



# COMPREHENSIVE DEEP LEARNING NOTES

# INDEX

## ANN

Exploring neural networks and backpropagation.

## RNN

Modeling sequential data with RNNs, LSTMs.

## TRANSFORMER

Transforming NLP with attention mechanisms.

## CNN

Revolutionizing image processing and recognition.

## ENCODER DECODER

Sequence-to-sequence learning and translation.

## Deep learning

- ① Deep learning is a subfield of AI and ML that is inspired by the structure of a human brain
- ② DL algorithm uses multiple layers to progressively extract high level feature from the raw input

### ③ DL vs ML

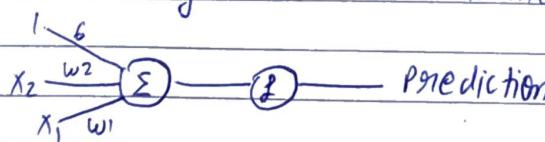
- (i) Data Dependency - DL is more Data Hungry
- (ii) Hardware Dependency - DL use GPU, ML use CPU
- (iii) Training Time - DL training time is high
- (iv) Feature Selection - DL does not require pre build features
- (v) Interpretability - DL models are not interpretable

### ④ Type

- (i) MLP (Multi layer Perceptron)
- (ii) CNN (Convolutional Neural Network)
- (iii) RNN (Recurrent Neural Network)
- (iv) Auto encoders
- (v) GAN (Generative Adversarial Networks)

### ⑤ Perceptron

- (i) It is a algorithm used for supervised machine learning problem
- (ii) we can also say it is a mathematical model / function
- (iii)



$x_1, x_2 = \text{Input}$

$w_1, w_2, b = \text{weight and bias}$

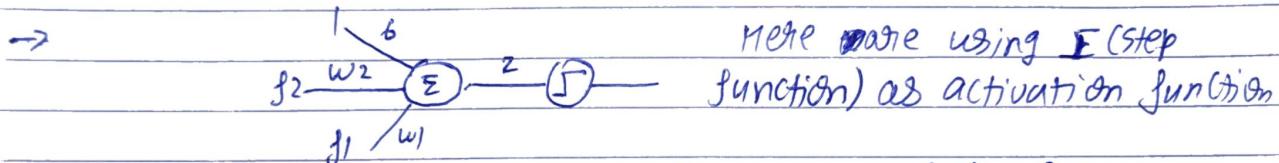
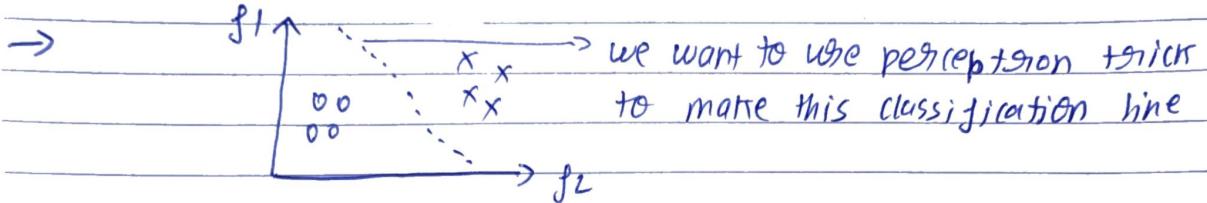
$E = w_1x_1 + w_2x_2 + b$

$f = \text{activation function} = \text{Scale down } w_1x_1 + w_2x_2 + b \text{ in a range of values}$

(2)

## (IV) Perceptron Trick

1. A Algo used to train perceptron
2. Let's suppose we have a data with 2 input features  $f_1, f_2$  and one output feature ( $f_3$ ) = categorical column. This is a Binary classification problem



$$z = f_1 w_1 + f_2 w_2 + b$$

we need to find the value of  $w_1, w_2$ , and  $b$

$$\text{Step function } f = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$$

### 3. Algo

- 3.1 Guesses a random value of  $w_1, w_2, b$
- 3.2 Question Random Point - (that the point correctly positioned or not)
  - 3.2.1 If (Yes) correctly positioned  $\rightarrow$  Do nothing to line, ~~Iteration another Random Point~~
  - 3.2.2 If (No) - make changes in bias and weight so that so that now line is correctly positioned for that point
- 3.3 Do this in a loop for 1000 or 10000 Epoch

Note :- Transformation maths

① Example |  $P_1 \rightarrow$  wrongly classified point  $(4, 5)$

$$\text{Line} = 2x + 3y + 5 = 0$$

Step 1 - Add 1 at the end of the point  $= (4, 5, 1)$

Step 2 - let learning rate  $= 1 \Rightarrow (4 * 1), (5 * 1), (1 * 1)$

Step 3 - if (we want to move the point in (+ve region)) line weight - Points

if (we want to move the point in (-ve region)) line weight + Points

(3)

It can also be written as =  $w_n = w_0 + n(y_i - \hat{y}_i)x_i$

$$\begin{matrix} 2 & 3 & 5 \\ -4 & 5 & 1 \\ -2 & -2 & 4 \end{matrix}$$

here is new line =  $-2x - 2y + 4 = 0$

$w_n = n \text{ new weights of line}$

$w_0 = \text{current weight} + \text{bias of line}$

$x_i = \text{coordinates of that Point } S+1$

cross  
OP  
Ball

## (V) Perceptron Loss Function

### 1. Problem with Perceptron trick

(a) It is possible the model never converge

(b) Perceptron trick does not use loss function, so we are not sure about our model accuracy.

### 2. Loss Function

$$L = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i)) \quad n = \text{number of rows}$$

$$f(x) = w_1 x_1 + w_2 x_2 + b$$

### 3. We use gradient descent

$$L = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i(w_1 x_{i1} + w_2 x_{i2} + b))$$

$$w_1 = w_1 - n \frac{\partial L}{\partial w_1} \quad w_2 = w_2 - n \frac{\partial L}{\partial w_2} \quad b = b - n \frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f(x_i)} \times \frac{\partial f(x_i)}{\partial w_1}, \text{ similarly for } \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial f(x_i)} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases}, \quad \frac{\partial f(x_i)}{\partial w_1} = x_{i1}$$

$$\frac{\partial L}{\partial w_1} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i1} & \text{if } y_i f(x_i) < 0 \end{cases}$$

$$\frac{\partial L}{\partial w_2} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i2} & \text{if } y_i f(x_i) < 0 \end{cases}$$

$$\frac{\partial L}{\partial b} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases}$$

### 3. Perceptron is very flexible

Loss function	Activation function	Output
① Hinge loss	Step function	Perceptron / binary classifier
② log loss (Binary cross entropy)	Sigmoid function	Logistic regression
③ Categorical cross entropy (sparse)	Softmax	Multi class classification Softmax regression
④ Mean square error	Linear	Linear regression

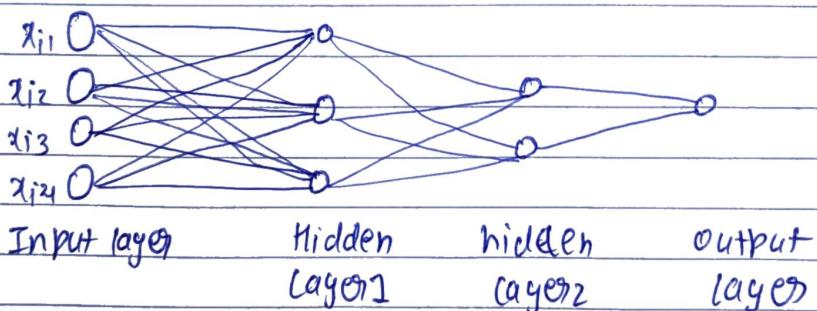
### 4. Why we need multi layer Perceptron and Issue with Perceptron

4.1 Perceptron does not work for Non-linear Data, that's why we use multi layer perceptron

#### 5) Multi layer Perceptron

5.1 Finding the number of Trainable Parameters in a Neural network

(a) Trainable Parameters = weights and bias



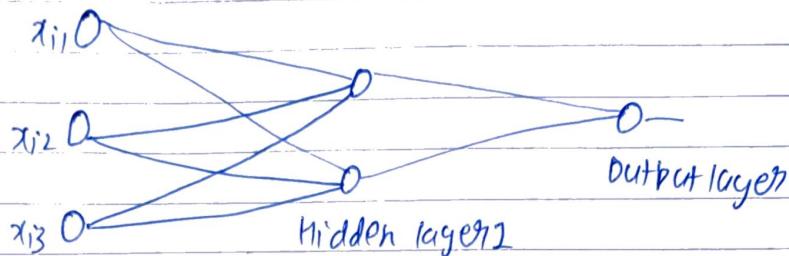
$$\text{Blw Input layer} = \text{number of weights} = 4 \times 3 = 12 \\ \text{and Hidden layer 1} \quad \text{number of bias} = 3 \quad (\text{3 nodes in layer 1})$$

$$\text{Blw Hidden layer 1 and 2} = \text{number of weights} = 3 \times 2 = 6 \\ \text{number of Bias} = 2$$

$$\text{Blw hidden layer 2 and OLP layer} = \text{number of weights} = 2 \times 1 = 2 \\ \text{number of Bias} = 1$$

$$\text{Total Trainable Parameters} = 12 + 3 + 6 + 2 + 1 = 26$$

## 6.2 Notation of weights and Bias and output



Bias =  $b_{ij}$  =  $i$  = layer number  
 $j$  = node number

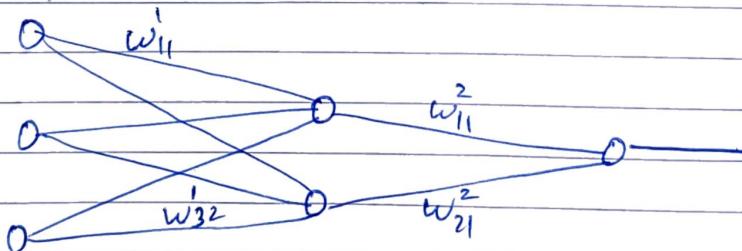
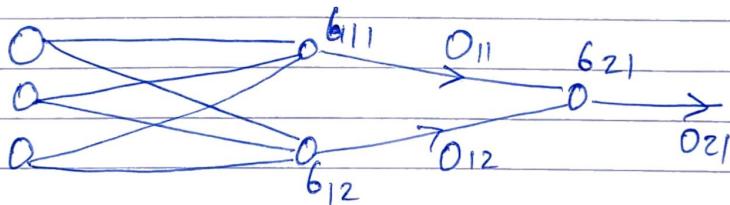
Output =  $O_{ij}$  =  $i$  = layer number

