

1)

$$\hat{y}_c = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_d x_{1d}$$

$$J(\beta) = \frac{1}{2} \|y - X\beta\|^2 = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

a)

Here regression is finding a relation b/w input variable x & target y , it is done by fitting a best fit 2D st. line or 3D plane. Goal is to find the set of β . Why minimizing squared errors? Because it is easy to compute derivatives and it is smooth parabola. Also, this makes it to prioritize large errors to bring them closer.

b)

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

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

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

since both are scalar.

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

$$\frac{\partial J}{\partial \beta} = -2 y^T X + 2 X X^T \beta = 0$$

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

$$\beta = (X^T X)^{-1} X^T Y$$

- (i) finding inverse is very computationally expensive. Also meeting the conditions for invertibility also becomes difficult.
- iterative method \rightarrow it is easy because the time it would take would be $n \times n$ one step, much less than the direct one.

- (2) a) core idea of back propagation:
- prediction using forward pass
 - calculating how error changes with prediction (backpropagation)
 - error propagation
 - update weight

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

$$= (a_2 - y)(a_1)$$

$$\frac{\partial L}{\partial b_2} = (a_2 - y)(1)$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial x} = (a_2 - y) \cdot (a_1) \cdot \frac{\partial a_1}{\partial x} = (a_2 - y) \cdot (a_1) \cdot (1 - a_1) \cdot x$$

$$\frac{\partial L}{\partial x} = (a_2 - y) \cdot (a_1) \cdot (1 - a_1) \cdot x$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_1}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

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

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

c) gradient descent algorithm

$w_1 \rightarrow w_1 - \eta \frac{\partial L}{\partial w_1}$ mit $\eta \rightarrow$ learning rate,

$$b_1 \rightarrow b_1 - \eta \frac{\partial L}{\partial b_1}$$

$$w_2 \rightarrow w_2 - \eta \frac{\partial L}{\partial w_2}$$

je nachdem welche Variable abweichen darf

$$b_2 \rightarrow b_2 - \eta \frac{\partial L}{\partial b_2}$$

separat vom anderen Bereich

3) a) ANN vs RNN

ANN \rightarrow processes input independently.
They just need to maintain the previous
do not inputs.

RNN \rightarrow processes inputs sequentially.
They have to save the previous
memory also.

b) RNN struggle with long term dependencies due to vanishing of gradient as gradients shrink exponentially during backpropagation through time so it becomes hard to learn from early information too small to learn from

c) Forget gate: \rightarrow decides the irrelevant info.
Input gate: \rightarrow decide which new info. is to be taken as input
Output gate: \rightarrow decides the result that's passed onto next layer.

d)

Using gated cells.

e)

ANN → Image classification

RNN → Speech Recognition

LSTM → Language Translation

4) a) The flights ~~that~~ that passengers booked with Indigo were ~~the~~ ~~the~~ ~~not~~ affected.
were → flights

Yes, RNN would struggle due to vanishing gradient problem.

b)

Suppose the task is to process a paragraph ~~so~~. So in that after full stop it is important to wipe out the previous ~~existing~~ memory for model to work effectively.