Ans 1.1)

$$V = g(\theta) = 1 - \theta$$

$$P(n|V) = \theta^n (1-\theta)^{1-n}$$

$$= (1-v)^n v^{1-n} \qquad \left[ \begin{array}{l} V = 1 - \theta \\ \theta = 1 - V \end{array} \right]$$

$$\therefore P(n|v) = (1-v)^n v^{1-n} \qquad n \in \{0, 1\}$$

Joint probability

$$P(n_1, n_2, \cdots, n_n | v)$$

$$= \prod_{i=1}^{n} P(n_i | v)$$

$$= \prod_{i=1}^{n} (1-v)^{n_i} v^{1-n_i}$$

$$= (1-v)^{\sum_{i=1}^{n} n_i} \; v^{n - \sum_{i=1}^{n} n_i}$$

Log likelihood

$$\mathcal{L}(v) = \log P(n_1, n_2, \cdots, n_n | v)$$

$$= \log (1-v)^{\sum_{i=1}^{n} n_i} \; v^{n - \sum_{i=1}^{n} n_i}$$

$$\mathcal{L}(v) = \left( \sum_{i=1}^{n} n_i \right) (\log (1-v)) + v \left( n - \sum_{i=1}^{n} n_i \right)$$

For maximising $\mathcal{L}(v)$,

$$\frac{\partial \mathcal{L}(v)}{\partial v} = 0$$

$$\therefore \left( \sum_{i=1}^{n} x_i \right) \times \frac{1}{1-v} \times (-1) + \frac{1}{v} \times \left( n - \sum_{i=1}^{n} x_i \right) = 0$$

$$\therefore \quad \frac{\sum_{i=1}^{n} x_i}{1-v} = \frac{n - \sum_{i=1}^{n} x_i}{v}$$

$$\therefore \quad v \sum_{i=1}^{n} x_i = n - nv + \left( \sum_{i=1}^{n} x_i \right) v - \sum_{i=1}^{n} x_i$$

$$\therefore \quad nv = n - \sum_{i=1}^{n} x_i$$

$$\therefore \quad v = 1 - \frac{\sum_{i=1}^{n} x_i}{n}$$

$$\boxed{v_{MLE} = 1 - \frac{\sum_{i=1}^{n} x_i}{n}}$$

This is so as

$$v_{MLE} = g(\theta_{MLE}) = 1 - \theta_{MLE}$$

$$= 1 - \frac{\sum_{i=1}^{n} x_i}{n}$$

---

**An: 1.2)** Joint distribution

$$P(x_1, x_2, \cdots, x_{2n} \mid \theta)$$

$$= \prod_{i=1}^{2n} P(x_i \mid \theta)$$

$$= \prod_{i=1}^{2r} \left( \frac{1}{2} e^{-|x_i - \theta|} \right)$$

$$= \frac{1}{2^{2n}} e^{-\sum_{i=1}^{2n} |n_i - \theta|}$$

Log-likelihood

$$\mathcal{L}(\theta) = \log P(n_1, n_2, \cdots n_{2n} \mid \theta)$$

$$= \log \frac{1}{2^{2n}} e^{-\sum_{i=1}^{2n} |n_i - \theta|}$$

$$\mathcal{L}(\theta) = -\log 2^{2n} - \sum_{i=1}^{2n} |n_i - \theta|$$

For maximising $L(\theta)$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = 0$$

$$-\sum_{i=1}^{2n} \frac{\partial |n_i - \theta|}{\partial \theta} = 0$$

$$\therefore -\sum_{i=1}^{2n} \text{sign}(n_i - \theta) \times (-1) = 0 \qquad \left. \begin{array}{l} \frac{\partial |n|}{\partial n} \\ = \text{sign}(n) \\ = 1 \quad n > 0 \\ = 0 \quad n = 0 \\ = -1 \quad n < 0 \end{array} \right\}$$