$j$  = node number

Weight =  $w_{ij}^k$  =  $i$  = input layer, node number

$k$  = receiving layer number

$j$  = receiving layer, node number



## 6.3 Geometric Intuition of Perceptron

(a) without weight =  $\frac{P_1}{P_2} \frac{L}{L} + = \# = L \rightarrow$  we are able to capture non linearity

(b) with weight

$P_1 \frac{L}{L} \times w=10$   $P_2 \frac{L}{L} \times w=5 + = \# =$  Here  $P_1$  shape more Dominating.

(6)

## 6.4 Math Intuition

### (a) Pure ReLU

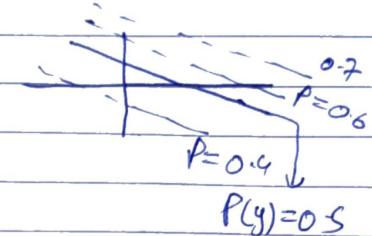
→ Here we are using Sigmoid as Activation Function

→ log loss as loss function

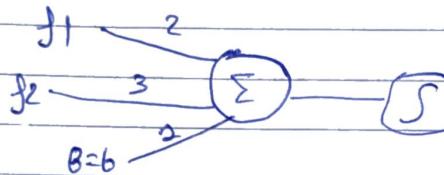
→ Sigmoid → return Probability =  $\frac{1}{1 + e^{-z}}$

Probability at the line = 0.5

$$P(\text{Probability of yes}) = 0.5$$

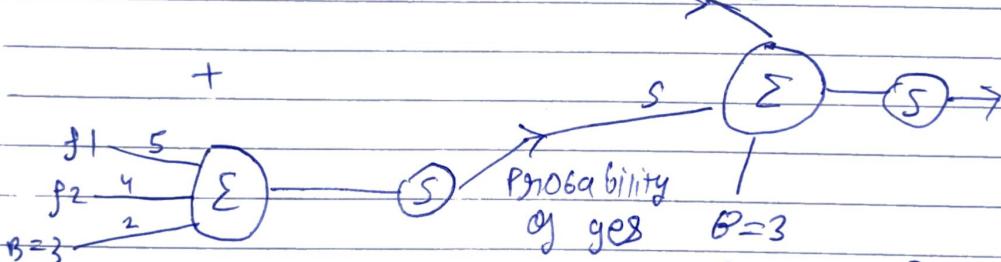


### (b) working



$$\text{Perception 2} = 2x_1 + 3x_2 + 6$$

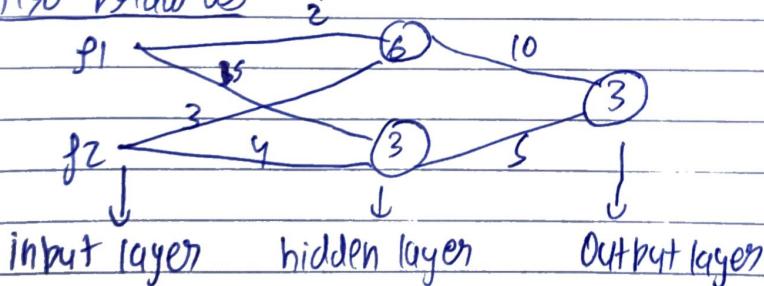
Probability  
of yes



$$\text{Perception 2} = 5x_1 + 4x_2 + 3$$

$$\text{Perception 3} \quad (10x_1 + 5x_2 + 3 = 0)$$

Also Draw as



## 6.5 How we can change architecture of neural network

① Increase number of nodes in Hidden layer

② Increase hidden layers

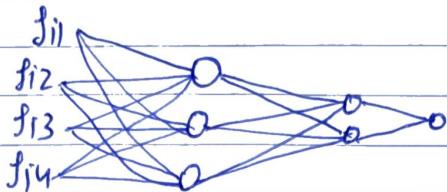
③ Increase input layer nodes

④ Increase output layer nodes (used for multi class classification problem)

## 6.6 Forward Propagation

(9) Let suppose we have a data with 4 input features and 1 output feature

(6) Let suppose we have a neural network like this



$$\text{Total Trainable feature} = 4 \times 3 + 3 + 3 \times 2 + 2 + 2 \times 1 + 1 = 26$$

$$(c) \text{ Prediction} = \sigma(W^T X + B)$$

Here  $\sigma$  = Activation Function

$W$  = weight metrics

$W^T$  = transpose of weight metrics

$X$  = value of features

$B$  = value of Bias

(d) Example - For layer 1

$$G \left( \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix}^T \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \\ b_{14} \end{bmatrix} \right) = \sigma \left( \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{11} \\ w_{21} & w_{22} & w_{23} & w_{21} \\ w_{31} & w_{32} & w_{33} & w_{31} \\ w_{41} & w_{42} & w_{43} & w_{41} \end{bmatrix} \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \\ b_{14} \end{bmatrix} \right)$$

$$\Rightarrow \sigma \begin{pmatrix} w_{11}x_{i1} + w_{12}x_{i2} + w_{13}x_{i3} + w_{14}x_{i4} + b_{11} \\ w_{21}x_{i1} + w_{22}x_{i2} + w_{23}x_{i3} + w_{24}x_{i4} + b_{12} \\ w_{31}x_{i1} + w_{32}x_{i2} + w_{33}x_{i3} + w_{34}x_{i4} + b_{13} \end{pmatrix} = \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix}$$

For layer 2

$$\sigma \left( \begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix}^T \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \\ b_{23} \end{bmatrix} \right) = \sigma \left( \begin{bmatrix} w_{11}^2 O_{11} + w_{12}^2 O_{12} + w_{13}^2 O_{13} + b_{21} \\ w_{21}^2 O_{11} + w_{22}^2 O_{12} + w_{23}^2 O_{13} + b_{22} \\ w_{31}^2 O_{11} + w_{32}^2 O_{12} + w_{33}^2 O_{13} + b_{23} \end{bmatrix} \right) = \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix}$$

For layer 3

$$\sigma \left( \begin{bmatrix} w_{11}^3 & w_{12}^3 \\ w_{21}^3 & w_{22}^3 \end{bmatrix} \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix} + \begin{bmatrix} b_{31} \end{bmatrix} \right) = \sigma \left( \begin{bmatrix} w_{11}^3 O_{21} + w_{12}^3 O_{22} + b_{31} \end{bmatrix} \right) = O_{31} = \hat{y}$$

## 6.7 Loss Function in Deep learning

(a) Regression = MSE, MAE, huber loss  
 (b) ~~category~~

(B) Classification = Binary cross entropy, hinge loss, categorical cross entropy

C) GAN = Discriminative loss, Min Max gan loss

D) Autoencoder = KL divergence

E) Object Detection = Focal loss

F) Embedding = triplet loss

(i) MSE (mean square error)

(a) LOSS Function =  $(y_i - \hat{y}_i)^2$

(b) COST Function =  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

(c) Output layer Activation Function = linear

(d) Advantages → ① easy to understand ② 1 minima ③ Differentiable

(e) Disadvantages → ① Not robust to outliers

(ii) MAE (mean Absolute Error)

(a) LOSS Function =  $|y - \hat{y}|$

(b) COST Function =  $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$

(c) Advantages → ① Robust to outliers ② Error unit is same as output

(d) Disadvantages → Not differentiable

(iii) Huber loss

(a) LOSS Function =  $\begin{cases} \frac{1}{2}(y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta|y - \hat{y}| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}$

$\delta$  = Hyper Parameter

(b) works good if data has too much outliers

(9)

## (iv) Binary Cross Entropy (log loss)

(a) used for binary class classification problem

$$(b) L = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

(c) Activation function for output layer = Sigmoid

## (v) Categorical Cross Entropy

(a) used for multi class classification problem

$$(b) L = - \sum_{j=1}^K y_j \log(\hat{y}_j)$$

K = number of classes in data

$$(c) \text{Cost Function} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij})$$

