log P(\mathbf{x}) \\
 &= -2 \log 2\sigma - \frac{1}{\sigma}|x_1 - \mu_{i1}| - \frac{1}{\sigma}|x_2 - \mu_{i2}| + \log P(C_i) \\
 g(\mathbf{x}) &= g_1(\mathbf{x}) - g_2(\mathbf{x}) \\
 &= -|x_1 - 1| - |x_2 - 1| + |x_1 - 3| + |x_2 - 5|
 \end{aligned}$$

Then let  $g(\mathbf{x}) = 0$ , we get:

$$\begin{aligned}
 &-|x_1 - 1| - |x_2 - 1| + |x_1 - 3| + |x_2 - 5| = 0 \\
 \Rightarrow x_2 &= \begin{cases} 4, & \text{for } x_1 < 1 \\ -x_1 + 5 = 0, & \text{for } 1 \leq x_1 \leq 3 \\ 2, & \text{for } x_1 > 3 \end{cases}
 \end{aligned}$$

