

2) a) For a given set of linear Equations —

$$Ax = b, \text{ where}$$

x = vector

A = matrix and

b = vector.

A and b are the known quantities, whereas x is the unknown for which we find the solution.

If A is a square matrix, then number of equations is equal to number of unknowns and if A is an invertible matrix then solution of ax can be found out as $x = A^{-1}b$. Thus in both the above cases calculating x is easier.

But for non square A matrices, calculating inverse is not possible.

So, Pseudo Inverse of a matrix (Moore-Penrose pseudo inverse) helps us to approximately invert the matrix and find a best fit for x that comes closer to solving the equation.

If A is invertible, then pseudo inverse of the matrix (A^+) is equal to A^{-1} and if A is not invertible, then also A^+ is defined.

$$A^+ A^T (A A^T)^{-1} = I \quad \leftarrow \quad \|x - d\|$$

$$\text{Result above holds if } A^+ A^T (A A^T)^{-1} \text{ exists}$$

2) a) i) Underdetermined system of equations :-

here, number of equations are less than no. of unknowns and thus, infinitely many solution will exist for x (where x is the unknown vector in a given system of linear equation, $Ax = b$). Also A is a short fat matrix and we cannot find its inverse.

There are not enough measurements in b to uniquely determine a single unique solution. We can pick one of these solutions by finding the smallest i.e; we will minimize x .

$$x = A^T (A A^T)^{-1} b \quad \text{where,}$$

$(A^T (A A^T)^{-1} b)$ is called pseudo-inverse.

ii) Overdetermined system of equations :-

here, no. of equations are more than the unknowns and we have no solution for x . Here A is a tall matrix and thus we cannot find its inverse.

We try to find an x which will minimize the error between the fit of Ax and b .

$$\| b - Ax \|^2 \Rightarrow x = (A^T A)^{-1} A^T b$$

where $(A^T A)^{-1} A^T$ is called pseudo inverse.

2) b)

$$\begin{aligned} n_1 + 3n_2 &= 17 \\ 5n_1 + 7n_2 &= 19 \\ 11n_1 + 13n_2 &= 23 \end{aligned}$$

this is an overdetermined system of equations

representing in matrix form -

$$Ax = y, \text{ where}$$

$$A = \begin{bmatrix} 1 & 3 \\ 5 & 7 \\ 11 & 13 \end{bmatrix}, x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, y = \begin{bmatrix} 17 \\ 19 \\ 23 \end{bmatrix}$$

here A is a non-invertible matrix, so we use pseudo inverse to find out x as follows -

$$x = A^+ y \text{ where,}$$

$$A^+ = \text{pseudo inverse} = (A^T A)^{-1} A^T$$

$$\text{so } x = (A^T A)^{-1} A^T y$$

$$A = \begin{bmatrix} 1 & 3 \\ 5 & 7 \\ 11 & 13 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 1 & 5 & 11 \\ 3 & 7 & 13 \end{bmatrix}$$

$$A^T A = 181 \begin{bmatrix} 1 & 5 & 11 \\ 3 & 7 & 13 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 5 & 7 \\ 11 & 13 \end{bmatrix}$$

$$= \begin{bmatrix} 1 + 25 + 121 & 3 + 35 + 143 \\ 3 + 35 + 143 & 9 + 49 + 169 \end{bmatrix}$$

$$= \begin{bmatrix} 147 & 181 \\ 181 & 227 \end{bmatrix}$$

$$A^T Y = \begin{bmatrix} 801 & 185 & 1180 \\ 803 & 17 & 13 \end{bmatrix} \begin{bmatrix} 17 \\ 19 \\ 23 \end{bmatrix}$$

$$= \begin{bmatrix} 17 + 95 + 253 \\ 51 + 133 + 299 \end{bmatrix}$$

$$= \begin{bmatrix} 365 \\ 483 \end{bmatrix} = \begin{bmatrix} 18 \\ 58 \end{bmatrix}$$

