

ML Assignment - 1

Ritisha-Gupta, MT22056

Quesy) g) What is psuedo inverse of a matrix?

Solⁿ - If A is a square matrix ~~of~~ then its inverse exists and A is known as an invertible matrix.

If A is not an invertible matrix, then we find the Moore Penrose psuedo inverse. However, Moore Penrose psuedo inverse is defined for invertible matrix also but in that case both are equal.

Let $A \in \mathbb{R}^{m \times n}$. & let $B \in \mathbb{R}^{n \times m}$ be the psuedo-inverse of A if it satisfies all four conditions:

- 1) $ABA = A \leftarrow B$ is generalized inverse of A .
- 2) $BAB = B \leftarrow A$ is generalized inverse of B .
- 3) $(AB)^T = AB \leftarrow AB$ is symmetric.
- 4) $(BA)^T = BA \leftarrow BA$ is symmetric.

The psuedo inverse of a matrix $A \in \mathbb{R}^{m \times n}$ always exists & is unique. Denoted as (A^+) .

The symmetric form of defⁿ implies $B = A^+$ & $A = B^+$ & thus.
 $A = (A^+)^+$.

(i) Undetermined System of Equation.

- In such system we have fewer equations than no of variables.
- It cannot have unique solⁿ.
- In this case either ∞ many or no solⁿ's (consistent).

~~$$x = (\theta^T \theta^T \theta^T)^{-1} y$$~~

where $\theta^T (\theta^T \theta)^{-1}$ = psuedo inverse.

$$x = \theta^T (\theta \theta^T)^{-1} y$$

(ii) Overdetermined.

- more equation than no. of variables
- also ∞ many or no solⁿ.
- may have a unique solⁿ.

$$x = (O^T O)^{-1} O^T y$$

Where $\odot \odot$ pseudo inverse.

b) $x_1 + 3x_2 = 17$

$5x_1 + 7x_2 = 19$

$11x_1 + 13x_2 = 23$

$A\vec{x} = \vec{B}$ { sys of linear equation }

(Coeff matrix) $A = \begin{bmatrix} 1 & 3 \\ 5 & 7 \\ 11 & 13 \end{bmatrix}_{m \times n}$

$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{n \times 1}$

$\vec{B} = \begin{bmatrix} 17 \\ 19 \\ 23 \end{bmatrix}_{m \times 1}$

\vec{x}, \vec{B} = vector
 A = matrix

Here $m > n$ ($m=3$ & $n=2$) so ~~over~~ eq^n we have over-determined system of linear eq. because we have more ~~not~~ no. of equations than variables.

$AX = B$

Multiply both sides with A^T .

$A^T A \vec{x} = A^T \vec{B}$

Multiply both sides with $(A^T A)^{-1}$

$(A^T A)^{-1} (A^T A) \vec{x} = (A^T A)^{-1} A^T \vec{B}$

We know that $A A^{-1} = I$ (Identity matrix)

So, $I \cdot \vec{x} = (A^T A)^{-1} A^T \vec{B}$

$\vec{x} = (A^T A)^{-1} A^T \vec{B}$
 $\vec{x} = A^+ \vec{B}$

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

First we'll find pseudo inverse. $A^+ = (A^T A)^{-1} A^T$

$A^T A = \begin{bmatrix} 1 & 5 & 11 \\ 3 & 7 & 13 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 5 & 7 \\ 11 & 13 \end{bmatrix} = \begin{bmatrix} 147 & 181 \\ 181 & 227 \end{bmatrix}$

$$(A^T A)^{-1} = \begin{bmatrix} +0.3733 & -0.2976 \\ -0.2976 & 0.2417 \end{bmatrix} = \frac{\text{adj}(A^T A)}{|A^T A|}$$

$$(A^T A)^{-1} A^T = \begin{bmatrix} 0.3733 & -0.2976 \\ -0.2976 & 0.2417 \end{bmatrix}_{2 \times 2} \begin{bmatrix} 1 & 5 & 11 \\ 3 & 7 & 13 \\ 4 & 8 & 14 \end{bmatrix}_{3 \times 3} = \begin{bmatrix} 1.5 & 11 \\ 3 & 7 & 13 \end{bmatrix}_{2 \times 3}$$

$$A^+ = (A^T A)^{-1} A^T = \begin{bmatrix} -0.5195 & -0.2167 & 0.2375 \\ 0.4275 & 0.2039 & -0.1315 \end{bmatrix}$$

