

Question 1.

$$\log \frac{P(y=c_1 | x; w)}{P(y=c_2 | x; w)} = \log \frac{\frac{\exp(w c_1 \cdot x)}{\sum_{y=1}^C \exp(w y \cdot x)}}{\frac{\exp(w c_2 \cdot x)}{\sum_{y=1}^C \exp(w y \cdot x)}} = \log \exp(w c_1 - w c_2) \cdot x = (w c_1 - w c_2) \cdot x.$$

Where $(w c_1 - w c_2)$ is a constant given any c_1 and c_2 .

Thus the log odd is a linear function of x .

When $C=2$,

$$P(y=1 | x; w) = \frac{\exp(w_1 \cdot x)}{\exp(w_1 \cdot x) + \exp(w_2 \cdot x)}$$

To find an alternative setting:

$$\begin{cases} P(y=1 | x; w) + P(y=2 | x; w) = 1 \\ P(y=1 | x; w) = P(y=2 | x; w) \cdot \exp(w_1 - w_2) \cdot x \end{cases}$$

$$\begin{aligned} \Rightarrow P(y=1 | x; w) &= \frac{\exp(w_1 - w_2) \cdot x}{1 + \exp(w_1 - w_2) \cdot x} = \frac{1}{1 + \exp(-(w_1 - w_2) \cdot x)} \\ &= \sigma[(w_1 - w_2) \cdot x] \\ &= \sigma(v \cdot x) \text{ where } v = w_1 - w_2 \end{aligned}$$

\Rightarrow There exists a single parameter $v = (w_1 - w_2)$ such that the above equation holds.



Question 2;

For any $C_k \in \{1, \dots, C\}$

$$\begin{aligned}
 P(y=k | X; w) &= \frac{\exp(w_k \cdot X)}{\sum_{c=1}^C \exp(w_c \cdot X)} = \frac{\exp(w_k \cdot X)}{\exp(w_1 \cdot X) + \exp(w_2 \cdot X) + \dots + \exp(w_k \cdot X) + \dots + \exp(w_C \cdot X)} \\
 &= \frac{1}{\exp((w_1 - w_k) \cdot X) + \dots + \underbrace{\exp((w_k - w_k) \cdot X)}_{=1} + \dots + \exp((w_C - w_k) \cdot X)} \\
 &= \frac{1}{1 + \underbrace{\sum_{c=1}^{k-1} (w_c - w_k) \cdot X + \sum_{c=k+1}^C (w_c - w_k) \cdot X}_{C-1 \text{ number of parameters}}} \quad (*)
 \end{aligned}$$

setting (*) yields the same $p(y|x)$ for every x .

This also shows we only need $C-1$ parameters because there is only $C-1$ parameter in equation (*), where we just subtract category k from each category.



Question 3.

For logistic regression.

$$L(w) = \arg \min \log(P(y|x)) = - \sum_{i=1}^n y_i \log \sigma(w \cdot x_i) + (1-y_i) \log(1-\sigma(w \cdot x_i)) \quad \text{where } \sigma(w \cdot x_i) = \frac{1}{1+e^{-w \cdot x_i}}$$

$$\frac{\partial \log \sigma(w \cdot x_i)}{\partial w_j} = \frac{\partial -\log(1+e^{-w \cdot x_i})}{\partial w_j} = + \frac{x_{ij} e^{-w \cdot x_i}}{1+e^{-w \cdot x_i}} = x_{ij}(1-\sigma(w \cdot x_i))$$

$$\frac{\partial \log(1-\sigma(w \cdot x_i))}{\partial w_j} = \frac{-x_{ij} e^{-w \cdot x_i}}{1-\frac{1}{1+e^{-w \cdot x_i}}} = -x_{ij}\sigma(w \cdot x_i)$$

$$\begin{aligned} \text{F.O.C: } \frac{\partial L(w)}{\partial w_j} &= - \sum_{i=1}^n y_i x_{ij}(1-\sigma(w \cdot x_i)) - (1-y_i) x_{ij} \sigma(w \cdot x_i) \\ &= - \sum_{i=1}^n y_i x_{ij} - y_i x_{ij} \sigma(w \cdot x_i) - x_{ij} \sigma(w \cdot x_i) + y_i x_{ij} \sigma(w \cdot x_i) \\ &= \sum_{i=1}^n x_{ij} \sigma(w \cdot x_i) - y_i x_{ij} \\ &= \sum_{i=1}^n x_{ij} [\sigma(w \cdot x_i) - y_i] \end{aligned}$$

