

1.

Show $\min(a, b) \leq \sqrt{ab}$.

$$\textcircled{1} \text{ if } a < b, \quad \min(a, b) = a = \sqrt{a} \cdot a < \sqrt{ab}$$

$$\textcircled{2} \text{ if } a > b, \quad \min(a, b) = b = \sqrt{b} \cdot b < \sqrt{ab}$$

$$\textcircled{3} \text{ if } a = b, \quad \min(a, b) = a = b = \sqrt{a} \cdot a = \sqrt{b} \cdot b = \sqrt{ab}$$

So $\min(a, b) \leq \sqrt{ab}$ in all cases.

$$P(\text{error}) = \int \mathbb{I}_1 P(X=x, Y=2) dx + \int \mathbb{I}_2 P(X=x, Y=1) dx$$

 \mathbb{I}_i , indicator function
 $= \begin{cases} 1 & \text{when predict } Y=i \\ 0 & \text{else} \end{cases}$
 $\pi_i: P(Y=i)$

$$= \int \mathbb{I}_1 P(X=x|Y=2) \cdot \pi_2 dx + \int \mathbb{I}_2 P(X=x|Y=1) \pi_1 dx$$

$$= \int [(\mathbb{I}_1=1) \cdot P(X=x|Y=2) \cdot \pi_2 + (\mathbb{I}_1=0) \cdot P(X=x|Y=2)] dx$$

$$+ \int [(\mathbb{I}_2=1) \cdot P(X=x|Y=1) \pi_1 + (\mathbb{I}_2=0) \cdot P(X=x|Y=1)] dx$$

$$= \int \mathbb{I}_1 \cdot \min_{c \in \{1,2\}} (P(X=x|Y=c) \cdot \pi_c) dx + \int \mathbb{I}_2 \cdot \min_{c \in \{1,2\}} (P(X=x|Y=c) \cdot \pi_c) dx$$

Since when $\mathbb{I}_1=1$, predict $Y=1$, $P(X=x|Y=2)$ must be smaller than $P(X=x|Y=1)$
 similarly when $\mathbb{I}_2=2$, predict $Y=2$, $P(X=x|Y=1)$ must be smaller than $P(X=x|Y=2)$.

$$= \int \min_{c \in \{1,2\}} (P(X=x|Y=c) \cdot \pi_c) dx$$

$$\leq \int \sqrt{P(X=x|Y=1) \pi_1 \cdot P(X=x|Y=2) \pi_2} dx$$

$$= \sqrt{\pi_1 \pi_2} \int \sqrt{P(X=x|Y=1) \cdot P(X=x|Y=2)} dx$$

where $\pi_1 = P(Y=1)$ $\pi_2 = P(Y=2)$ Since $\pi_1 + \pi_2 = 1$, $\sqrt{\pi_1 \pi_2} \leq \frac{\pi_1 + \pi_2}{2} = \frac{1}{2}$

$$\leq \frac{1}{2} \int \sqrt{P(X=x|Y=1) \cdot P(X=x|Y=2)} dx$$

X is continuous,

$$\text{So } P(\text{error}) \leq \sqrt{P(Y=1) P(Y=2)} \int \sqrt{f(x|Y=1) f(x|Y=2)} dx$$

$$< \frac{1}{2} \int \sqrt{f(x|Y=1) f(x|Y=2)} \cdot dx$$

2. (a). $f(x|Y=k) \sim N(u_k, \sigma^2) \quad k=1,2$

Assume $u_1 > u_2$.

Decision boundary: At boundary x_0 , $f(x_0|Y=1)\pi_1 = f(x_0|Y=2)\pi_2$

$$\frac{f(x_0|Y=1)}{f(x_0|Y=2)} = \frac{\pi_2}{\pi_1}$$

$$\frac{f(x_0|Y=1)}{f(x_0|Y=2)} = \exp\left(-\frac{1}{2\sigma^2}((x_0-u_1)^2 - (x_0-u_2)^2)\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2}(x_0^2 + u_1^2 - 2x_0u_1 - x_0^2 - u_2^2 + 2x_0u_2)\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2}(2x_0(u_2 - u_1) + u_1^2 - u_2^2)\right)$$

$$= \frac{\pi_2}{\pi_1}$$

$$(2x_0(u_2 - u_1) + u_1^2 - u_2^2) = -\left(\log \frac{\pi_2}{\pi_1}\right)2\sigma^2$$

$$x_0 = \frac{-\left(\log \frac{\pi_2}{\pi_1} \cdot 2\sigma^2\right) - u_1^2 + u_2^2}{2(u_2 - u_1)}$$

$$= \frac{u_2^2 - u_1^2 - 2\sigma^2 \log \frac{\pi_2}{\pi_1}}{2(u_2 - u_1)}$$

So the decision boundary is $x_0 = \frac{u_2^2 - u_1^2 - 2\sigma^2 \log \frac{\pi_2}{\pi_1}}{2(u_2 - u_1)}$