(d) Activation Function for output layer = Softmax Regressor

## (vi) Sparse Categorical Cross Entropy

(a) used for multi class classification problem

(b) Faster than categorical cross entropy, almost similar to categorical cross entropy

## 6.8 Back Propagation

(a) Algo use to train neural network using gradient descent

Note → Back Propagation algo  
use memorization

(b) Algo :- Epochs = 10

for i in range(Epochs):

    for j in range(X.shape[0]):

        → select 1 row (random)

        → Predict (using forward propagation)

        → calculate loss

        → update weight and bias using GD

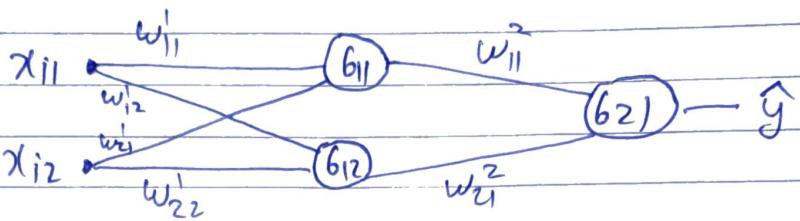
$$w_{\text{new}} = w_{\text{old}} - n \frac{\partial L}{\partial w}, b_{\text{new}} = b_{\text{old}} - n \frac{\partial L}{\partial b}$$

(c) Example

(i) Initialize all weight ( $w=1$ ) and all bias ( $b=0$ )

(ii) Here we are using Activation Function = linear and  
Loss Function = MSE (mean squared error)

(iii) Network



For gradient descent we need to calculate these 9 derivative

$$\frac{\partial L}{\partial w_{11}^2}, \frac{\partial L}{\partial w_{12}^2}, \frac{\partial L}{\partial b_{21}}, \frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial b_{11}}, \frac{\partial L}{\partial w_{21}^2}, \frac{\partial L}{\partial w_{22}^2}, \frac{\partial L}{\partial b_{12}}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_{11}^2} = -2(y_i - \hat{y}_i) \times \frac{\partial}{\partial w_{11}^2} (O_{11}w_{11}^2 + O_{12}w_{12}^2 + b_{11}) = -2(y_i - \hat{y}_i)O_{11}$$

$$\frac{\partial L}{\partial w_{12}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_{12}^2} = -2(y_i - \hat{y}_i) O_{12} \quad \left| \frac{\partial L}{\partial b_{21}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial b_{21}} = -2(y_i - \hat{y}_i) \right.$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_{11}} = -2(y_i - \hat{y}_i) \times w_{11}^2 \times \frac{\partial}{\partial w_{11}^2} (x_{11}w_{11} + x_{12}w_{12} + b_{11}) = -2(y_i - \hat{y}_i)w_{11}^2 x_{11}$$

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{12}} \times \frac{\partial O_{12}}{\partial w_{12}} = -2(y_i - \hat{y}_i) w_{12}^2 x_{12} \quad \left| \frac{\partial L}{\partial b_{11}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial b_{11}} = -2(y_i - \hat{y}_i) w_{11}^2 \right.$$

$$\frac{\partial L}{\partial w_{21}^2} = -2(y_i - \hat{y}_i) w_{21}^2 x_{11}$$

$$\frac{\partial L}{\partial w_{22}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{12}} \times \frac{\partial O_{12}}{\partial w_{22}^2} = -2(y_i - \hat{y}_i) w_{21}^2 x_{12}$$

Note

$$\frac{\partial L}{\partial b_{12}} = -2(y_i - \hat{y}_i) w_{21}^2 \quad \left| \begin{array}{l} \hat{y}_i = O_{11}w_{11}^2 + O_{12}w_{12}^2 + b_{11} \\ O_{11} = x_{11}w_{11}^2 + x_{12}w_{12}^2 + b_{11} \\ O_{12} = x_{11}w_{12}^2 + x_{12}w_{22}^2 + b_{12} \end{array} \right.$$

Note :- A PRO tip about gradient descent]

if  $\frac{\partial L}{\partial w} = (+)$  = mean increase in weight lead to increase in  $L$  (Loss)

if  $\frac{\partial L}{\partial w} = (-)$  = mean decrease in weight lead to decrease in  $L$  (loss)

if  $\frac{\partial L}{\partial w} = \neq$

→ we need to decrease the value of  $w$  so that  $L$  also decrease

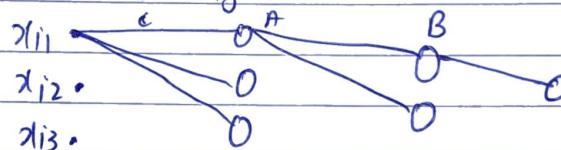
$$w_{\text{new}} = w_{\text{old}} - h \frac{\partial L}{\partial w}$$

if  $\frac{\partial L}{\partial w} = \neq$

→ we need to increase the value of  $w$  so that  $L$  decrease

$$w_{\text{new}} = w_{\text{old}} - h \left( -\frac{\partial L}{\partial w} \right) = w_{\text{old}} + h \frac{\partial L}{\partial w}$$

### (d) Example 2 (Multiple Hidden layers)



NOTE

calculate derivative for A ~~and~~ B and C

$$\hat{y} = O_{21} w_{11}^3 + O_{22} w_{12}^3 + b_{23}$$

$$O_{21} = O_{11} w_{11}^2 + O_{12} w_{12}^2 + O_{13} w_{13}^2 + b_{21}$$

$$O_{22} = O_{11} w_{12}^2 + O_{12} w_{22}^2 + O_{13} w_{23}^2 + b_{22}$$

$$O_{11} = x_{11} w_{11}^1 + x_{12} w_{21}^1 + x_{13} w_{31}^1 + b_{11}$$

$$O_{12} = x_{11} w_{12}^1 + x_{12} w_{22}^1 + x_{13} w_{32}^1 + b_{12}$$

$$O_{13} = x_{11} w_{13}^1 + x_{12} w_{23}^1 + x_{13} w_{33}^1 + b_{13}$$

we want to calculate  $\frac{\partial L}{\partial w_{11}^3}$ ,  $\frac{\partial L}{\partial w_{11}^2}$ ,  $\frac{\partial L}{\partial w_{11}^1}$

In maths if:-

$$x \xrightarrow{f(x)} \rightarrow h(f(x), g(x))$$

And then,

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial f(x)} \frac{\partial f(x)}{\partial x} + \frac{\partial h}{\partial g(x)} \frac{\partial g(x)}{\partial x}$$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_{11}^3}, \quad \frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{21}/O_{22}} \times \frac{\partial O_{21}/O_{22}}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_{11}^1}$$

Because → Both  $O_{21}$  and  $O_{22}$  depend on  $O_{11}$



$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial \hat{y}} \times \left[ \frac{\partial \hat{y}}{\partial O_{j1}} \times \frac{\partial O_{j1}}{\partial O_{i1}} \times \frac{\partial O_{i1}}{\partial w_{ij}} + \frac{\partial \hat{y}}{\partial O_{j2}} \times \frac{\partial O_{j2}}{\partial O_{i1}} \times \frac{\partial O_{i1}}{\partial w_{ij}} \right]$$

## 6.8 Types of Gradient Descent in Neural Network

### (a) Types - Batch Gradient Descent

- Stochastic Gradient Descent
- Mini Batch Stochastic Gradient Descent

### (b) Batch Gradient Descent

(i) For each epoch, we first do prediction on complete dataset, then calculate the loss and then update the weights and bias

(ii) Epoch = Number of time we calculate weights and bias

### (c) Stochastic Gradient Descent

(i) The Gradient Descent we have studied in Back propagation

### (d) Mini Batch Stochastic Gradient Descent

(i) Divide training data into small batch

(ii) For each epoch, shuffle data of each batch, then calculate the loss and update the weight and bias

(e) Speed =  $BGD > MBGD > SGD$

Converge =  $BGD < MBGD < SGD$

#### Note

(i) Why Batch Size is provided in multiple of 2? Ans = To use Ram effectively

(Ram is designed to handle binary value)

(ii) If row = 400, batch size = 150

then batches will be = 150, 150, 100

### (f) Vanishing Gradient Descent Problem

(i) When the gradient of the loss function with respect to the parameters of the earlier layers in the network become exceedingly small. This makes it difficult for the network to learn effectively

(i) causes

- (a) Very Deep Networks
- (b) Certain Activation Functions (Sigmoid, tanh)
- (c) Poor weight initialization

(iii) Is it Exclusive to Backpropagation

NO, It can affect any gradient-based optimization method.

