



# ASSIGNMENT - I

(a) (i) The model is made to find  $\hat{y}$  such that they are closest to true values, as squared errors make a much bigger difference in  $|\hat{y} - y|$ .

$$(b) J(\beta) = \frac{1}{2} \|y - X\beta\|^2 = \frac{1}{2} (y - X\beta)^T (y - X\beta)$$

$$J(\beta) = \frac{1}{2} (y^T y - y^T X\beta - (X\beta)^T y + (X\beta)^T (X\beta))$$

$$J(\beta) = \frac{1}{2} (y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta)$$

$$J(\beta) = \frac{1}{2} (y^T y - 2\beta^T X^T y + \beta^T X^T X\beta)$$

$$\nabla_{\beta} J(\beta) = \nabla_{\beta} \left[ \frac{1}{2} y^T y - \beta^T (X^T y) + \frac{1}{2} \beta^T (X^T X) \beta \right]$$

$$\nabla_{\beta} J(\beta) = 0 - X^T y + \frac{1}{2} \cdot 2 (X^T X) \beta$$

$$\nabla_{\beta} J(\beta) = -X^T y + X^T X \beta = 0$$

$$\therefore X^T y = X^T X \beta$$

$$\checkmark \hat{\beta} = (X^T X)^{-1} X^T y$$

normal equation.

(C) Computational costs are very high.  
If features are linearly dependent,  $X^T X \rightarrow$  singular  
and a sol<sup>n</sup> would not exist.

why iterative methods are preferred?  $\rightarrow$   
- faster than inversion.

Q.2 Core idea: - to compute the gradient of loss  
function ( $L$ ) w.r.t. every weight & bias in a  
neural network. These help in optimisation to  
reduce the loss.

Chain Rule:-  $\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial \text{output}} \times \frac{\partial \text{output}}{\partial \text{hidden layer}} \times \dots \times \frac{\partial \text{local sum}}{\partial w_{ij}}$

use of partial derivatives scales linearly with the no  
of layers & connections.

$$(a) \quad z_1 = w_1 x + b_1 \quad z_2 = w_2 a_1 + b_2 \\ a_1 = \sigma(z_1) \quad a_2 = \sigma(z_2) = \hat{y}$$

$$L = -(y \log a_2 + (1-y) \log (1-a_2))$$

$$\frac{\partial L}{\partial z_2} = a_2 - y$$

$$\sigma'(z) = \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$$

$$\frac{\partial a_2}{\partial z_2} = a_2(1 - a_2)$$

$$\frac{\partial h}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$$

$$z_2 = w_2 a_1 + b_2 \Rightarrow \frac{\partial z_2}{\partial w_2} = a_1$$

$$\frac{\partial h}{\partial w_2} = (a_2 - y) a_1 \quad \frac{\partial h}{\partial b_2} = \frac{\partial h}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2}$$

$$z_2 = w_2 a_1 + b_2$$

$$\Rightarrow \frac{\partial z_2}{\partial b_2} = 1 \quad \frac{\partial h}{\partial b_2} = a_2 - y$$

2. Propagated Error :-  $\frac{\partial h}{\partial a_1}$

$$z_2 = w_2 a_1 + b_2 \quad \frac{\partial z_2}{\partial a_1} = w_2 \Rightarrow \frac{\partial h}{\partial a_1} = (a_2 - y) w_2$$

$$\frac{\partial a_1}{\partial z_1} = a_1 (1 - a_1) = \sigma(z_1) \Rightarrow \frac{\partial h}{\partial z_1} = (a_2 - y) w_2 a_1 (1 - a_1)$$

$$\frac{\partial z_1}{\partial w_1} = n \quad \frac{\partial L}{\partial w_1} = (a_2 - y) w_2 a_1 (1 - a_1) x$$

$$\frac{\partial z_1}{\partial b_1} = 1 \quad \frac{\partial h}{\partial b_1} = (a_2 - y) w_2 a_1 (1 - a_1)$$

Q.3 (a) ANN input processing:-

→  $x_{i+1}$  is independent of  $x_i$

→ information flows in one dir<sup>n</sup>.

RNN sequence processing:-

→ information flows in a loop, there is updation in hidden state for any input.

(b) During BPTT, the gradients are repeatedly multiplied by weight matrices of the Recurrent layer. The gradient vanishes for weights connected to early time steps. These early weights receive minimal updates, thus the network is unable to learn or retain info from many steps back.

(c) Gates regulate flow of info into & out of the cell state

(d) LSTMs address it through cell state ( $C$ ).

(e) ANN: image classification.

RNN: speech recognition.

LSTM: machine translation.

Q.4 (a) Example - subject-verb agreement.  
→ vanishing gradient problem

(b) The memory cell & gating mechanism allow the LSTM to retain important info over long sequences by providing a path for flow.