## Ex 3.1 (a)

(i) Correlation $[\rho_{xy}] = \dfrac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 (y_i - \bar{y})^2}}$

$= \dfrac{\dfrac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{n^2}}}{\sqrt{\dfrac{\sum (x_i - \bar{x})^2}{n} \dfrac{\sum (y_i - \bar{y})^2}{n}}}$

$\sigma_x^2 = \dfrac{\sum (x_i - \bar{x})^2}{n} \Rightarrow \boxed{\sigma_x^2 = \dfrac{\sum x_i^2}{n}}$

$\rho_{xy} = \dfrac{\dfrac{\sum x_i y_i}{n}}{\sigma_x \, \sigma_y}$

$\Rightarrow \sum x_i y_i = n \, \sigma_x \sigma_y \rho_{xy}$

$$\boxed{a = (X^T X)^{-1} X^T Y}$$

$$X^T X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix}$$

$$= \sum x_i^2 = n\sigma_x^2$$

$$(X^T X)^{-1} = \frac{1}{n\sigma_x^2}$$

$$(X^T Y) = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix}$$

$$= \sum x_i y_i$$

$$= n\sigma_x \sigma_y \rho_{xy}$$

$$\Rightarrow a = \frac{n\sigma_x \sigma_y \rho_{xy}}{n\sigma_x^2}$$

$$\boxed{a = \frac{\sigma_y \rho_{xy}}{\sigma_x}}$$

(ii) Given $X' = sX$ and $Y' = tY$

$$\Rightarrow \text{ if } \bar{X} = \bar{Y} = 0$$

$$\Rightarrow \overline{X'} = \overline{Y'} = 0$$

$$\sigma_{x'} = \sqrt{\frac{\sum x_i'^2}{n}} = s\sqrt{\frac{\sum x_i^2}{n}} = s\sigma_x$$

Similarly $\rightarrow \sigma_{y'} = t\sigma_y$

$$f_{x'y'} = \frac{\sum x_i' y_i'}{n \sigma_{x'} \sigma_{y'}}$$

$$= \frac{st \sum x_i y_i}{n \, s\sigma_x \, t\sigma_y}$$

$$= \frac{\sum x_i y_i}{n \sigma_x \sigma_y}$$

$$= f_{xy}$$

$$\Rightarrow \boxed{f_{x'y'} = f_{xy}}$$

Here $a' = \dfrac{\sigma_{y'} \, f_{x'y'}}{\sigma_{x'}}$

$$= \dfrac{t \, \sigma_y \, f_{xy}}{s \, \sigma_x}$$

$$= \dfrac{t}{s} \, \dfrac{\sigma_y \, f_{xy}}{\sigma_x}$$

$\Rightarrow$ $\boxed{a' = \dfrac{t}{s} \, a}$

b) MSE $\quad f(x,y:w) = \frac{1}{n}\sum_{i=1}^{n}(y_i - \langle w, x_i\rangle)^2$

Adding the noise term:

$$f(x,y:w) = \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}(y_i - \langle w,x_i\rangle)^2 - \right.$$
$$\left. 2\epsilon \cdot w_i(y_i - \langle w, x_i\rangle) + \sum_{i=1}^{n} w_i\, \epsilon_i \epsilon_j \sum_{j=1}^{d} w_i\right]$$

Applying linearity of expection:

$$= \frac{1}{n}\sum_{i=1}^{n}\left((y_i - \langle w,x_i\rangle)^2 - 2\mathbb{E}[\epsilon]w_i(y_i - \langle w,x_i\rangle) + \mathbb{E}[\epsilon_i\epsilon_j]\sum_{i=1}^{d}w_i^2\right)$$

Given $\quad \mathbb{E}[\epsilon] = 0$ and $\mathbb{E}[\epsilon_i\epsilon_j] = \sigma^2$

$$f(w) = \frac{1}{n}\sum_{i=1}^{n}\left((y_i - \langle w,x_i\rangle)^2 + \sigma^2\sum_{i=1}^{d}w_i^2\right)$$

Therefore we can write:

$$\mathbb{E}\left[(y_i - \langle w, x_i + \epsilon_i \rangle)^2\right] =$$

$$\mathbb{E}\left[(y_i - \langle w, x_i \rangle)^2\right] + \sigma^2 \sum_{i=1}^{d} w_i^2$$

## Ex 3.3

a) The softmax function is susceptible to two vulnerabilities.

Case I : When very small numbers are passed through softmax, it is rounded up to zero; this is known as underflo-wing.