(iv) Symptoms - If loss is not decreasing, it may be because of VGP

- (v) Solutions - (i) Reduce hidden layers
- (ii) Use ReLU Activation Function
- (iii) Do proper weight initialization
- (iv) Use Batch normalization
- (v) Residual network

⑧ Exploding Gradient Problem

(i) when the gradient grow exponentially large the updates to the network weights become excessively large, causing the model to become unstable and the loss function to diverge

- (i) cause - (a) Very deep Networks
- (b) Poor weight initialization
- (c) Long sequence in RNNs

(iii) Symptoms - Loss increase rapidly instead of decreasing

- (iv) Solution - (a) Gradient clipping
- (b) Weight initialization
- (c) Use small learning rates
- (d) Use Batch normalization
- (e) Residual Network

## (9) How to improve the performance of a Neural Network

### 9.1 Problems we face

- (i) Vanishing and exploding gradient descent
- (ii) NOT enough data

Solution - Transfer learning

unsupervised pre training

### (iii) Slow training

Solution - Optimizer

learning rate schedules

### (iv) Overfitting

Symptoms - (i) High Training Accuracy But low validation Accuracy

(ii) High validation loss But low training loss

Cause - (i) Deep Networks (ii) Insufficient training Data  
 (iii) High epoch

Solution - (i) Dropouts (ii) Regularization (iii) Add more data  
 (iv) Reduce layers (v) Early stopping

### 9.2 How to Decide value of these Parameters

(i) Number of Hidden layers - stop adding layer when model start overfitting

(ii) Number of neurons per layer - sufficient/more than what is required

(iii) Batch size → Large Batches - Fast But not good result

→ Small Batches - Slow But good result

(iv) Epochs - very high epoch + Early Stopping

### 9.3 Early Stopping

(i) Early Stopping is used to stop training when performance stops improving

(ii) Issues with wrong Epoch - large Epoch = lead to overfitting

less Epoch = ~~lead to underfitting~~  
 model will not converge

## 9.4 Data Scaling

(i) For unnormalize data  $\rightarrow$  training will be slow, reason = normalization improves the gradient decent convergence.

## 9.5 Dropout layers

- (i) used to reduce overfitting
- (ii) dropping out a fraction of the neurons during training
- (iii) During Prediction, dropout layer does not applied, All the neurons are used for prediction, The weight of each neurons calculated like this =  $W(1-P)$
- (iv) For example - A layer have dropout = 0.2, then weight of the layer =  $W(1-0.2) = 0.75W$

### (v) Tip

- (a) Effect of P - very low value of P = lead to overfitting
  - very high value of P = lead to underfitting
  - it must be b/w 0.2 to 0.5

(b) Start applying dropout from last layer and if you see improvement then move to other layers

(c) For CNN - Range of P = 0.4 to 0.5

For ANN - Range of P = 0.1 to 0.5

For RNN - Range of P = 0.2 to 0.5

## 9.6 Regularization

$$(i) L1 \rightarrow \text{Penalty term} = \lambda \sum_{i=1}^n \|w_i\|_1 \quad \lambda = \text{Hyper parameter}$$

$$(ii) L2 \rightarrow \text{Penalty term} = \frac{\lambda}{2n} \sum_{i=1}^n \|w_i\|^2 \quad n = \text{number of rows}$$

(iii) used to reduce overfitting.

## 9.7 Activation Function

- (i) Need - Activation function introduce non-linearity into the model
- (ii) ideal activation function characteristics
  - (a) ~~the~~ Activation function Should be Non-Linear
  - (b) Activation function Should be differentiable (to allow for gradient-based optimization methods)
  - (c) Activation function Should be non-saturating (saturating activation function effectively cause the gradient vanishing)
  - (d) Activation function Should be zero centred (mean=0, normalized)
  - (e) computational efficiency

### (i) Sigmoid

$$(a) \text{Formula} = \frac{1}{1 + e^{-x}} \quad (b) \text{Shape}$$

(c) output range = 0, 1      (d) saturating function that's why never used in hidden layer

### (iv) tanh

$$(a) \text{Formula} = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (b) \text{Output range} = -1, 1$$

### (v) softmax

$$(a) \text{Formula} = \frac{e^{x_i}}{\sum_j e^{x_j}} \quad (b) \text{Range} = 0, 1$$

### (vi) ReLU (Rectified Linear Unit)

$$(a) \text{Formula} = \max(0, x) \quad (b) \text{Shape} =$$

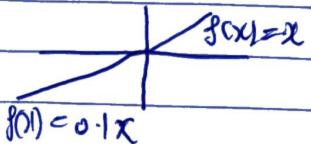
(c) variants  $\rightarrow$  leaky ReLU

linear  $\rightarrow$  parametric ReLU

non linear  $\rightarrow$  Selu  
 $\rightarrow$  ELU

### (vii) leaky ReLU

$$(a) \text{Formula} = f(z) = \begin{cases} 0.01z & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases} \quad (b) \text{Shape}$$



(viii) Parametric ReLU

(a) Formula =  $f(x) = \begin{cases} x & \text{if } x > 0 \\ ax & \text{otherwise} \end{cases}$

(b) shape

$$f(x) = x$$

$$f(x) = ax$$

a = trainable parameter

(ix) Exponential Linear Unit (ELU)

(a) Formula =  $f(x) = \begin{cases} x & \text{if } x \geq 0 \\ a(e^x - 1) & \text{if } x < 0 \end{cases}$

(b) shape

$$f(x) = x$$

(x) Scaled Exponential Linear unit (SELU)

(a) Formula =  $f(x) = \begin{cases} x & \text{if } x \geq 0 \\ a e^x - a & \text{if } x < 0 \end{cases}$

$$a = 1.6732$$

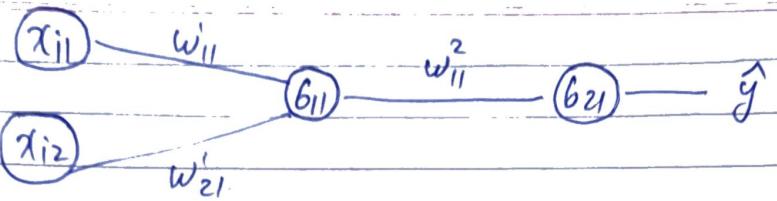
$$b = 1.0507$$

Activation Function	Differentiable	Non-linear	Non-saturating Activation	zero centred	computational Efficiency	Dying Relu problem
Sigmoid	✓	✓	✗	✗	✗	
tanh	✓	✓	✗	✓	✗	
softmax	✓	✓	✗	✗	✗	
gelu	✗	✓	✓	✗	✓	✓
leaky gelu	✓	✓	✓	close	✓	✗
Parametric elu	✓	✓	✓	close	✓	✗
elu	✓	✓	✓	close	✗	✗
selu	✓	✓	✓	✓	✗	✗

9.8 Dying Relu Problem

(i) Relu Neurons that output zero for any input, leading to stop contributing to the learning process

(ii) cause - if  $z < 0$  of a gelu neuron



$$z_{11} = w_{11}' x_{11} + w_{12}' x_{12} + b_{11}$$

$$z_{21} = w_{11}^2 z_{11} + b_{21}$$

$$o_{11} = \text{relu}(z_{11})$$

$$o_{21} = \text{linear}(z_{21}) = \hat{y}$$

Let assume  $z_{11} < 0$ , that make  $o_{11} = 0$

Gradient decent

$$w_{11}' = w_{11}' - \eta \frac{\partial L}{\partial w_{11}'}$$

$$w_{21}' = w_{21}' - \eta \frac{\partial L}{\partial w_{21}'}$$

$$\frac{\partial L}{\partial w_{11}'} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_{21}} \frac{\partial z_{21}}{\partial o_{11}} \frac{\partial o_{11}}{\partial z_{11}} \frac{\partial z_{11}}{\partial w_{11}'} = 0 \quad \left\{ \begin{array}{l} o_{11} = 0 \rightarrow \frac{\partial o_{11}}{\partial z_{11}} = 0 \\ \frac{\partial z_{11}}{\partial w_{11}'} = 0 \end{array} \right.$$

$$\text{Similarly } \frac{\partial L}{\partial w_{21}'} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_{21}} \frac{\partial z_{21}}{\partial o_{11}} \frac{\partial o_{11}}{\partial z_{11}} \frac{\partial z_{11}}{\partial w_{21}'} = 0 \quad \left\{ \begin{array}{l} o_{11} = 0 \rightarrow \frac{\partial o_{11}}{\partial z_{11}} = 0 \\ \frac{\partial z_{11}}{\partial w_{21}'} = 0 \end{array} \right.$$

$$w_{11}' = w_{11}' - 0 = \{ \text{this lead to no learning} \}$$

$$w_{21}' = w_{21}' - 0$$

poor

- (iii) causes
  - (a) ~~positive~~ weight initialization
  - (b) high negative Bias
  - (c) large learning ~~rate~~ Rate

(iv) Once a relu neuron  $z$  become  $< 0$ , it is considered dead and not recoverable

- (v) Solution — (a) Set low learning rate  
 (b) use  $g_{ReLU}$  variant  
 (c) set bias 0.01

## 9.9 Weight initialization technique

(i) choosing wrong weight initialization technique lead to these problem

- (a) vanishing gradient problem (b) exploding gradient problem
- (c) slow convergence

(ii) wrong weight initialization technique

- (a) zero initialization (b) Non-zero constant initialization
- (c) Random initialization with small weight
- (d) Random initialization with large weight

(iii) what to do : - Random initialization But with correct variance

(iv) Techniques (a) Xavier / Glorot weight initialization (Default in Keras)  
 (b) He weight initialization

TanH and Sigmoid activation function  $\rightarrow$  Xavier / Glorot

Tanh activation function  $\rightarrow$  He weight initialization work better.

## 9.10 Batch Normalization

(i) Batch normalization consists of normalizing activation vector from hidden layer using mean and covariance of the current batch,

(ii) This normalization step is applied right before or right after the non-linear Function.

(iii) Adv - (a) Speed up training (b) Regularizes the model  
 (c) Decrease the importance of initial weights.

