

Q.1.

$\Rightarrow$  Let  $X = \{x_1, x_2, x_3, \dots, x_n\}$  be a set of i.i.d (independent & identically distributed) samples drawn from a normal distribution s.t.

$$X_i \sim \mathcal{N}(\mu, \sigma^2) \quad (1)$$

$\mu$  = mean of the distribution  
 $\sigma$  = variance of the distribution

The probability density function of a normal distribution is:

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (2)$$

Estimate  $\theta_{MLE} = (\mu_{MLE}, \sigma^2_{MLE})$

$$L(\theta) = f(X_1, X_2, \dots, X_n | \theta) = \prod_{i=1}^n f(X_i | \theta) \quad (3)$$

The MLE of  $\theta$  is the value that maximizes  $L(\theta)$ . i.e.  $\theta_{MLE} = \arg \max_{\theta} L(\theta)$

$$L(\mu, \sigma^2) = \prod_{i=1}^n f(x_i | \mu, \sigma^2) \quad (4)$$

Substituting the Gaussian p.d.f.:

$$L(\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \quad (5)$$

Taking the log-likelihood to make differentiation easier,

$$\begin{aligned} \log L(\mu, \sigma^2) &= \sum_{i=1}^n \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \right] \\ &= \sum_{i=1}^n \left[ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x_i-\mu)^2}{2\sigma^2} \right] \\ &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned} \quad (6)$$

Taking the derivative w.r.t  $\mu$

$$\begin{aligned} \frac{\partial}{\partial \mu} \log L &= \frac{\partial}{\partial \mu} \left( -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right) \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n \frac{\partial}{\partial \mu} (x_i - \mu)^2 \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(x_i - \mu) \\ &= -\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \end{aligned} \quad (7)$$

Setting this derivative to zero:

$$\frac{\partial}{\partial \mu} \log L = -\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0.$$

$$\Rightarrow \sum_{i=1}^n (x_i - \mu) = 0.$$

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

(8)

Taking the derivative w.r.t  $\sigma^2$ :

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} \log L &= \frac{\partial}{\partial \sigma^2} \left( -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right) \\ &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned}$$

(9)

Setting this derivative to zero:

$$-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\Rightarrow \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = \frac{n}{2\sigma^2}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

(10)

$\therefore$  (8) and (10) make up the MLE for univariate Gaussian.

Q.2.

=> Given a linear model

$$y(x, w) = w_0 + w_1 x_1 \quad (1)$$

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{t_n - y(x_n, w)\}^2 \quad (2)$$

Substituting (1) into (2)

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{t_n - (w_0 + w_1 x_n)\}^2 \quad (3)$$

To minimize  $E(w)$ , take partial derivatives of  $E(w)$  w.r.t  $w_0$  &  $w_1$  and set to 0.

$$\begin{aligned} \text{A) } \frac{\partial E(w)}{\partial w_0} &= \sum_{n=1}^N \{t_n - (w_0 + w_1 x_n)\} (-1) = 0 \\ &= \sum_{n=1}^N \{t_n - w_0 - w_1 x_n\} = 0 \end{aligned} \quad (4)$$

Simplifying,

$$\sum_{n=1}^N t_n - w_0 N - w_1 \sum_{n=1}^N x_n = 0$$

$$\Rightarrow \sum_{n=1}^N t_n = w_0 N + w_1 \sum_{n=1}^N x_n \quad (5)$$

$$\begin{aligned} \text{B) } \frac{\partial E(w)}{\partial w_1} &= \sum_{n=1}^N \{t_n - (w_0 + w_1 x_n)\} (-x_n) = 0 \\ &= \sum_{n=1}^N \{t_n - w_0 - w_1 x_n\} x_n = 0 \end{aligned}$$

Simplifying,

$$\sum_{n=1}^N t_n x_n - w_0 \sum_{n=1}^N x_n - w_1 \sum_{n=1}^N x_n^2 = 0$$

$$\Rightarrow \sum_{n=1}^N t_n x_n = w_0 \sum_{n=1}^N x_n + w_1 \sum_{n=1}^N x_n^2 \quad (6)$$

The two equations are:

$$w_0 N + w_1 \sum_{n=1}^N x_n = \sum_{n=1}^N t_n$$

$$w_0 \sum_{n=1}^N x_n + w_1 \sum_{n=1}^N x_n^2 = \sum_{n=1}^N t_n x_n$$

In matrix form,

$$\begin{bmatrix} N & \sum_{n=1}^N x_n \\ \sum_{n=1}^N x_n & \sum_{n=1}^N x_n^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} \sum_{n=1}^N t_n \\ \sum_{n=1}^N t_n x_n \end{bmatrix} \quad (7)$$