$$(A^T A)^{-1} = \frac{\text{adj}(A^T A)}{|A^T A|}$$

$$\text{adj}(A^T A) = \begin{bmatrix} 227 & -181 \\ -181 & 227 \end{bmatrix}$$

$$|A^T A| = \begin{vmatrix} 147 & 181 \\ 181 & 227 \end{vmatrix}$$

$$= 147 \times 227 - 181 \times 181 = 470$$

$$= 608$$

$$(A^T A)^{-1} = \begin{bmatrix} 227/608 & -181/608 \\ -181/608 & 147/608 \end{bmatrix}$$

multiplying, $(A^T A)^{-1}$ with $A^T y$ to get x

$$= \begin{bmatrix} 227/608 & -181/608 \\ -181/608 & 147/608 \end{bmatrix} \begin{bmatrix} 365 \\ 483 \end{bmatrix}$$

$$= \begin{bmatrix} -7.51 \\ 8.12 \end{bmatrix}$$

$$\text{so, } \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -7.51 \\ 8.12 \end{bmatrix} \quad \underline{\text{ans}}$$

$$(A^T A)^{-1} = -(A^T A)$$

$$|A^T A|$$

$$\begin{bmatrix} 181 & 181 \\ 181 & 227 \end{bmatrix} = |A^T A|$$

2)c);) The estimated model of linear regression can be written as -

$$Y = X\theta + \varepsilon \text{, where}$$

if m = number of rows

n = number of columns then

Y is a $m \times 1$ and Y is a matrix of size $m \times 1$

Y is a $(m \times 1)$ matrix of target variables,

X is a $(m \times n)$ matrix of features,

θ is a $(n \times 1)$ matrix of parameters and

ε is a $(m \times 1)$ matrix of residuals.

$$\varepsilon = Y - X\theta$$

To find the least square estimator, we write sum of square of residuals as follows -

$$S(\theta) = \sum \varepsilon_i^2 = \varepsilon^T \varepsilon = (Y - X\theta)^T (Y - X\theta)$$

$$= Y^T Y - Y^T X\theta - \theta^T X^T Y + \theta^T X^T X\theta$$

minimizing $S(\theta)$ ie setting derivative to 0,

$$-2X^T Y + 2X^T X\theta = 0$$

$$\Rightarrow \theta = (X^T X)^{-1} X^T Y; \text{ this is the}$$

normal equation

2) c) ii) We prefer iterative methods like gradient Descent rather than using closed form solutions because gradient descent is computationally cheaper. For example, if number of parameters is 10000 then we have to compute 10000^2 values for $\mathbf{x}^T \mathbf{x}$ in closed form equation, which becomes very expensive to compute parameters.

2) d) ii) RMSE :-

RMSE stands for root mean squared error.

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2} \text{ where}$$

N = total number of data points

\hat{y}_i = predicted target value

y_i = actual target value

It measures the distance between the predicted and actual values i.e. the residuals. It measures the standard deviation of the residuals i.e. how spread out the residuals are or how concentrated the data is around the line of best fit. Gives high weightage to large errors hence should be used when large errors are undesirable. Lower value means better fit of the model. Range from 0 to ∞ .

MAE :-

MAE stands for mean absolute error.

$$\text{MAE} = \frac{1}{N} \sum_{j=1}^N |y_j - \hat{y}_j| \text{ where}$$

N = total number of data points

\hat{y}_j = predicted target value

y_j = actual target value

Its value ranges from 0 to ∞ .

It measures absolute average distance between real data and predicted data. Here, all errors are given equal weightage; MAE fails to give weightage to large errors. Hence must be used when large errors are not present.

MSE :-

MSE stands for mean squared error.

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

where

N = total number of data points

y_i = actual target value

\hat{y}_i = predicted target value

It ranges between 0 to ∞ . MSE puts large weightage on outlier errors and squaring is done if it magnifies the error. Hence must be used when the model has less or no outliers.

However if we want to ignore outliers then MSE will be a bad choice.