$$\therefore \sum_{i=1}^{2n} \text{sign}(n_i - \theta) = 0$$

This is Each

As the sign function takes only $+1, -1$ and $0$ values,

no of terms with value $-1$ = no of terms with value $1$

That is, $\theta$ is greater than exactly half of $n_1, n_2, \cdots, n_{2n}$ and less than half

of $n_1, n_2, \ldots, n_{2n} \longrightarrow$ Sorted

Let $\boxed{n_1', n_2', \ldots, n_{2n}'}$ be the

" sequence obtained on sorting

$n_1, n_2, \ldots . n_{2n}$

So, $\displaystyle\sum_{i=1}^{2n} \text{sign}(n_i' - \theta) = \theta$

$\theta > $

So, $\theta$ is greater than $n_1', n_2', \ldots n_n'$

$\theta < $

and $\theta$ is less than $n_{n+1}', n_{n+2}', \ldots, n_{2n}'$

So, $\theta_{MLE}$ is any value in

the Interval $\left[ n_n', n_{n+1}' \right]$

One particular value for $\theta_{MLE}$ is the

$\boxed{\text{median}}$ of $n_1, n_2, \ldots, n_{2n}$

One value
of $\theta_{MLE} = \boxed{\dfrac{n_n' + n_{n+1}'}{2}}$ where $n_1', n_2', \ldots, n_{2n}$

are in $\boxed{\text{sorted}}$ order

Also, $\theta_{MLE}$ can be anything in

the Interval $\boxed{\left[ n_n', n_{n+1}' \right]}$

$\theta$ MLE is any $\boxed{\text{median}}$ of $u_1, u_2, \cdots, u_{2n}$.

---

Ans. 2)

2.1)

Likelihood of a single sample $i$

$$= P(y_i \mid u_i, \theta, d_1, d_2)$$

$$= 2(d_1, d_2) \frac{e^{d_1(y_i - \theta^T u_i)}}{\left(d_1 e^{2(y_i - \theta^T u_i)} + d_2\right)^{\frac{d_1 + d_2}{2}}}$$

Log likelihood of single sample $i$

$$= L_i(\theta) = \log\left(P(y_i \mid u_i, \theta, d_1, d_2)\right)$$

$$= \log\left(\frac{2(d_1, d_2) e^{d_1(y_i - \theta^T u_i)}}{\left(d_1 e^{2(y_i - \theta^T u_i)} + d_2\right)^{\frac{d_1 + d_2}{2}}}\right)$$

$$L_i(\theta) = \log 2(d_1, d_2) + d_1\left(y_i - \theta^T u_i\right)$$

$$- \# \frac{d_1 + d_2}{2} \times \log\left(d_1 e^{2(y_i - \theta^T u_i)} + d_2\right)$$

$$\therefore \frac{\partial L_i(\theta)}{\partial \theta} = 0 + d_1 \times (-u_i)$$

$$- \# \frac{d_1 + d_2}{2} \times \frac{1 \times d_1 \times e^{2(y_i - \theta^T u_i)} \times (-2u_i)}{d_1 e^{2(y_i - \theta^T u_i)} + d_2}$$

$$= -d_1 u_i + \frac{(d_1 + d_2) d_1 u_i \, e^{2(y_i - \theta^T u_i)}}{d_1 e^{2(y_i - \theta^T u_i)} + d_2}$$

$$\therefore \frac{\partial \ell_i(\theta)}{\partial \theta} = -d_1 u_i + \frac{(d_1+d_2)d_1 u_i}{d_1 + d_2 \, e^{-2(y_i - \theta^T u_i)}}$$

$$\therefore \boxed{\frac{\partial \ell_i(\theta)}{\partial \theta} = -d_1 u_i \left(1 - \frac{d_1 + d_2}{d_1 + d_2 \cdot e^{-2(y_i - \theta^T u_i)}}\right)}$$

<u>Ans</u>

Ans: 2·2) From the plots it is seen that the learning rate of $7e-2$ $7e-2$ in the first case is high which causes drastic updates at each step and causes divergent behaviour and the loss never converges. In the second case, the learning rate of $1e-3$ is the best as the loss decreases exponentially and converges swiftly. In the third case, the learning rate of $1e-2$ is too low as the loss decreases linearly and the convergence is very slow and the final value of loss is much longer than the second case for the same number of updates. Thus, a learning rate of $1e-3$ is optimal, $7e-2$ is high and $1e-6$ is low.