$$\text{F.O.C gives the gradient} = X^T [\sigma(w \cdot x_i) - y_i]$$

$$\begin{aligned} \text{S.O.C: } \frac{\partial^2 L(w)}{\partial w_j \partial w_k} &= \sum_{i=1}^n x_{ij} x_{ik} \sigma(w \cdot x_i) [1-\sigma(w \cdot x_i)] \\ &= z_j^T \beta z_k \end{aligned}$$

$$\begin{aligned} \text{Where } z_j^T &= [x_{1j}, \dots, x_{nj}] \quad \beta = \text{diag}(\sigma(w \cdot x_1)[1-\sigma(w \cdot x_1)], \dots, \sigma(w \cdot x_n)[1-\sigma(w \cdot x_n)]) \\ z_k &= \begin{bmatrix} x_{1k} \\ \vdots \\ x_{nk} \end{bmatrix} \quad \sigma(w \cdot x_i) [1-\sigma(w \cdot x_i)] = \frac{1}{1+e^{-w \cdot x_i}} \frac{e^{-w \cdot x_i}}{1+e^{-w \cdot x_i}} \in [0,1] \quad (*) \end{aligned}$$

Because z_j is the j^{th} column of X , z_k is the k^{th} column of X

$z_j^T \beta z_k$ is the $j-k^{\text{th}}$ entry of H , thus $H = X^T \beta X$

Because all entries in $\beta \geq 0$ (proved by (*)), we can re-write H as

$$H = X^T \beta^{\frac{1}{2}} \beta^{\frac{1}{2}} X = (\beta^{\frac{1}{2}} X)^T (\beta^{\frac{1}{2}} X) \Rightarrow H \text{ is positive definite} \Rightarrow L \text{ is convex}$$

Taylor's expansion

$$L(w) \approx L(a) + g^T(w-a) + \frac{1}{2}(w-a)^T H(w-a)$$

$$= \frac{1}{2} w^T H w + b^T w + c$$

$$\text{where } b = g - H a$$

$$0 = H w + b \Rightarrow w^* = H^{-1} b = -H^{-1} g + a \Rightarrow \text{Therefore the minimum is unique.}$$



Question 4

$$W = \begin{bmatrix} w_{10} & w_{11} & \dots & w_{1d-1} \\ w_{20} & w_{21} & \dots & w_{2d-1} \\ \vdots & \vdots & \ddots & \vdots \\ w_{c0} & w_{c1} & \dots & w_{cd-1} \end{bmatrix}_{c \times d} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d-1} \\ 1 & x_{21} & x_{22} & \dots & x_{2d-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nd-1} \end{bmatrix}_{n \times d} \quad x_i$$

an sample x_i belongs to class c , has 1 in the c^{th} position, and 0 in the rest of cells.

for each row (sample), $\hat{p}(y_i=c|x_i, W) = \frac{\exp(w_i \cdot x_i)}{\sum_{y=1}^c \exp(w_y \cdot x_i)}$

$$W^* = \arg \min L_i(x_i, y_i, W) = -\frac{1}{N} \sum_{i=1}^n t_i \log \hat{p}(y_i|x_i, W) + \lambda \|W\|^2$$

$$= -\frac{1}{N} \sum_{i=1}^n t_i \log \hat{p}_i + \lambda \|W\|^2$$

$$\frac{\partial L_i}{\partial w_{ij}} = \sum_{i=1}^n \frac{\partial L}{\partial w_i x_i} \frac{\partial w_i x_i}{\partial w_{ij}} + \lambda \sum_i \sum_j \frac{\partial w_{ij}^2}{\partial w_{ij}}$$

$$= \frac{1}{N} \sum_{i=1}^N [\hat{p}_i - t_i] x_{ij} + 2\lambda \sum_{i=1}^n w_{ij}$$

$$= \frac{1}{N} \sum_{i=1}^N (\hat{p}_i - t_i) x_i + 2\lambda W_i$$

Stochastic gradient. $W_{t+1} = \eta \frac{1}{N} (X^T (\hat{p} - t)) + 2\lambda W_t$