3) c) tanh n is defined as

3) c) tanh n is defined as \rightarrow

$$\tanh n = \frac{e^n - e^{-n}}{e^n + e^{-n}} = g(n)$$

$$\text{let } h_\theta(n) = g(\theta^\top n) = \frac{e^{\theta^\top n} - e^{-\theta^\top n}}{e^{\theta^\top n} + e^{-\theta^\top n}}$$

where θ = parameters of the model

n = independent features of the model

calculating the derivative of $\tanh(z) \rightarrow$

$$\frac{d}{dz}(g(z)) = \frac{d}{dz} \left(\frac{e^z - e^{-z}}{e^z + e^{-z}} \right)$$

$$\Rightarrow (e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z}) \\ (e^z + e^{-z})^2$$

$$= 1 - \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$= 1 - \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$= 1 - \tanh^2(z)$$

$$= 1 - (g(z))^2 \dots \textcircled{1}$$

in logistic regression, we predict the probability of any of the class, the ^{probability of} other class can be calculate with the 1st one.

$$P(y=1 | x; \theta) = h_{\theta}(x) = \frac{e^{\theta^T x}}{e^{\theta^T x} + e^{-\theta^T x}}$$

$$\text{then } P(y=0 | x; \theta) = 1 - P(y=1 | x; \theta)$$

$$= 1 - \left(\frac{e^{\theta^T x} - e^{-\theta^T x}}{e^{\theta^T x} + e^{-\theta^T x}} \right)$$

combining the two equations —

$$P(y | x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

probability of y given x parameterised by θ

$$L(\theta) = \text{likelihood of } \theta = L(\theta)$$

maximizing $L(\theta)$ so that data points best fit the model

$$\max(L(\theta)) = \max \left(P(\vec{y} | x; \theta) \right)$$

$$[m = \text{no. of instances}] = \prod_{i=1}^m P(y^i | x^i; \theta)$$

$$= \prod_{i=1}^m h_{\theta}(x^i)^{y^i} (1 - h_{\theta}(x^i))^{1-y^i}$$

maximising $l(\theta)$ = maximising $l(\theta)$ curves
 $\lambda(\theta) = \log$ likelihood function thus

$$\max \lambda(\theta) = \max \log l(\theta) =$$

$$\sum_{i=1}^m \log (h_\theta(x_i))^{y_i} + \sum_{i=1}^m \log (1 - h_\theta(x_i))^{1-y_i}$$

here, $\theta = \theta_0 + \alpha \nabla_{\theta_0} l(\theta)$, we are
 maximising $\lambda(\theta)$, i.e. gradient ascent.

$$\nabla_{\theta_j} l(\theta_j)$$

$$= \frac{\partial}{\partial \theta_j} (y_i \log (h_\theta(x_i))) + (1-y_i) \log (1 - h_\theta(x_i))$$

$$= \frac{y_i}{h_\theta(x_i)} \frac{\partial}{\partial \theta_j} h_\theta(x_i) - \frac{1-y_i}{1-h_\theta(x_i)} \frac{\partial}{\partial \theta_j} h_\theta(x_i)$$

$$= \left(\frac{y_i}{h_\theta(x_i)} - \frac{1-y_i}{1-h_\theta(x_i)} \right) \frac{\partial}{\partial \theta_j} h_\theta(x_i)$$

$$= \left(\frac{y_i}{h_\theta(x_i)} - \frac{1-y_i}{1-h_\theta(x_i)} \right) [1 - h_\theta(x_i)^2] x_j$$

$$= \frac{y_i - y_i h_\theta(x_i)}{h_\theta(x_i) \times (1 - h_\theta(x_i))} \times$$

$$[1 - h_\theta(x_i)^2] \times \frac{y_i}{x_j}$$

$$= \frac{y^i - h_\theta(x^i)}{h_\theta(x^i) \times (1 - h_\theta(x^i))} \times (1 + h_\theta(x^i)) \times (1 - h_\theta(x^i)) \times \frac{x}{x^i} \times \frac{u_j^i}{u_j}$$

$$= \frac{y^i - h_\theta(x^i)}{h_\theta(x^i)} \times (1 + h_\theta(x^i)) \times \frac{u_j^i}{u_j}$$

thus the update rule for logistic regression =

$$\boxed{\theta_j = \theta_j + \alpha \left(\frac{y^i - h_\theta(x^i)}{h_\theta(x^i)} \times u_j^i \times (1 + h_\theta(x^i)) \right)}$$

→ update rule for logistic regression using tanh function as the decision boundary