Ans: 2·3) 3)

Ans: 3·1)

Cumulative distribution function $F_{logistic}$

$$= \frac{1}{1 + e^{-(x-\mu)/s}} \, , \qquad$$

where $\mu$ is mean

and $s$ is standard deviation

For our case,

$$\mu = 0, \quad s = \sigma_\epsilon$$

So, $F_{logistic} = \dfrac{1}{1 + e^{-\frac{u}{\sigma_\epsilon}}} = F_{\epsilon_i}$

$P(y_i = 1 \mid \theta, u_i)$

$$= P(\theta^T u_i + \epsilon_i \geq 0)$$

$$= 1 - P(\theta^T u_i + \epsilon_i \leq 0)$$

$$= 1 - P(\epsilon_i \leq -\theta^T u_i)$$

$$= 1 - F_{logistic}(-\theta^T u_i) \quad \begin{bmatrix} \text{Definition} \\ \text{of cdf} \end{bmatrix}$$

$$= 1 - \dfrac{1}{1 + e^{-\frac{(-\theta^T u_i)}{\sigma_\epsilon}}}$$

$$= 1 - \dfrac{1}{1 + e^{\frac{\theta^T u_i}{\sigma_\epsilon}}}$$

$$= \dfrac{e^{\frac{\theta^T u_i}{\sigma_\epsilon}}}{1 + e^{\frac{\theta^T u_i}{\sigma_\epsilon}}}$$

$$= \dfrac{1}{1 + e^{-\frac{\theta^T u_i}{\sigma_\epsilon}}}$$

$$\therefore \boxed{P(y_i = 1, \theta, x_i) = \text{logistic}\left(\frac{\theta^T x_i}{\sigma_\epsilon}\right)}$$

Ans $\div$ 3.2)          Given, $\sigma_\epsilon = 1$

Now,
$$P(y_i = 0 \mid \theta, x_i)$$

$$= 1 - P(y_i = 1 \mid \theta, x_i)$$

$$\boxed{\sigma_\epsilon = 1} \quad = 1 - \text{logistic}\left(\frac{\theta^T x_i}{\sigma_\epsilon 1}\right) = 1 - \text{logistic}(\theta^T x_i)$$

$y_i$ takes only binary values, 0 or 1

$$y_i \in [0, 1]$$

Case 1    $y_i = 1$

$$RHS = \left(\text{logistic}(\theta^T x_i)\right)^1 \times \left(1 - \text{logistic}(\theta^T x_i)\right)^{1-1}$$

$$= \text{logistic}(\theta^T x_i) \times \left(1 - \text{logistic}(\theta^T x_i)\right)^0$$

$$= \text{logistic}(\theta^T x_i)$$

$$= P(y_i = 1 \mid \theta, x_i)$$

$$= LHS$$

Case 2    $y_i = 0$

$$RHS = \left(\text{logistic}(\theta^T x_i)\right)^0 \times \left(1 - \text{logistic}(\theta^T x_i)\right)^1$$

$$= 1 - \text{logistic}(\theta^T x_i)$$

$$= P(y_i = 0 \mid \theta, x_i)$$

$$= LHS$$

For all cases, RHS = LHS

Hence, $P(y_i \mid \theta, x_i) = \left(\text{logistic}(\theta^T x_i)\right)^{y_i} \left(1 - \text{logistic}(\theta^T x_i)\right)^{(1-y_i)}$

