

Ans - I.

(i) (a) Regression builds relationship between input features and output. It aims to predict the output using linear relationship $y = XB + \epsilon$ in case of Linear regression.

- Due to squaring, there is more impact on residuals.
- Cost function becomes smooth function, which ensures a global minimum point.

$$(6) J(B) = \frac{1}{2} \|y - XB\|^2$$

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

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

$$\frac{\partial J}{\partial B} = 0 \rightarrow \nabla J(B) = X^T XB - X^T y$$

$$\begin{cases} \nabla(y^T XB) = X^T y \\ \nabla(B^T X^T XB) = 2X^T XB \end{cases} \rightarrow \nabla J(B) = 0$$

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

(c) Direct inversion involves ~~calculation~~ high computational cost calculations like calculating $X^T X$ and then inverse of it, so computationally expensive.

Gradient descent (iterative method) is preferred as it minimises cost function iteratively which is good for real world datasets. Here no matrix inversions are required.

(Q2.) Backpropagation is an algorithm which iteratively
 (i) updates weights and biases through gradient descent to minimize the cost function.

As neural network is composition of functions using chain rule we first obtain the gradient of loss from output layer then moving backwards layer by layer.

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

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \times \frac{\partial z_2}{\partial w_2} \quad (z_2 = w_2 a_1 + b_2)$$

$$\equiv (a_2 - y) \times (a_1)$$

$$\rightarrow \boxed{\frac{\partial L}{\partial w_2} = a_1(a_2 - y)}$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \times \frac{\partial z_2}{\partial b_2} = (a_2 - y) (1)$$

$$\rightarrow \boxed{\frac{\partial L}{\partial b_2} = a_2 - y}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial z_2} = a_2 - y, \quad \frac{\partial z_2}{\partial a_1} = w_2, \quad \frac{\partial z_1}{\partial w_1} = x$$

$$\frac{\partial a_1}{\partial z_1} = \frac{1}{(1+e^{-z_1})} \times e^{-z_1} = a_1(1-a_1)$$

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

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \times \frac{\partial z_1}{\partial b_1}$$

~~$\frac{\partial z_1}{\partial b_1} = -1$~~

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

(c) Using Backpropagation, I would start updating parameters of network from output layer iteratively.

$$w_n = w_n - \eta \frac{\partial L}{\partial w_n}$$

$$B_n = B_n - \eta \frac{\partial L}{\partial B_n}$$

$$w_1 = w_1 - \eta (a_2 - y) w_2 \times a_1 (1 - a_1)$$

$$w_2 = w_2 - \eta a_1 (a_2 - y)$$

$$b_1 = b_1 - \eta (a_2 - y) w_2 a_1 (1 - a_1)$$

$$b_2 = b_2 - \eta (a_2 - y)$$

$\eta \rightarrow$ learning rate : It is the step size for updation of parameters

If it is too large, we may skip minimum Loss function value and if too small, we will require very long time for convergence.

Proper and optimal choice of η should be taken for good learning.

(Q3) (a) ANN processes inputs in independent form as there is no memory of previous inputs, and output is obtained from present input.

RNN has ability to retain memory of previous step. It has hidden state that transfers information from previous to next steps sequentially. Output depends on current input and hidden state.

(b) RNN face problem of vanishing gradient during BPTT (Backpropagation through time). Also the gradient terms are multiplied so it gets shrink slowly so RNN loses information from previous steps.

(c) LSTMs are more concise than RNN as it has 3 gate as forget gate : It loses unnecessary information
 Input gate : decides what to store
 Output gate : decides result which is passed in next layer
 So LSTM are preferred in long term sequences.

(d) They have ~~no~~ cell state.

(e) ANN : Image recognition

RNN : Speech recognition

LSTM : Language Translation

(Q.4) (a) "When mobile dropped on floor, it broke down"

Here, to understand "it", model must remember that even though "floor" is close to "it", but "it" resembles more to the "mobile". And these 2 words have a large gap between them, so model must have long-term memory to remember.

In this case RNN would struggle as RNN has short term memory but they also due to vanishing gradient problem during backpropagation.

The gradients shrinks down gradually and model loses memory.

(b) LSTM has long term memory due to its cell state which carries information ~~back~~ without significant change in it through many steps.

Writing mechanism properly ensures relevant information to be transformed.

Forget \rightarrow discard irrelevant information

Input \rightarrow decides what to be stored

Output \rightarrow decides what to be passed on

In LSTM, due to cell state, gradients do not shrink, the information is carried over many steps.

If forget gate close to 0, the model must forget information of the previous step.

Eg: "Yesterday, I went to hotel. Today, I am going to my home".

Here model should forget information of previous day and focus on new day (Today).