## 9.11 Optimizer

(i) challenges in using gradient decent as optimizer

(a) Deciding the value of learning rate

Slow value  $\rightarrow$  lead to slow convergence

large value  $\rightarrow$  lead to miss the minima point

(b) Same learning rate for all the dimensions

(c) Converge at local minima

(d) Saddle Point

It is a point where slope of 1 direction is going up and slope of another direction is going down wards, this make a big region where slope remain same

which mean  $\frac{\partial L}{\partial w} = 0$ ,  $w_n = w_0 - n \frac{\partial L}{\partial w} \Rightarrow [w_n = w_0]$

(i) Gradient decent alternative

(a) momentum (b) adagrad (c) NAG (d) RMS prop (e) Adam

(ii) EWMA

(a) Exponentially weighted moving average or exponential weighted average

(b) a technique used to find trend in time series based data.

### (C) Formula

$$V_t = \beta V_{t-1} + (1-\beta)x_t$$

$\beta$  = Hyper Parameter ( $0-1$ )

$x_t$  = Observation / data

(d) EWMA assign exponentially decreasing weight to older observation and give more importance to recent observation

$$V_t = \beta V_{t-1} + (1-\beta)x_t$$

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta)x_1 = (1-\beta)x_1$$

$$V_2 = \beta V_1 + (1-\beta)x_2 = \beta(1-\beta)x_1 + (1-\beta)x_2$$

$$V_3 = \beta^2(1-\beta)x_1 + \beta(1-\beta)x_2 + (1-\beta)x_3$$

here for  $V_3 \rightarrow x_1 \times \beta^2, x_2 \times \beta, x_3 \times \beta$

### (e) Impact of $\beta$

High  $\beta$  (close to 1): The EWMA place more weight on past values and less weight on the most recent observation. The EWMA react slowly to recent changes in the data.

Low  $\beta$  (close to 0): The EWMA place more weight on the most recent observation and less weight on past values. The EWMA react quickly to recent changes in the data.

### (ii) Convex and non-convex optimization

$$\text{Convex optimization} = \boxed{\cup}, \text{Non-convex opti.} = \boxed{\cap}$$

Issue with non-convex graph = local minima, saddle point, high curvature.

## (v) Momentum

(a) Problem Solve - local minima, Saddle Point, high curvature.

(b) The main idea momentum optimization is to acceleration the gradient vector in the right direction, leading to faster converging towards the minima and reducing oscillations

## (c) Formula

$$w_{t+1} = w_t - v_t \quad \left\{ \begin{array}{l} v_t = \beta * v_{t-1} + \eta \nabla w_t \\ \text{momentum term} \quad \text{gradient of } w_t \end{array} \right\}$$

Hyperparameter -  $\eta$ ,  $\beta$

if  $\beta = 0$ , then momentum = gradient decent

## (d) Disadvantages

Momentum, Because of momentum the gradient decent make oscillation on global minima, wasting some time.

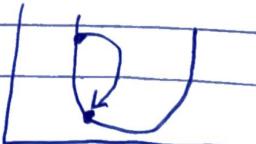
## (vi) Nesterov Accelerated Gradient

(a) mostly perform better then momentum. Builds on the concept

(b) The main idea behind NAG is to anticipate the future position of the parameters and compute the gradient not at the current position but at this anticipated future Position

## (c) Formula

→ First we apply momentum term and make the jump to new point (look ahead term)



→ Now we calculate the gradient at the new point

→ Now on the basis of this new gradient we decide where to move.

$$w_{ea} = w_t - \beta v_{t-1} \quad (\text{step 1})$$

$\downarrow$   
look ahead

$$v_t = \beta v_{t-1} + \eta \nabla w_{ea} \quad (\text{calculating slope at new point})$$

$$w_{t+1} = w_t - v_t \quad (\text{update rule})$$

(d) Disadvantages : convergence at local minima

### (vii) Adaptive Gradient Algorithm

(a) Adagrad adapts the learning rate for each parameter individually making larger updates for infrequent parameters and smaller updates for frequent parameters.

(b) works better for :- if input feature values have different range  
if input feature values are sparse

### (c) Formula

$$w_{t+1} = w_t - \frac{\eta \nabla w_t}{\sqrt{v_t + \epsilon}} \quad \left( \begin{array}{l} \epsilon \rightarrow \text{non-zero very small number} \end{array} \right)$$

$$v_t = v_{t-1} + (\nabla w_t)^2$$

(d) Div - learning rate Decay (that why not used in neural network)

### (viii) Root mean Squared Propagation (RMSprop)

(a) RMSprop modifies Adagrad by introducing a decaying avg of past squared gradients which help stabilize the learning rate over time.

(6) It is one of the optimizing technique for neural networks

(c) Formula

$$w_{t+1} = w_t - \frac{\eta \nabla w_t}{\sqrt{v_t + \epsilon}} \quad \epsilon = \text{very small non zero number.}$$

$$v_t = \beta v_{t-1} + (1-\beta) (\nabla w_t)^2$$

(ix) Adam

(a) Adam (Adaptive Moment Estimation) is an optimization algorithm that combines the advantages of 2 other popular optimization method: Adagrad, RMSprop

(b) Best Optimizing technique.

(c) Formula

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) \nabla w_t$$

$$v_t = \beta_2 v_{t-1} + (1-\beta_2) (\nabla w_t)^2$$

Bias correction

$$\hat{m}_t = (m_t) / (1 - \beta_1^t) \quad \begin{matrix} \rightarrow \text{this } t \text{ mean} \\ \text{epochs number.} \end{matrix}$$

$$\hat{v}_t = (v_t) / (1 - \beta_2^t)$$

$\beta_1$  and  $\beta_2$  = trainable Parameter ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.99$  generally)

$$w_{t+1} = w_t - \frac{\eta \hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$$

(10) convolutional Neural Networks

10.1 Special kind of neural network for processing data that has a grid like topology like time series data (1D) or image (2D)

## 10.2 CNN Architecture

### (a) Input layer

### (b) convolutional layer

- contain multiple [filter (kernels) + relu]
- Output = feature map
- we can also do batch normalization after convolution layer before the pooling layer

### (c) Pooling layer

- Output = tensor

### (d) Flatten layer

- Convert high dimensional data into 1D

### (e) Fully connected layer

- like ANN

## 10.3 Why we not use ANN on image data

- overfitting
- High computational cost
- loss of arrangement of pixels

## 10.4 Basics of images

### (a) Types of images

→ Grey Scaled images  
→ Colored images

### (b) Grey Scaled images

- Store image like 2d array in 1 channel
- value varies from 0 (black) to 255 (white)

### (C) Colored image

- ↳ Store image like 2d array in 3 channel
- ↳ Shape  $628, 228, 3$

### 10.5 Convolutional layer

(a) Core Building Block of a CNN

(b) It applies a convolution operation to the input, passing the result to the next layer.

### (c) Intuition

1 Convolutional layer = Detect primitive features

2 Convolutional layer = Detect slightly more complex feature

3 Convolutional layer = Detect advance features

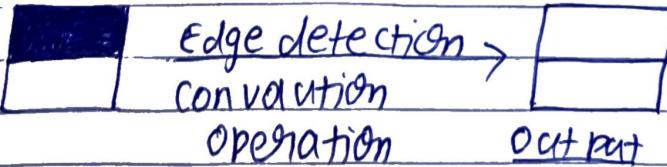
### (d) Steps

- Initialize the number filter and size of filter (generally  $3 \times 3$ )
- Slide the filter over the input image
- Perform element wise multiplication and summation

Note - The value of filter in CNN are not manually decided. Instead they are learned automatically during the training process through backpropagation. Initially filters are initialized with small random value.

### (e) Example

Input image =



- Image =  $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 255 & 255 & 255 & 255 & 255 & 255 \\ 255 & 255 & 255 & 255 & 255 & 255 \\ 255 & 255 & 255 & 255 & 255 & 255 \end{bmatrix}$   $\times$   $\begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$  → Kernel filter

Size (6x6)      convolution operation

$$= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} * \begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} = 0 \times 1 + 0 \times -1 + 0 \times -1 + 0 \times 1 + 0 \times 1 = 0$$

$$\Rightarrow \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} * \begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} = 0$$

→ repeat this = feature map  $\begin{bmatrix} 0 & 0 & 0 & 0 \\ 255 & 255 & 255 & 255 \\ 255 & 255 & 255 & 255 \\ 0 & 0 & 0 & 0 \end{bmatrix}$

$= \boxed{\quad}$  ↙

(8) Formula If image size ( $n \times n$ ), if filter size ( $m \times m$ )

$\Rightarrow$  Feature map size =  $(n - m + 1), (n - m + 1)$

like image size = (6,6), filter size (3x3)

feature map size =  $(6 - 3 + 1), (6 - 3 + 1) = (4, 4)$

## 10.6 Padding

### (a) Need

→ Loss of data because of convolution operation

→ The border pixels will be part of less convolution as compare to centre pixels

- (b) Padding refers to adding extra pixels around the borders of the input image
- (c) In Keras are 2 types of Padding  $\rightarrow$  valid (no padding)  
 $\rightarrow$  same (Applying padding)

#### (d) Intuition

$\rightarrow$  Add zero row and column at the borders of image  
 $\rightarrow$  image size =  $4 \times 4$

$$\begin{array}{ccccccccc} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ & | & & & & & & & \\ & 0 & 1 & 2 & 3 & 4 & 0 & 0 & 0 \\ & | & & & & & & & \\ & 0 & 5 & 6 & 7 & 8 & 0 & 0 & 0 \\ & | & & & & & & & \\ & 0 & 9 & 10 & 11 & 12 & 0 & 0 & 0 \\ & | & & & & & & & \\ & 0 & 13 & 14 & 15 & 16 & 0 & 0 & 0 \\ & | & & & & & & & \\ & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \rightarrow \text{Result}$$

$\text{image}$

(e) Formula if image size  $(n \times m)$ , if filter size  $(m \times m)$   
 $p = \text{padding}$

$$\text{Feature map size} = (n + 2p - m + 1), (n + 2p - m + 1)$$

#### 10.7 Strides

(a) Strides determine the step size by which the filter moves across the input image. The stride affects the spatial dimensions of the output feature map.