---

Ans 3.3)

$$\log P(y_i \mid \theta, x_i)$$

$$= \cancel{y_i} \log \left( \left(\text{logistic}(\theta^T x_i)\right)^{y_i} \left(1 - \text{logistic}(\theta^T x_i)\right)^{1-y_i} \right)$$

$$= y_i \log \left(\text{logistic}(\theta^T x_i)\right) + (1-y_i) \log\left(1 - \text{logistic}(\theta^T x_i)\right)$$

$$= y_i \log \left(\frac{1}{1 + \exp(-\theta^T x_i)}\right) + (1-y_i) \log\left(1 - \frac{1}{1 + \exp(-\theta^T x_i)}\right)$$

$$= y_i \log\left(\frac{\exp(\theta^T x_i)}{1 + \exp(\theta^T x_i)}\right) + (1-y_i) \log\left(1 - \frac{\exp(\theta^T x_i)}{1 + \exp(\theta^T x_i)}\right)$$

$$= y_i \log\left(\frac{\exp(\theta^T x_i)}{1 + \exp(\theta^T x_i)}\right) + (1-y_i) \log\left(\frac{1}{1 + \exp(\theta^T x_i)}\right)$$

$$= y_i \theta^T x_i - y_i \log\left(1 + \exp(\theta^T x_i)\right) + (1-y_i) \times -\log\left(1 + \exp(\theta^T x_i)\right)$$

$$= y_i \theta^T x_i - y_i \log\left(1 + \exp(\theta^T x_i)\right) - \log\left(1 + \exp(\theta^T x_i)\right)$$

$$+ y_i \log\left(1 + \exp(\theta^T x_i)\right)$$

$$\boxed{\log P(y_i | \theta, x_i) = y_i \theta^T x_i - \log (1 + \exp(\theta^T x_i))}$$

Ans. 3.4)

$$l_{MLE}(\theta) \underset{n}{=} \log P(y_1, \cdots y_n | \theta, x_1, x_2 \cdots x_n)$$

$$= \log P(y|x, \theta) = \log \prod_{i=1}^{IT} P(y_i | \theta, x_i)$$

$$= \sum_{i=1}^{n} \log P(y_i | \theta, x_i)$$

$$l_{MLE}(\theta) = \sum_{i=1}^{a} \left( y_i \theta^T x_i - \log (1 + \exp(\theta^T x_i)) \right] \qquad \text{①}$$

Now,

$$\theta^T x_i = \sum_{j=1}^{d} \theta_j x_{ij}$$

$$= \sum_{j=1}^{d} x_{ij} \theta_j$$

$$\boxed{\theta^T x_i = (x \theta)_i} \qquad \left[ \begin{array}{l} \text{Multiplication} \\ \text{definition} \end{array} \right.$$

From ①

$$\therefore l_{MLE} = \sum_{i=1}^{n} y_i (x\theta)_i - \sum_{i=1}^{n} \log \left( 1 + \exp ((x\theta)_i) \right)$$

$$= \sum_{i=1}^{n} y_i (x\theta)_i - \sum_{i=1}^{n} \log \left( 1_{nx1_i} + \exp(x\theta)_i \right)$$

$$= y^T x \theta - \sum_{i=1}^{n} \log \left( 1_{ax1} + \exp(x\theta) \right)_i$$

$$\left[ x^T y = \sum_{i=1}^{n} x_i y_i \right]$$

$$= y^T x \theta \quad - \sum_{i=1}^{n} \left(\underline{1}_{n \times 1}\right)_i \log\left(\underline{1}_{n \times 1} + \exp(x\theta)\right)_i$$

$$\left[\begin{array}{l} \left(\underline{1}_{n \times 1}\right)_i = 1 \\ \text{multiplication by 1} \end{array}\right]$$

$$= y^T x \theta \quad - \underline{1}_{n \times 1}^T \cdot \log\left(\underline{1}_{n \times 1} + \exp(x\theta)\right)$$