Comparing (9) w.r.t  $\sum_{j=0}^1 A_{ij} w_j = T_i$

$$A_{00} = N$$

$$A_{01} = \sum_{n=1}^N x_n$$

$$A_{10} = \sum_{n=1}^N x_n$$

$$A_{11} = \sum_{n=1}^N x_n^2$$

$$T_0 = \sum_{n=1}^N t_n$$

$$T_1 = \sum_{n=1}^N t_n x_n$$

$\therefore$  The normal equations for the least squares solutions are:

$$\sum_{j=0}^1 A_{ij} w_j = T_i \quad \text{where}$$

$$A_{ij} = \sum_{n=1}^N (x_n)^{i+j}$$

$$T_i = \sum_{n=1}^N (x_n)^i t_n$$

Q. 3.

=> For a binary classification problem, cross-entropy loss function is:

$$L_{CE} = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]; \quad (1)$$

$y \in \{0,1\}$  is true label  
 $\hat{y}$  is predicted probability.

$$\begin{aligned} \frac{\partial L_{CE}}{\partial \hat{y}} &= -\left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}\right) \\ &= -\frac{\hat{y} - y}{\hat{y}(1-\hat{y})} \end{aligned} \quad (2)$$

The least square error function is defined as:

$$L_{LSE} = \frac{1}{2} |y - \hat{y}|^2; \quad (3)$$

$y$  is true label  
 $\hat{y}$  is the predicted probability.

$$\frac{\partial L_{LSE}}{\partial \hat{y}} = \hat{y} - y \quad (4)$$

The sigmoid function is defined as:

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}} \quad (5)$$

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)[1 - \sigma(z)] = \hat{y}(1-\hat{y}) \quad (6)$$

Differentiating (2) w.r.t  $z$ ,

$$\begin{aligned} \frac{\partial L_{CE}}{\partial z} &= \frac{\partial L_{CE}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \\ &= -\frac{\hat{y} - y}{\hat{y}(1-\hat{y})} \cdot \hat{y}(1-\hat{y}) \\ &= \hat{y} - y \end{aligned} \quad (7)$$

This shows that when using sigmoid function  $\hat{y} = \sigma(z)$ , the derivative of CE loss and least squares error are the same.

Q. 4.

=> for a discrete random variable  $X$  that can take  $k$  different categories, the multinomial p.m.f is:

$$P(x_1, x_2, \dots, x_k | n, p) = \frac{n!}{x_1! x_2! \dots x_k!} \prod_{k=1}^N p_k^{x_k} ; \text{ where} \quad (1)$$

$x_k$  : number of times category  $k$  is observed

$p_k$  : probability of category  $k$

$n$  :  $\sum_{k=1}^N x_k$  is total number of trials

$p$  :  $(p_1, p_2, \dots, p_k)$  vector of probabilities

When  $n=1$ , the distribution becomes a Categorical distribution

For a single observation  $X$ , which has  $k$  categories,  $x = (x_1, x_2, \dots, x_k)$  is a one-hot vector where only one of the  $x_i$  is 1 and rest are 0.

Substituting  $n=1$  into (1)

$$P(X=x | p) = \frac{1!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k} \quad (2)$$

Since, only one  $x_i = 1$  and rest are 0,

$$P(X=x | p) = p_i ; \text{ where } x_i = 1 \quad (3)$$

The likelihood function is given by:

$$L(p | x_1, x_2, \dots, x_m) = \prod_{j=1}^m p_{i_j} ; \text{ } i_j \text{ is the observed category for sample } j.$$

$$L(p | X) = \prod_{j=1}^m \prod_{i=1}^k p_i^{x_{ji}} ; \text{ summing over all samples gives total count of each category } i \text{ as } x_i = \sum_{j=1}^m x_{ji}$$

$$L(p | X) = \prod_{i=1}^k p_i^{x_i}$$

Q. 5.

5.

$\Rightarrow$  Given,

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

Let  $y = \sigma(x) = \frac{1}{1+e^{-x}}$  (1)

Solving for  $x$ ,

$$y(1+e^{-x}) = 1$$

$$\text{or } y + ye^{-x} = 1$$

$$\text{or } ye^{-x} = 1-y$$

Taking Natural log on both sides,

$$\ln(ye^{-x}) = \ln(1-y)$$

$$\text{or } \ln y + \ln e^{-x} = \ln(1-y)$$

$$\text{or } \ln y + (-x) = \ln(1-y)$$

$$\text{or } x = \ln y - \ln(1-y)$$

$$\text{or } \boxed{x = \ln \left( \frac{y}{1-y} \right)} \quad \text{span style="color: red;">(2)}$$

$$\text{logit}(y) = \ln \left( \frac{y}{1-y} \right)$$

$$\text{From (2), } x = \ln \left( \frac{y}{1-y} \right) = \text{logit}(y)$$

$\therefore$  The inverse of the sigmoid function is the logit function.

# References

1.2 - Maximum Likelihood Estimation | STAT 415

20\_mle\_annotated, Jerry Cain

Gaussian Distribution and Maximum Likelihood Estimate Method (Step-by-Step) | by Anel Music | The Startup | Medium

Lecture 6: The Method of Maximum Likelihood for Simple Linear Regression, CMU

Lecture 8: Properties of Maximum Likelihood Estimation (MLE), Purdue

Maximum Likelihood Estimation for Gaussian Distributions - Programmatically

Maximum likelihood estimation for the univariate Gaussian | The Book of Statistical Proofs