

1. Linear regression is a supervised machine learning algorithm use to predict a continuous variable by using the training data that is fed into the model.

Linear regression assumes a linear relation. It fits the input features  $x \in \mathbb{R}^d$  and output  $y$ , goal is to find parameters  $\beta$  such that  $\hat{y}$  (prediction) are closest to  $y$  (outputs).

- \* Squaring penalizes large errors more than small ones
- \* Ensure the fn is differentiable, hence able to take gradient
- \* Direct summation may ~~lead to~~ allow  $+ve$  &  $-ve$  terms to cancel each other. squaring ensures all errors are  $pos$

$$(b). J = \frac{1}{2} \|y - XB\|^2$$

$$= f(y - XB)^T (y - XB)$$

$$= \frac{1}{2} (y^T - B^T x^T)(y - XB)$$

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

$$\therefore \cancel{(y^T X B)^T} = B B^T x^T y = \text{scalar}$$

$$\nabla J(\beta) = \nabla \left( \frac{1}{2} (y^T y - 2 B^T x^T y + B^T X^T X B) \right)$$

$$\frac{\partial (X^T \alpha)}{\partial x} = 0$$

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

$$\Rightarrow -X^T y + X^T X B = 0$$

$$X^T X B = X^T y$$

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

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad x = R^{n \times d+1}$$

$$\beta = R^{d+1 \times 1}$$

(d). Direct convolution of  $X^T X$  in the normal equation can be problematic as taking inverse of a high dimensional matrix is computationally intensive as it takes both <sup>long</sup> time & memory.

Gradient descent takes small steps rather than take inverse of whole matrix at once, updates the other parameters of model constantly without starting calculation from start.

2.

(a). Backpropagation is an algorithm used to calculate gradients by propagating the error backward from output to the input layer.

(b). Chain rule helps to break down complex derivative of loss function into simple steps.

(b).

$$\frac{\partial L}{\partial z_2} = q_2 - y$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} \quad : \quad z_2 = w_2 a_1 + b_2 \\ \frac{\partial z_2}{\partial w_2} = a_1 \\ = (q_2 - y) a_1$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2}$$

$$z_2 = w_2 a_1 + b_2$$

$$\frac{\partial z_2}{\partial b_2} = 1$$

$$\frac{\partial L}{\partial b_2} = q_2 - y$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \leftarrow \frac{\partial q_1}{\partial z_1} = q_1(1-q_1) \\ = (q_2 - y) w_2$$

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

$$= (\alpha_2 - y) w_2 \alpha_1 (1 - \alpha_1) x$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

$$= (\alpha_2 - y) w_2 \alpha_1 (1 - \alpha_1)$$

- ④ To update the weights we subtract a fraction of the gradient from the current value

$$w_1' = w_1 - \alpha \frac{\partial L}{\partial w_1}$$

$$w_2' = w_2 - \alpha \frac{\partial L}{\partial w_2}$$

$$b_1' = b_1 - \alpha \frac{\partial L}{\partial b_1}$$

$$b_2' = b_2 - \alpha \frac{\partial L}{\partial b_2}$$

$\alpha \rightarrow$  controls the step size

If  $\alpha$  is too small  
↳ convergence very slow

If  $\alpha$  is too large  
↳ may overshoot the minimum

### Q3.

- ① (a). RNN  $\rightarrow$  RNN processes inputs sequentially if maintains a hidden state (memory) the output at step  $t$  depends both on current input & step  $t-1$ .

- (b). ANN  $\rightarrow$  processes inputs independently

- ② Similar RNN's struggle with long term dependencies due to the problem of vanishing gradient.

- ③ LSTMs uses gates to control information flow unlike RNN
- Forget gate  $\rightarrow$  decides what irrelevant info to throw away
  - Input gate  $\rightarrow$  decides what new important info to store
  - Output gate  $\rightarrow$  decides what info is relevant for prediction

(ii). LSTM uses cell state to address the problem of vanishing gradient

- (Q) ANN → House price prediction  
RNN → Time series forecasting  
LSTM → Machine Translation

Q4.

Example of long range dependency

(a). The ~~boy~~ boys who play Soccer in the rain were scores

The model must connect the plural subject boys at the start of the plural verb were at the end

\* Standard RNN would struggle, by the time it processes the sentence the vanishing gradient problem ~~problem~~ would cause <sup>it</sup> to forget that the subject was "boys". It might incorrectly predict "ies".

(b). LSTM use a memory cell that carries relevant context across time steps. Gates control flow, -  
input gate decides what new <sup>into</sup> to store, forget gate decides what old <sup>into</sup> to loss thru output gate, -  
decides what to give as output. This gating prevents vanishing gradients & enables long term dependency learning.