Hence

$$\boxed{\mathcal{L}_{MLE}(\theta) = y^T x \theta - \underline{1}_{n \times 1}^T \cdot \log\left(\underline{1}_{n \times 1} + \exp(x\theta)\right)}$$

Ans. 3.5) We have,

$$\theta^T u_i = \sum_{k=1}^{d} \theta_k x_{ik}$$

$$\therefore \quad \frac{\partial \theta^T u_i}{\partial \theta_j} = \sum_{k=1}^{d} \frac{\partial(\theta_k x_{ik})}{\partial \theta_j}$$

$$\frac{\partial(\theta^T u_i)}{\partial \theta_j} = x_{ij} \qquad \text{———①}$$

Also,
$$\theta^T u_i = \sum_{k=1}^{d} \theta_k x_{ik}$$

$$= \sum_{k=1}^{d} x_{ik} \times \theta_k \qquad \qquad \text{Definition of}$$

$$\theta^T u_i = (x\theta)_i \qquad \text{———②} \qquad \left[\begin{array}{l} \text{Matrix} \\ \text{multiplication} \end{array}\right]$$

$$\mathcal{L}_{MLE}(\theta) = \sum_{i=1}^{n} \log P\left(y_i \mid \theta u_i\right)$$

$$\mathcal{L}_{MLE}(\theta) = \sum_{i=1}^{n} \left(y_i \theta^T u_i - \log\left(1 + \exp(\theta^T x_i)\right)\right)$$

$$\therefore \frac{\partial l_{MLE}(\theta)}{\partial \theta_j} = \sum_{i=1}^{n} \left( \frac{\partial(y_i \theta^T x_i)}{\partial \theta_j} - \frac{\partial \log(1 + \exp(\theta^T x_i))}{\partial \theta_j} \right)$$

$$= \sum_{i=1}^{n} \left( y_i \frac{\partial \theta^T x_i}{\partial \theta_j} - \frac{1 \times \exp(\theta^T x_i)}{1 + \exp(\theta^T x_i)} \times \frac{\partial \theta^T x_i}{\partial \theta_j} \right)$$

$$= \sum_{i=1}^{n} \left( y_i x_{ij} - \frac{1}{1 + \exp(-\theta^T x_i)} \times x_{ij} \right)$$
$$[\text{Using } ①]$$

$$= \sum_{i=1}^{n} \left( y_i x_{ij} - logistic(\theta^T x_i) \times x_{ij} \right)$$

$$= \sum_{i=1}^{n} \left( y_i x_{ij} - logistic\left( (x\theta)_i \right) \times x_{ij} \right)$$
$$\left[ \text{Using } ② \right]$$

$$= \sum_{i=1}^{n} x_{ij} \left( y_i - (logistic(x\theta))_i \right)$$

$$= \sum_{i=1}^{n} (x^T)_{ji} \times \left( y - logistic(x\theta) \right)_i$$

$$\frac{\partial l_{MLE}(\theta)}{\partial \theta_j} = \left( x^T \left( y - logistic(x\theta) \right) \right)_j$$

$$\left[ \begin{array}{l} \text{Definition of} \\ \text{Matrix multiplication} \end{array} \right]$$

$$\therefore \left( \frac{\partial l_{MLE}(\theta)}{\partial \theta} \right)_j = \left( x^T (y - logistic(x\theta)) \right)_j$$

$$\boxed{\therefore \frac{\partial l_{MLE}(\theta)}{\partial \theta} = x^T \left( y - logistic(x\theta) \right)}$$

Ans: 4.1.2) SGD is slower than GD because n-samples number of python loop iterations of their updates in case of SGD are replaced by a single thto update in case of GD which uses faster numpy ~~operations~~ vectorised operation. Since numpy vector operations in case of GD such as matrix multiplications and subtraction are faster than python loops in SGD, SGD is slower than GD and takes more time for the same number of epochs.