Predict $\hat{Y} = 1$ when $x \geq x_0$
 $\hat{Y} = 2$ when $x < x_0$

Bayes Loss: $L(\hat{Y}, Y) = E(I_{\{\hat{Y} \neq Y\}})$

$$= P(\hat{Y} \neq Y) \cdot 1 + P(\hat{Y} = Y) \cdot 0$$

$$= P(\hat{Y} \neq Y)$$

$I_{\{\hat{Y} \neq Y\}}$ is the indicator function

$$= \begin{cases} 1 & \hat{Y} \neq Y \\ 0 & \hat{Y} = Y \end{cases}$$

$$= P(x > x_0 | Y=2)\pi_2 + P(x < x_0 | Y=1)\pi_1$$

$$= P\left(\frac{x - u_2}{\sigma} > \frac{x_0 - u_2}{\sigma}\right) \cdot \pi_2 + P\left(\frac{x - u_1}{\sigma} < \frac{x_0 - u_1}{\sigma}\right) \cdot \pi_1$$

$$= [1 - \Phi\left(\frac{x_0 - u_2}{\sigma}\right)]\pi_2 + \Phi\left(\frac{x_0 - u_1}{\sigma}\right) \cdot \pi_1$$

is the Bayes Loss

(b).

$$P(\text{error}) = P(\hat{Y} \neq Y) = [1 - \Phi(\frac{x_0 - u_2}{\sigma})] \pi_2 + \Phi(\frac{x_0 - u_1}{\sigma}) \pi_1$$

As $\sigma \rightarrow 0$, $\Phi(\frac{x_0 - u_2}{\sigma}) \rightarrow 1$ because $u_1 > x_0 > u_2$, $x_0 - u_2 > 0$, $\frac{x_0 - u_2}{\sigma} \rightarrow +\infty$

$\Phi(\frac{x_0 - u_1}{\sigma}) \rightarrow 0$ because $u_1 > x_0 > u_2$, $x_0 - u_1 < 0$, $\frac{x_0 - u_1}{\sigma} \rightarrow -\infty$

$$\text{So } P(\hat{Y} \neq Y) \Rightarrow (1-1) \pi_2 + 0 \cdot \pi_1$$

$$= 0$$

So $P(\text{error}) \rightarrow 0$ as $\sigma \rightarrow 0$.

(c).

For σ fixed, as $\pi_1 \rightarrow 0$, $x_0 = \frac{u_2^2 - u_1^2 - 2\sigma^2 \log \frac{\pi_2}{\pi_1}}{2(u_2 - u_1)} \rightarrow +\infty$ since $\frac{\pi_2}{\pi_1} \rightarrow +\infty$ & $u_2 < u_1$

For small π_1 , we can always predict $\hat{Y} = 2$.

$$P(\text{error}) = P(\hat{Y} \neq Y) = [1 - \Phi(\frac{x_0 - u_2}{\sigma})] \pi_2 + \Phi(\frac{x_0 - u_1}{\sigma}) \pi_1$$

$$\rightarrow (1-1) \pi_2 + 1 \cdot \pi_1 \quad \text{as } x_0 \rightarrow +\infty$$

$$= \pi_1$$

here π_1 is small, so low error rate guaranteed.

(d).

Loss of individual example: $L(h, x) = \sum_{k=1}^K L_{k, h(x)} P(k|x)$

$$\begin{aligned} \text{Expectation of the individual loss: } E(L(h, x)) &= \int_{h(x)} f(x=x) \cdot L(h, x) dx \\ &= \int_{h(x)} f(x=x) \cdot \sum_{k=1}^K L_{k, h(x)} P(k|x) dx \\ &= \sum_{l=1}^K \sum_{k=1}^K \int_{h(x)=l} L_{k, l} P(k|x) f(x=x) dx \\ &= \sum_{k=1}^K \sum_{l=1}^K \int_{h(x)=l} P(x, k) L_{k, l} dx \\ &= L(h) \end{aligned}$$

So to minimize $L(h)$, we can minimize $L(h, x)$ for each given x .

$$\min L(h, x) = \sum_{k=1}^K L_{k, h(x)} P(k|x) \quad \text{by definition of } h(x).$$

So $h(x) = \arg \min_{j=1, \dots, K} \sum_{k=1}^K L_{k, j} P(k|x)$ gives the lowest loss.