(b) why to use  $\rightarrow$  To reduce computing power  
 $\hookrightarrow$  To capture high level features

(c) Not so much important now days

(d) Formula image size  $(n \times n)$ , filter size  $(m \times m)$ ,  $p = \text{padding}$   
 $s = \text{Strides}$

$$\text{Feature map size} = \left( \frac{n + 2p - m}{s} + 1, \frac{n + 2p - m}{s} + 1 \right)$$

### (e) Effect of strides

$\text{Stride}(1,1)$  = The output feature map is larger and closer in size to the input

$\text{Stride}(3,3)$  = The output feature map is smaller, reducing the computational load, information loss.

## 10.8 Pooling layer

### (a) Problem with convolution operation

- Memory issue = convolution operation consume a lot of memory
- Translation variance = the convolution operation will become location dependent.

(b) Types - max pooling, min pooling, avg pooling, L2 pooling, global pooling

(c) Parameters - size (2,2), stride, type

(d) Adv - memory reduced, translation invariance, max pooling → enhance feature

(e) div - loss of information

- image segmentation required translation variance.

(f) Implementation - feature map

$$\begin{matrix} & \boxed{3} & \boxed{1} \\ & \boxed{2} & \boxed{5} \end{matrix} \quad \begin{matrix} & \boxed{1} & \boxed{3} \\ & \boxed{0} & \boxed{2} \end{matrix} \rightarrow \begin{matrix} \boxed{5} & \boxed{3} \\ \boxed{7} & \boxed{4} \end{matrix}$$

Pooling = size =  $2 \times 2$

stride = 0

type = max pooling

## 10.9 Data Augmentation

- (a) Data augmentation is a technique used to increase the diversity and quantity of training data without actually collecting new data.
- (b) why it is used  
→ increase data size  
→ Reduce overfitting  
→ Enhance generalization
- (c) we do not apply data augmentation on test data set.
- (d) common data augmentation techniques
- Geometric Transformations  
Rotation, translation, scaling, Flipping
  - Color Transformations  
Brightness adjustment, contrast adjustment, saturation adjustment
  - Occlusion  
Random cropping, cutout (randomly using masking out square regions of the input during training)
  - Elastic Distortions  
Applying random elastic distortions to the image.
  - Noise injection  
Introducing random noise to the image.

## 10.10 Pretrained Models

(a) Pretrained models are neural network model that have been previously trained on a large dataset and can be reused by new, related tasks.

- (b) Benefits
- Reduced Training Time
  - Improved performance
  - Resource efficiency

(c) Dis - If pretrained model does not trained on the our current target task (Solved by transfer learning)

## (d) Commonly used Pretrained Models

- Image Model :- ResNet, VGG16, VGG19, Inception, MobileNet EfficientNet
- NLP Model :- GPT, BERT, TS, ROBERTa, Transformer.

## 10.11 Transfer learning

(a) A technique where a pretrained model is used as a starting point for a new task, leveraging the knowledge gained from the original task.

- (b) Adv
- Reduced Training time
  - Improved Performance
  - Resource Efficiency

## (c) Type

- Feature extraction
- Fine - Tuning

### (d) Feature Extraction

- Use the pretrained model as a fixed feature extractor
- In this approach we freeze the weights of the pretrained model and only train the final layers specific to the new task
- In case of CNN we pre-train fully connected layer and freeze the convolutional layer
- used when your model are already trained on kind a similar problem - using dog vs cat classifier on lion vs tiger classifier.

### (e) Fine Tuning

- Unfreeze some or all of the layers of the pretrained model and jointly train them with the new layers on the new task. This allows the model to adapt the pretrained weight to the specifics of the new task
- In case of CNN we pre-train fully connected layer and some last convolution layers also
- used when your model are not already trained on kind a similar problems

## II Functional API

### (a) Assumptions of sequential models

- i) 1 input
- ii) 1 output
- iii) All layers are linearly attached

(b) The Function API allow us to build more complex neural network architecture as compared to Sequential API

(c) The Function API allow us to build model with multiple input, multiple output, shared layers or non-linear topology.

## 12 RNN

### 12.1 RNN Basics

(a) Designed to process sequential data

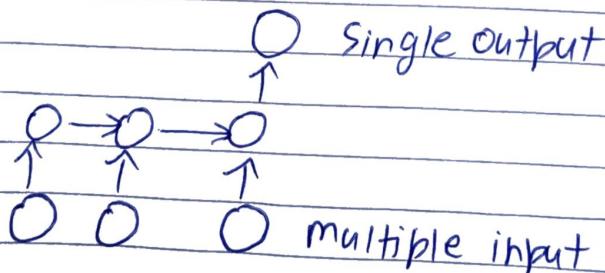
(b) RNN have connections that form directed cycles, enabling them to maintain a memory of previous inputs in the sequence

(c) Good For - Time Series analysis, language modeling.

### (d) Types

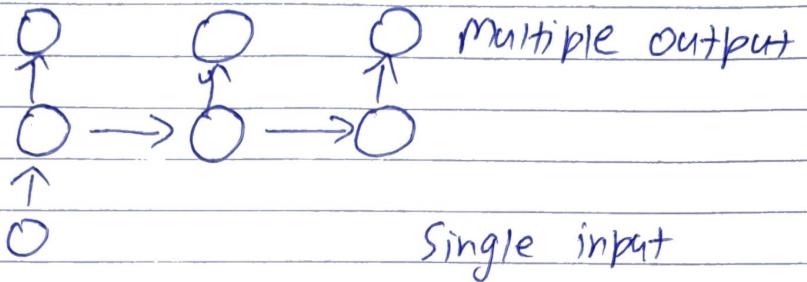
- (i) Many to one RNN
- (ii) One to many RNN
- (iii) Many to many RNN
- (iv) One to One RNN

### (i) Many to one RNN



Application - Sentiment analysis, Rating prediction

### (ii) One to Many RNN



Application - Image captioning

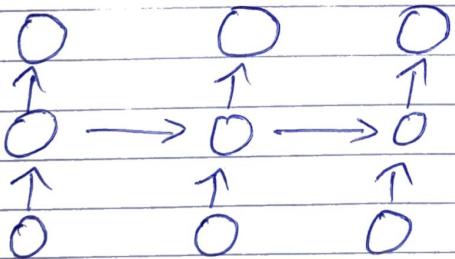
### (iii) Many to Many RNN

1. Also known as Seq2Seq models

2. Types

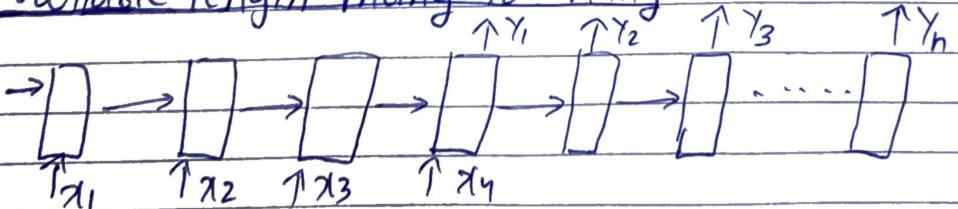
- Same length many to many
- Variable length many to many

3. Same length many to many



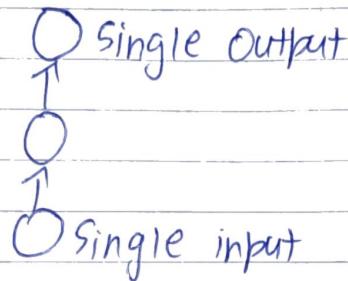
Application - Part of speech tagging  
- Name Entity recognition

4. Variable length many to many



Application - Machine translation

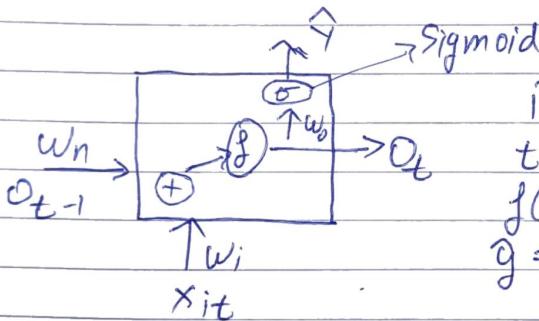
### (iv) Single to Single RNN / One-to-One RNN



