



Assignment 1

1.(a) Linear Regression tells relationship between input features x and output y in linear function. Minimizing squared error function gives stable solution because it is smooth and differentiable and penalizes large errors more than smaller ones.

$$\hat{y} = X\beta \text{ (Linear Regression)}$$

$$\frac{1}{2} (y_i - \hat{y}_i)^2 \text{ (Squared Error Loss)}$$

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

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

$$J(\beta) = \frac{1}{2} (y - X\beta)^T (y - X\beta) \quad \Rightarrow \quad J'_{\beta} = -X^T (y - X\beta) \text{ (w.r.t } \beta)$$

$$X^T (y - X\beta) = 0 \quad (\text{gradient to 0})$$

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

$$\beta = (X^T X)^{-1} X^T y \quad (\text{normal equation})$$

(c) Direct inversions of $X^T X$ in normal equation may be problematic because -

- It needs a lot of memory & computation when dataset is large
- It does not work well for very large dimensional data
- Small calculation errors can become big, so that result is not reliable
- Matrix inversion takes a lot of time



Iterative methods are usually preferred because

- no matrix inversion
- works with small batches
- memory efficient
- scales to large datasets

2. (a) Backpropagation is a method that helps a neural network learn from its mistakes.

Forward pass; Network makes a prediction using the current weights

Backward pass; Checks how wrong the prediction is and sends error backword

Using this, network slightly changes weights so that next prediction is better.

As the output depends on many layers, each layer depends on previous one so the loss is not direct function it is chain of functions.

"Chain Rule" lets us break complex derivatives into small, simple derivatives

$$(b)$$

$$\begin{aligned} z_1 &= w_1 x + b_1 & a_1 &= \sigma(z_1) \\ z_2 &= w_2 a_1 + b_2 & a_2 &= y = \sigma(z_2) \\ L &= -[y \log a_2 + (1-y) \log(1-a_2)] \end{aligned}$$

given



$$\rightarrow \frac{\partial L}{\partial w_2} = \frac{\partial z_2}{\partial w_2} \frac{\partial L}{\partial z_2} \quad (\text{by chain Rule})$$

$$\boxed{\frac{\partial L}{\partial w_2} = a_1 (a_2 - y)} \quad \left(\begin{array}{l} z_2 = w_2 a_1 + b_1 \\ \frac{\partial z_2}{\partial w_2} = a_1 \end{array} \right)$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial z_2}{\partial b_2} \frac{\partial L}{\partial z_2} \quad (\text{chain Rule})$$

$$= 1 \times 1$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial z_2}{\partial b_2} \frac{\partial L}{\partial z_2} \quad (\text{Chain Rule})$$

$$\boxed{\frac{\partial L}{\partial b_2} = 1 (a_2 - y)} \quad \left(\begin{array}{l} z_2 = w_2 a_1 + b_2 \\ \frac{\partial z_2}{\partial b_2} = 1 \end{array} \right)$$

$$\rightarrow \frac{\partial L}{\partial w_1} = \frac{\partial z_2}{\partial w_1} \frac{\partial L}{\partial z_2} \quad (\text{Chain Rule})$$

$$\boxed{\frac{\partial L}{\partial w_1} = w_2 x a_1 (1-a_1) (a_2 - y)} \quad \left(\begin{array}{l} z_2 = w_2 a_1 + b_2 \\ a_1 = \frac{1}{1+e^{-(w_1 x + b_1)}} \end{array} \right)$$

$$z_2 = \frac{w_2}{1+e^{-(w_1 x + b_1)}} + b_2$$

$$\frac{\partial z_2}{\partial w_1} = \frac{w_2 (-x)}{(1+e^{-(w_1 x + b_1)})^2} e^{-(w_1 x + b_1)}$$

$$= w_2 a_1 (1-a_1) x$$

$$\rightarrow \frac{\partial L}{\partial b_1} = \frac{\partial z_2}{\partial b_1} \frac{\partial L}{\partial z_2} \quad (\text{Chain Rule})$$

$$z_2 = w_2 a_1 + b_2$$

$$a_1 = \frac{1}{1+e^{-(w_1 x + b_1)}}$$

$$z_2 = \frac{w_2}{1+e^{-(w_1 x + b_1)}} + b_2$$



$$\frac{\partial z_2}{\partial b_1} = -w_2 e^{-(w_1 a_1 + b_1)} \\ \frac{\partial L}{\partial b_1} = (1 + e^{-(w_1 a_1 + b_1)})^2$$

$$\frac{\partial z_2}{\partial b_1} = w_2 a_1 (1 - a_1)$$

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

$$(c) \quad w_1 \rightarrow w_1 - \eta \frac{\partial L}{\partial w_1}$$

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

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

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

where use these equations to update w_1, w_2, b_1, b_2 in gradient descent.

η = learning rate; it controls how big or small step to take while updating parameters.
 large $\eta \Rightarrow$ fast learning but may lead to unstable learning
 small $\eta \Rightarrow$ slow learning but stable

3 (a) ANN (Artificial Neural Network) processes inputs independently and has no memory of past inputs whereas RNN (Recurrent Neural Network) is designed for sequence data with memory of past inputs.



- (b) Simple RNN struggles with long term dependencies as it suffers from vanishing gradients and cannot remember long term sequence dependencies.
- (c) Role of gates in LSTMs are
- forget gate → it decides what information to remove from memory
 - Input gate → decides what new information to share
 - Output gate → decides what information to pass next step.

Gates help to preserve long term information as memory cell carries information forward with very little change. Gates control updates smoothly instead of overwriting.

- (d) LSTMs fix vanishing gradients as memory cell allows constant error flow and gates regulate information flow. In LSTM,

- (e) Example Tasks:

ANN → House price prediction

RNN → Language Speech Recognition

LSTM → Machine Translation

4.(a) "Manav met Dev after the meeting, and he was exhausted"

The word 'he' is a pronoun, and to understand correctly, we must know 'he' refers to - Manav or Dev.

This creates long term dependency because the meaning of 'he' depends on information that occurred several steps earlier.

Yes, RNN would struggle as in a standard RNN, information about earlier words gradually fades due to vanishing gradient problem.

By the time "he" it may no longer retain clear information about Manav or Dev.

This can lead to incorrect or ambiguous pronoun resolution.

(b) The LSTM stores entity information (Manav or Dev) in its memory cell while reading the sentence. This information is carried forward across multiple words. While processing intermediate words, the forget gate remains close to 1 and prevents loss of entity information.

If the forget gate were close to 0 too early, the entities would have been forgotten.