(e). Expected loss $L(h) = \sum_{k=1}^2 \sum_{l=1}^2 \int_{h(x)=l} p(x,k) L_{k,l} dx$

$$= \sum_{k=1}^2 \left(\int_{h(x)=1} p(x,k) L_{k,1} dx + \int_{h(x)=2} p(x,k) L_{k,2} dx \right)$$

$$= \int_{h(x)=1} (p(x,1) L_{1,1} + p(x,2) L_{2,1}) dx + \int_{h(x)=2} (p(x,1) L_{1,2} + p(x,2) L_{2,2}) dx$$

$h_B(x) = \underset{j=1,2}{\operatorname{argmin}} \sum_{k=1}^2 L_{k,j} P(k|x)$

decision boundary: $(L_{1,1} P(1|x_0) + L_{2,1} P(2|x_0) = L_{1,2} P(1|x_0) + L_{2,2} P(2|x_0))$

$$(L_{1,1} - L_{1,2}) P(1|x_0) = (L_{2,2} - L_{2,1}) P(2|x_0)$$

$$P(1|x_0) = \frac{P(x_0|1) \cdot \pi_1}{f_X(x_0)}$$

$$P(2|x_0) = \frac{P(x_0|2) \cdot \pi_2}{f_X(x_0)}$$

$$\frac{P(1|x_0)}{P(2|x_0)} = \frac{L_{2,2} - L_{2,1}}{L_{1,1} - L_{1,2}}$$

$$\text{LHS} = \frac{P(x_0|1) \cdot \pi_1}{P(x_0|2) \cdot \pi_2} = \exp\left(-\frac{1}{2\sigma^2} (2x_0(u_2 - u_1) + u_1^2 - u_2^2)\right) \cdot \frac{\pi_1}{\pi_2}$$

$$= \frac{L_{2,2} - L_{2,1}}{L_{1,1} - L_{1,2}}$$

$$-\frac{1}{2\sigma^2} (2x_0(u_2 - u_1) + u_1^2 - u_2^2) = \log\left(\frac{L_{2,2} - L_{2,1}}{L_{1,1} - L_{1,2}} \cdot \frac{\pi_2}{\pi_1}\right)$$

$$x_0 = \frac{-2\sigma^2 \cdot \log\left(\frac{L_{2,2} - L_{2,1}}{L_{1,1} - L_{1,2}} \cdot \frac{\pi_2}{\pi_1}\right) + u_2^2 - u_1^2}{2(u_2 - u_1)}$$

By definition of $L_{k,j}$, when $k=j$, prediction is correct, so we can let $L_{k,j}=0$ when $k=j$.

$$x_0 = \frac{2\sigma^2 \cdot \log\left(\frac{L_{1,2} \cdot \pi_1}{L_{2,1} \cdot \pi_2}\right) + u_2^2 - u_1^2}{2(u_2 - u_1)} \text{ is the decision boundary.}$$

Let $L_{1,2} = \frac{1}{\pi_1}$, $L_{2,1} = \frac{1}{\pi_2}$, x_0 becomes

$$= \frac{2\sigma^2 \cdot \log(1) + u_2^2 - u_1^2}{2(u_2 - u_1)}$$

$$= \frac{u_2^2 - u_1^2}{2(u_2 - u_1)}$$

$$= \frac{u_2 + u_1}{2} \text{ is independent of } \pi_1 \text{ and } \pi_2.$$

So this will remedy the effect of small π_1 or π_2 .

3. (a).

Newton's Iteration $\theta_{new} = \theta_{old} - \frac{\nabla J(\theta)}{H(\theta)}$