Application - image classification

\* Technically this is not RNN

### (e) RNN architecture

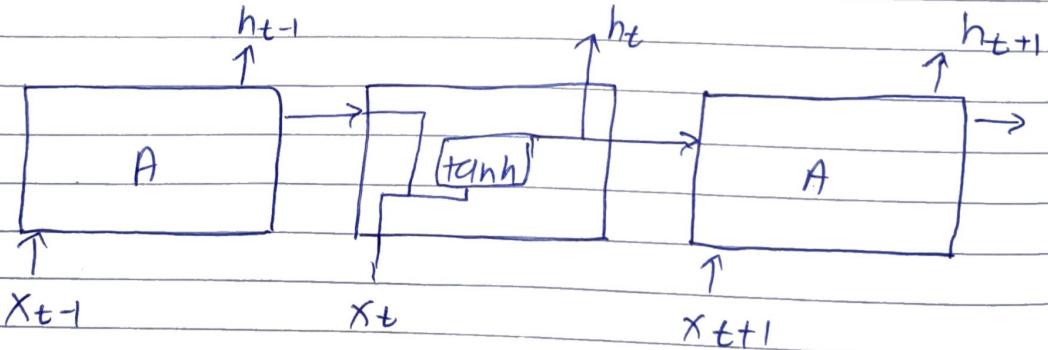


$i = \text{row number}$

$t = \text{time step}$

$$f(w_i x_{it} + O_{t-1} w_n) = O_t$$

$$g = \text{Sigmoid}(O_t w_0)$$



### (f) variation of RNN

- (i) Long Short Term Memory (LSTM)
- (ii) Gated Recurrent Unit (GRU)
- (iii) Bidirectional RNN (Bi-RNN)
- (iv) Echo State Networks (ESN)

(g) why ANN is not good for text data

- (i) The order of words in a sentence is important for meaning. ANN do not inherently account for the order of input features, treating each word as independent.
- (ii) ANN lack the ability to maintain and utilize context over long sequences.
- (iii) High dimensionality of text data.
- (iv) Variable length of text sequences.

(h) Difference between RNN and ANN

RNN

- (i) RNN process input sequence based on time steps.
- (ii) RNN are trained using Back propagation through time (BPTT)
- (iii) RNN maintain a hidden state that captures info about the sequence seen so far, which makes it capable of capturing dependencies over time.

ANN

- (i) ANN process input simultaneously in a feed forward manner.
- (ii) ANN are trained using standard backpropagation algorithm.
- (iii) Each input is processed independently without considering previous input in the same training instance.



## (i) Bidirectional RNN

(a) Capture information from both past and future states in a sequence

- (b) Application - Name entity recognition
- Part of speech tagging
  - Sentiment analysis
  - Machine translation

## (j) Problem with simple RNN

- (i) Vanishing gradient problem  
 (ii) Exploding gradient problem

## (iii) Different in learning Long-Term-Dependencies

(a) Simple RNN struggle to capture dependencies that span long sequences due to the vanishing gradient

### (b) Solution

- Use ReLU activation functions
- Better weight initialization
- Use LSTM / GRU

## (iv) Unstable training because of exploding gradient decent.

### 12.2 LSTM

(a) A type of RNN designed to better capture long-term dependencies in sequential data.

(b) They are particularly effective at addressing the vanishing gradient problem that simple RNN face, enabling them to learn from and remember important information over long sequences.

### (c) Components of LSTM

#### (i) Input

- Previous hidden state
- Previous cell state
- Input for the current time

#### (ii) Output

- Current cell state
- Current hidden state
- 

#### (iii) Processing

- calculate hidden state
- Add or remove information from cell state

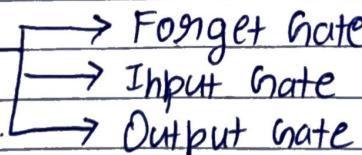
#### (iv) Cell State ( $c_t$ )

- It carrying relevant information through the entire sequence.
- Help in maintaining long-term dependencies.

#### (v) Hidden State ( $h_t$ )

- Output of the LSTM at each time step.

#### (vi) Gates



### (vi) Forget Gate

- Decides what to forget from the cell state based on current input and short term context

### (vii) Input Gate

- Decides what new information to add to cell state

### (viii) Output Gate

- Produces the new hidden state based on the updated cell state

## (d) Why LSTM are effective

(i) Mitigating vanishing gradients - The state allows gradients to follow unchanged over long sequences.

(ii) Selective Memory - The gates (input, output, forget) enable selective memory, allowing the network to retain or discard information as needed.

## 12.3 Gated Recurrent Unit (GRU)

(a) GRU is a type of RNN architecture designed to impact the performance of traditional RNN

### (b) Architecture

#### (i) Reset Gate

- Decides how much of the past info should be forgotten

#### (ii) Update Gate

- Decide how much of the previous hidden state should be passed to the current hidden state

### (iii) Input

- Previous hidden state
- Input for the current time

### (iv) State

- Combines the previous hidden state and the candidate activation based on the update gate.

### (c) Advantages of GRU

- Simple architecture as compared to LSTM
- Less computationally intensive as compared to LSTM
- Efficient memory use

### (d) GRU vs LSTM

Feature	GRU	LSTM
Gates	2 (Update, Reset gates)	3 (Input, Output, Forget gate)
Memory cell	No memory cell	have separate memory cell
Hidden state	Combine Hidden and cell state	Separate hidden and cell state
Update process	Fewer parameters (Simple)	Complex
Parameters	Fewer parameters	More Parameters
Time	Less training time	More Training Time
Performance	Comparable to LSTM	R robust
Complexity	Simple architecture	Complex architecture
Computation	Less intensive	more intensive

## 13 Encoder Decoder

### 13.1 Encoder

- (i) The encoder processes the input sequence and compresses the information into a context vector
- (ii) This typically involves several layers of LSTM/GRU

### 13.2 Decoder

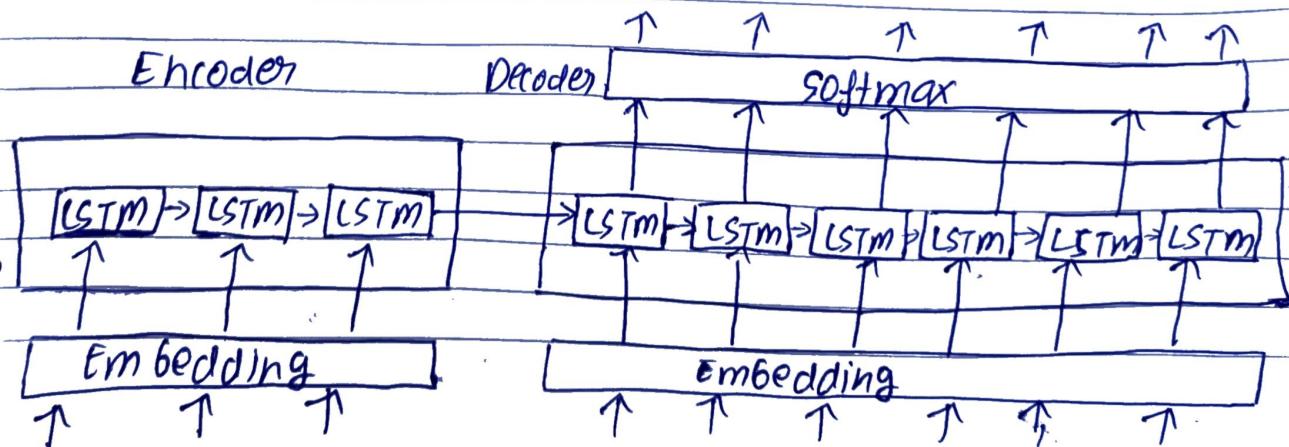
- (i) The decoder takes the context vector from the encoder and generates the output sequence
- (ii) Similar to encoder, it contains several layers of LSTM/GRU

### 13.3 Context vector

The context vector is a fixed size vector that encapsulates the information from the entire input sequence

### 13.4 Architecture

Text  $\rightarrow$  [Encoder]  $\xrightarrow{\text{context vector}}$  [Decoder]  $\rightarrow$  Output



## 13.5 challenges of vanilla Encoder-Decoder

- (i) Fixed length context vector lead to loss of information
- (ii) Gradient Vanishing / Exploding Problem
- (iii) Inability to capture long Range Dependency

## 13.6 Attention mechanism

- (i) It allows the decoder to focus on different parts of the input sequence at each step of the output generation
- (ii) The decoder get a weighted sum of all encoder hidden states where the weights are dynamically computed based on the decoder's current state.
- (iii) Encoder . Attention mechanism, decoder
- (iv) we can't use trans for learning for Encoder-Decoder with attention mechanism

Because Encoder-Decoder with attention mechanism train sequentially (which lead to slow training) that's lead to high cost and time for training. Hence we can't train it on very huge dataset.

## 14. Transformers

14.1 Transformers are a type of neural network architecture designed to handle sequential data

### 14.2 Architecture

(i) Encoder - The encoder is responsible for processing the input sequence and generating a set of encoder representations

(ii) Decoder - The decoder generates the output sequence, using the encoder representations from the encoder

(iii) Output layer - The final linear layer transforms the decoder's output into logits which are converted into probability using a softmax function to produce the final output tokens.