Case II : When the numbers are very large, Softmax approximations are interpreted as infinity. This is called overflowing.

b) The overflowing problem of softmax can be solved by subtracting the largest numbers and shifting all the inputs. So the equation becomes:

$$exp(x_i) - max(x)$$

$$\text{Softmax}(x)_i = \frac{\exp(x_i) \cdots \max(x)}{\sum_{j=1}^{n} \exp(x_j) - \max(x)}$$

It can be shown by exponential division rule that both equations are same while solving the numerical issues of Softmax.

c) For a function to be a valid probability distribution it needs to satisfy 3 attributes:

    ① The random variable should be associated with numerical.

    ② Sum of probabilities should be 1

    ③ Each probability should be in range of 0~1

As we know exp function returns values from $0 \sim +\infty$, the result of $\exp(x_i)$ will always return a positive value. The same applies for the denominator. The denominator normalizes the value of

$\exp(x_i)$ dividing by the sum of all $\exp(x_j)$ for all possible $n$. Since the probabilities are a ration of these two values, they'll always sum up to 1 as the largest value will always return the largest probability. No matter how large the upper value is the denominator will always normalize it and turn into a value between $0 \sim 1$. Thus it will always return a valid probabity distribution for all values of $n \in \mathbb{R}$.

d) Jacobian matrix of Softmax:

$$S_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$Js_i = \begin{bmatrix} \dfrac{\partial s_1}{\partial x_1} & \dfrac{\partial s_1}{\partial x_2} & \cdots & \cdots & \dfrac{\partial s_1}{\partial x_n} \\ \dfrac{\partial s_2}{\partial x_1} & \dfrac{\partial s_2}{\partial x_2} & \ddots & & \vdots \\ & & & & \end{bmatrix}$$

$$\begin{vmatrix} \vdots & \vdots & \cdots & \vdots \\ \dfrac{\partial S_n}{\partial x_1} & \dfrac{\partial S_n}{\partial x_2} & \cdots & \dfrac{\partial S_n}{\partial x_n} \end{vmatrix}$$

For easier computation we are considering the log of Softmax,

$$\log S_i = \log \left( \frac{e^{x_i}}{\sum\limits_{j} e^{x_j}} \right)$$

$$= x_i - \log \left( \sum\limits_{j}^{n} e^{x_j} \right)$$

For any $k$ th value;

$$\frac{\partial}{\partial x_k} \left( \log S_i \right) = \frac{\partial x_i}{\partial x_k} - \frac{\partial}{\partial x_k} \log \left( \sum\limits_{j} e^{x_j} \right)$$

here,

$$\frac{\partial x_i}{\partial x_k} = \begin{cases} 1 & ; \quad i = k \\ 0 & ; \quad i \neq k \end{cases}$$

$$\frac{\partial}{\partial x_k} \left( \log S_i \right) = 1\{i=k\} - \frac{1}{\sum\limits_{j} e^{x_j}} \cdot \left( \frac{\partial}{\partial x_j} e^{x_j} \right)$$

$$= 1\{i=j\} - \frac{e^{x_j}}{\sum_k e^{x_k}} \qquad \left[\frac{d}{dx} \log(x) = \frac{1}{x}\right]$$

Therefore, we can write

$$\frac{\partial}{\partial x_k} \log(s_i) = 1\{i=k\} - s_k$$

$$\frac{1}{s_i} \cdot \frac{\partial s_i}{\partial x_k} = 1\{i=k\} - s_k$$

$$\frac{\partial s_i}{\partial x_k} = s_i \left(1\{i=k\} - s_k\right)$$

Plugging this formula into Jacobian format we get:

$$J_s = \begin{bmatrix} s_1 \cdot (1-s_1) & -s_1 \cdot s_2 & -s_1 \cdot s_2 \cdots & -s_1 \cdot s_n \\ -s_2 \cdot s_1 & s_2 \cdot (1-s_2) & \cdots & -s_2 \cdot s_4 \\ \vdots & \vdots & \ddots & \vdots \\ -s_n \cdot s_1 & -s_n \cdot s_2 & \cdots & s_n \cdot (1-s_n) \end{bmatrix}$$

$$f(x, y; \omega) = \frac{1}{n} \| Y - X\omega \|^2 + \lambda \|\omega\|^2$$

$$f(x, y; \omega) = \frac{1}{n} (y - X\omega)^T (y - X\omega) + \lambda \omega^T \omega$$

for best $\omega^*$, we use:

$$\frac{df(x, y; \omega)}{d\omega} = 0$$

$$\Rightarrow -\frac{2X^T}{n}(y - X\omega) + 2\lambda\omega = 0$$

$$\Rightarrow X^T(y - X\omega) = n\lambda\omega \qquad \text{[We write } n\lambda \text{ as } \lambda\text{]}$$

$$\Rightarrow X^T y - X^T X\omega = \lambda\omega$$

$$\Rightarrow (X^T X + \lambda I)\omega = X^T y$$

$$\Rightarrow \boxed{\omega^* = (X^T X + \lambda I)^{-1} X^T y}$$