(pseudo-inverse)

$$\vec{x} = A^+ \vec{b}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -0.5195 & -0.2167 & 0.2375 \\ 0.4275 & 0.2039 & -0.1315 \end{bmatrix}_{2 \times 3} \begin{bmatrix} 17 \\ 19 \\ 23 \end{bmatrix}_{3 \times 1}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -7.4863 \\ 8.1171 \end{bmatrix} \quad \underline{\text{Ans}}$$

$$x_1 = -7.4863$$

$$x_2 = 8.1171$$

Ques 3) Gradient descent is an iterative method to find the local minima / maxima of a function. It is used to find the best parameters that minimizes the model's cost function. Normal eqⁿ is a closed form solution for Linear Regression. Using this Normal eqⁿ, we can directly solve for optimal value of parameter θ in just one step. Optimal parameter can be obtained by just formula.

$$\theta = (A^T A)^{-1} A^T \vec{y}$$

where θ = hypothesis parameter.

A = Input matrix

\vec{y} = output

- (ii) We prefer iterative methods like Gradient descent rather than normal eqⁿ because in closed form solution of 'N' independent variables requires to find inverse of matrix which takes app $O(n^3)$. Its computational complexity is very high, whereas in gradient descent, it has a time complexity that is linear i.e. $O(n)$.

Normal eqⁿ is very easy for univariate variables but not in case of multivariate variables.

Q: $\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}} = g(x)$

We know that $h_0(x) = \text{hypothesis}$

$$h_0(x) = g(\theta^T x) = \frac{e^{\theta^T x} - e^{-\theta^T x}}{e^{\theta^T x} + e^{-\theta^T x}}$$

For
logistic
regression.

Where $\theta = \text{parameters}$, $x = \text{input features}$

Derivative of $\tanh(z) = \frac{\partial}{\partial z} (g(z)) = \frac{\partial}{\partial z} \left(\frac{e^z - e^{-z}}{e^z + e^{-z}} \right)$

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

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

$$1 - \left(\frac{e^z - e^{-z}}{e^z + e^{-z}} \right)^2 = 1 - \tanh^2(z) = 1 - (g(z))^2 \quad \text{--- (1)}$$

In logistic regression, $P(y=1|x; \theta) = h_0(x) = \frac{e^{\theta^T x} - e^{-\theta^T x}}{e^{\theta^T x} + e^{-\theta^T x}}$

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 -

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

probability of y given x parameterized by θ

We know that likelihood is ~~defined~~ represented as $L(\theta)$.

maximising $L(\theta)$ so that we get best fit model.

$$\max(L(\theta)) = \max \left[P = (y|x; \theta) \right]$$

$$= \prod_{i=1}^m P(y^i/x^i; \theta) \quad m = \text{no of i/p features}$$

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

maximising $L(\theta) = \text{maximising } \ln L(\theta)$

Taking log on both sides to make our calculation easy & convert into summation.

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

here $\theta = \theta + \frac{\partial}{\partial \theta} L(\theta) \rightarrow$ maximising gradient-descent.

$$\frac{\partial}{\partial \theta_j} = (y^j \log(h_\theta(x^j)) + (1 - y^j) \log(1 - h_\theta(x^j)))$$

$$= \frac{y^j}{h_\theta(x^j)} \frac{\partial h_\theta(x^j)}{\partial \theta_j} - \frac{1 - y^j}{1 - h_\theta(x^j)} \frac{\partial (h_\theta(x^j))}{\partial \theta_j}$$

$$\left(\frac{y^j}{h_\theta(x^j)} - \frac{1 - y^j}{1 - h_\theta(x^j)} \right) \frac{\partial h_\theta(x^j)}{\partial \theta}$$

$$\left(\frac{y^j}{h_\theta(x^j)} - \frac{1 - y^j}{1 - h_\theta(x^j)} \right) [1 - h_\theta(x^j)^2] \times x_{j,k}^j$$

$$\left(\frac{y^j - y^j h_\theta(x^j) - h_\theta(x^j) + h_\theta(x^j) y^j}{h_\theta(x^j) \cdot (1 - h_\theta(x^j))} \right) (1 - h_\theta(x^j)^2) \cdot x_{j,k}^j$$

$$\frac{y^j - h_0(x^j)}{h_0(x^j) \cdot (1 - h_0(x^j))} \times (1 + h_0(x^j)) \cdot (1 - h_0(x^j))$$

$$= \frac{y^j - h_0(x^j)}{h_0(x^j)} \times (1 + h_0(x^j)) \times x_k^j$$

Thus the updated rule for logistic reg =

$$\left[\theta_k = \theta_k + \alpha \left(\frac{y^j - h_0(x^j)}{h_0(x^j)} \cdot x_k^j \times (1 + h_0(x^j)) \right) \right]$$

→ using tanh function as deen boundary.