To find $\nabla J(\theta)$ and $H(\theta)$:

$$\text{sigm}(y) = \frac{1}{1 + e^{-y}} = \frac{e^y}{e^y + 1}$$

$$P(y|x, \theta) = \prod_{i=1}^n (\text{sigm}(x_i \theta)^{y_i} (1 - \text{sigm}(x_i \theta))^{1-y_i}) \quad \text{where } x_i \theta = \theta_0 + \sum_{j=1}^d \theta_j x_{ij}$$

$$P(y|x, \theta) \in (0, 1)$$

n : number of data point
 θ : coefficients and intercept

The negative loglikelihood $J(\theta) = -\log P(y|x, \theta)$

$$\nabla J(\theta) = \frac{dJ(\theta)}{d\theta} = \sum_{i=1}^n x_i^T (\text{sigm}(x_i \theta) - y_i) = x^T (\text{sigm}(X\theta) - y) \quad \begin{matrix} x \in \mathbb{R}^{n \times d} \\ \theta \in \mathbb{R}^{d+1} \\ y \in \mathbb{R}^{n \times 1} \\ \nabla J(\theta) \in \mathbb{R}^{d+1} \end{matrix}$$

$$H(\theta) = D(\nabla J(\theta))$$

$$= D(x^T (\text{sigm}(X\theta) - y))$$

$$= x^T \text{diag}[\text{sigm}(x_i \theta) \cdot (1 - \text{sigm}(x_i \theta))] x \quad i=1, 2, \dots, n$$

$$\theta_{new} = \theta_{old} - \frac{\nabla J(\theta)}{H(\theta)}$$

$$= \theta_{old} - H^{-1}(\theta) \cdot \nabla J(\theta) \quad \text{Let } H(\theta) = x^T S x \quad \text{where } S = \text{diag}[\text{sigm}(x_i \theta) \cdot (1 - \text{sigm}(x_i \theta))]$$

$$= \theta_{old} - (x^T S x)^{-1} x^T (\pi - y) \quad \text{Let } \nabla J(\theta) = x^T (\pi - y) \quad \text{where } \pi = \text{sigm}(X\theta)$$

$$= (x^T S x)^{-1} [(x^T S x) \theta_{old} + x^T (y - \pi)]$$

$$= (x^T S x)^{-1} x^T [S X \theta_{old} + y - \pi]$$

θ_{new} is the solution of a weighted least square problem.

$$WLS(\theta, S) = \frac{1}{n} \sum_{i=1}^n s_i (y_i - x_i \theta)^2 \quad s_i = \text{sigm}(x_i \theta) (1 - \text{sigm}(x_i \theta))$$

$$\hat{\theta}_{WLS} = (x^T S x)^{-1} x^T [S X \theta_{old} + y - \pi]$$

$$= \theta_{new}$$

(b). If the maximum conditional Likelihood estimator exist, let it be $\hat{\theta}$, then $\nabla J(\hat{\theta}) = 0$.

$$x^T (\text{sigm}(X\hat{\theta}) - y) = 0$$

Let θ be such that $\begin{cases} y_i = 0, & x_i \theta < 0 \\ y_i = 1, & x_i \theta > 0 \end{cases}$

Since $x^T (\text{sigm}(X\hat{\theta}) - y) = 0$,

$$\theta^T x^T (\text{sigm}(X\hat{\theta}) - y) = 0$$

$x \in \mathbb{R}^{n \times d}$
 $\hat{\theta}, \theta \in \mathbb{R}^{d+1}$

$$\theta^T X^T (\text{sgm}(X\hat{\theta}) - y) = [\theta^T x_1^T, \theta^T x_2^T, \dots, \theta^T x_n^T] \begin{bmatrix} \text{sgm}(x_1 \hat{\theta}) - y_1 \\ \text{sgm}(x_2 \hat{\theta}) - y_2 \\ \vdots \\ \text{sgm}(x_n \hat{\theta}) - y_n \end{bmatrix}$$

$$= \sum_{y_i=0} \theta^T x_i^T \text{sgm}(x_i \hat{\theta}) + \sum_{y_i=1} \theta^T x_i^T (\text{sgm}(x_i \hat{\theta}) - 1)$$

$$< 0 \quad \text{Since } \theta^T x_i^T < 0 \text{ when } y_i=0, \text{sgm}(x_i \hat{\theta}) > 0, \sum_{y_i=0} \theta^T x_i^T \text{sgm}(x_i \hat{\theta}) < 0$$

$$\theta^T x_i^T > 0 \text{ when } y_i=1, \text{sgm}(x_i \hat{\theta}) - 1 < 0, \sum_{y_i=1} \theta^T x_i^T (\text{sgm}(x_i \hat{\theta}) - 1) < 0$$

Contradicts that $\theta^T X^T (\text{sgm}(X\hat{\theta}) - y) = 0$.

So the maximum conditional likelihood estimator doesn't exist.

References: Problem 3(a) referenced from www.cs.ox.ac.uk/people/nando.defreitas/machinelearning
www.stat.cmu.edu/~shalizi/mæeg/15/lectures/24
 Problem 3(b) referenced from Zihao Wang