### 14.3 Advantages

(i) Scalability (ii) Transfer learning (iii) multi mode input output  
 (iv) flexible architecture (v) Transfer learning

### 14.4 Disadvantages

(i) High computational cost (ii) Large amount of data requires  
 (iii) Prone to overfitting (iv) High energy consumption

### 14.5 Self Attention mechanism

(i) mechanism that take static embedding and generate contextual embedding.

## (i) Issue with static embedding

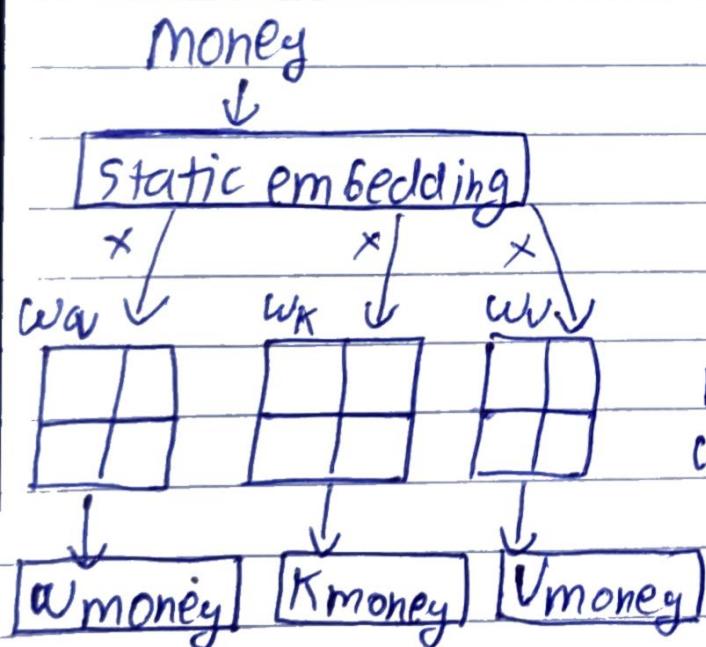
A word can have more than 1 meaning, But static embedding capture only one meaning.

## (ii) Working

Creating contextual embedding For line →

"Money bank grows"

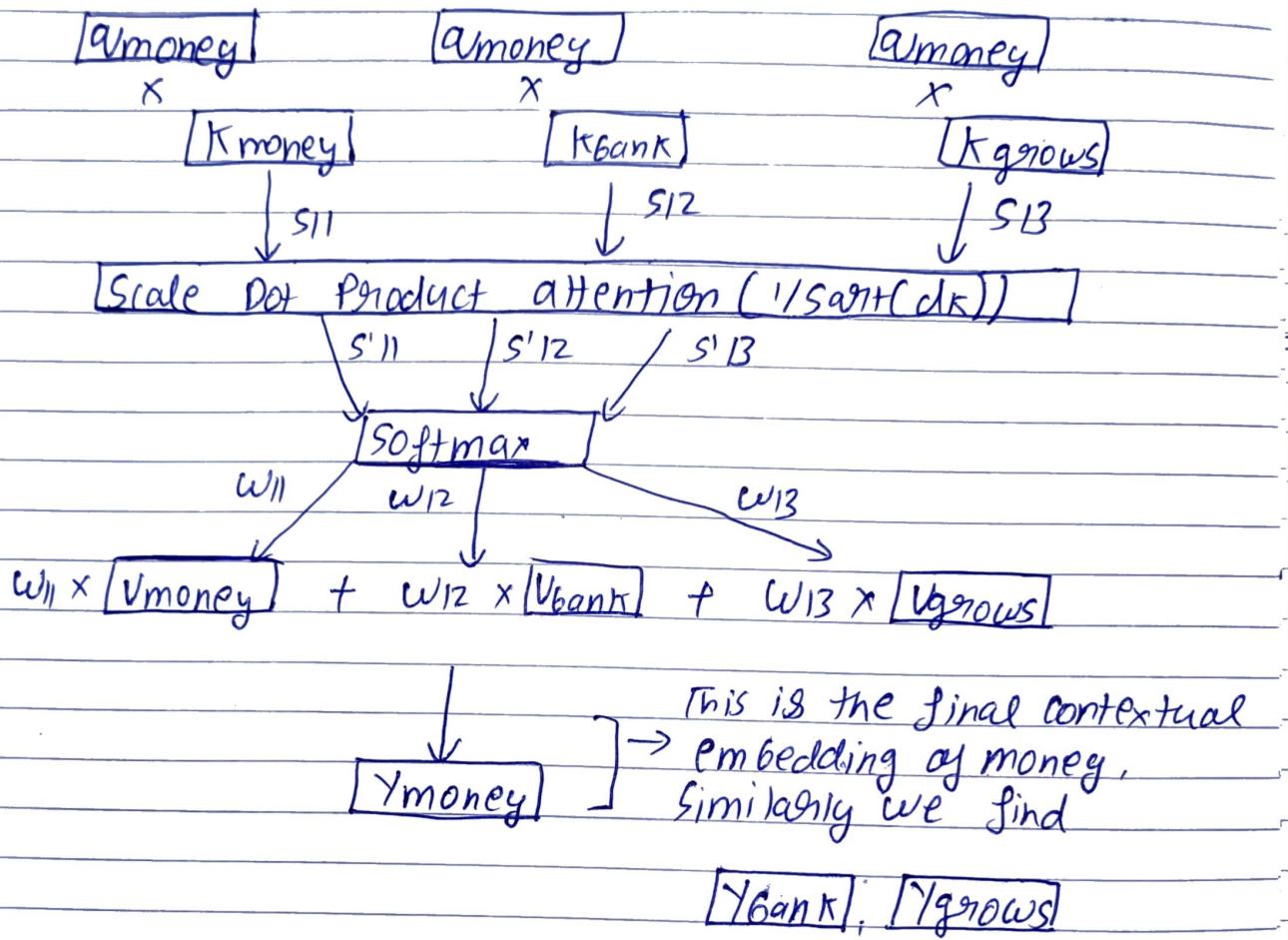
Step 1:- Making 3 embedding for each word (query, key, value)



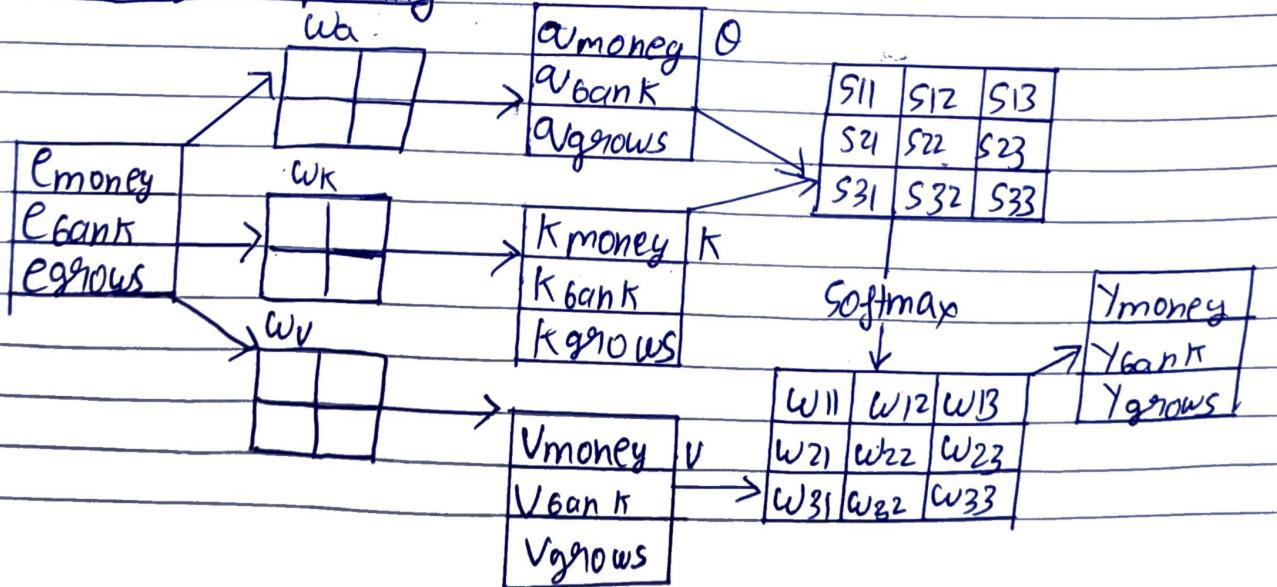
These ( $a$ ,  $k$ ,  $v$  matrices) are obtained by multiplying the word static embedding with 3 learned weighted matrix, initially they were initialized with random number and obtain correct value through training

Similarly we find  $a_{\text{Bank}}$   $k_{\text{Bank}}$   $v_{\text{Bank}}$ ,  $a_{\text{grows}}$

## Step 2: Finding Contextual Embedding



How we do this partially



### (iv) Scaled Dot Product Attention

To stabilize the gradient during training, the dot product are scaled by the square root of the dimension of the key vector ( $d_k$ )

$$\text{Scaled Score}(x_i, x_j) = \frac{\phi_i \cdot k_j^T}{\sqrt{d_k}}$$

$$(v) \text{Attention}(\phi, k, v) = \text{softmax}\left(\frac{\phi \cdot k^T}{\sqrt{d_k}}\right) \cdot v$$

(vi) Self Attention is called "self" because it operates on a single sequence of elements and relates each element at all other elements in the same sequence.

(vii) Problem with the self attention mechanism without multi head attention

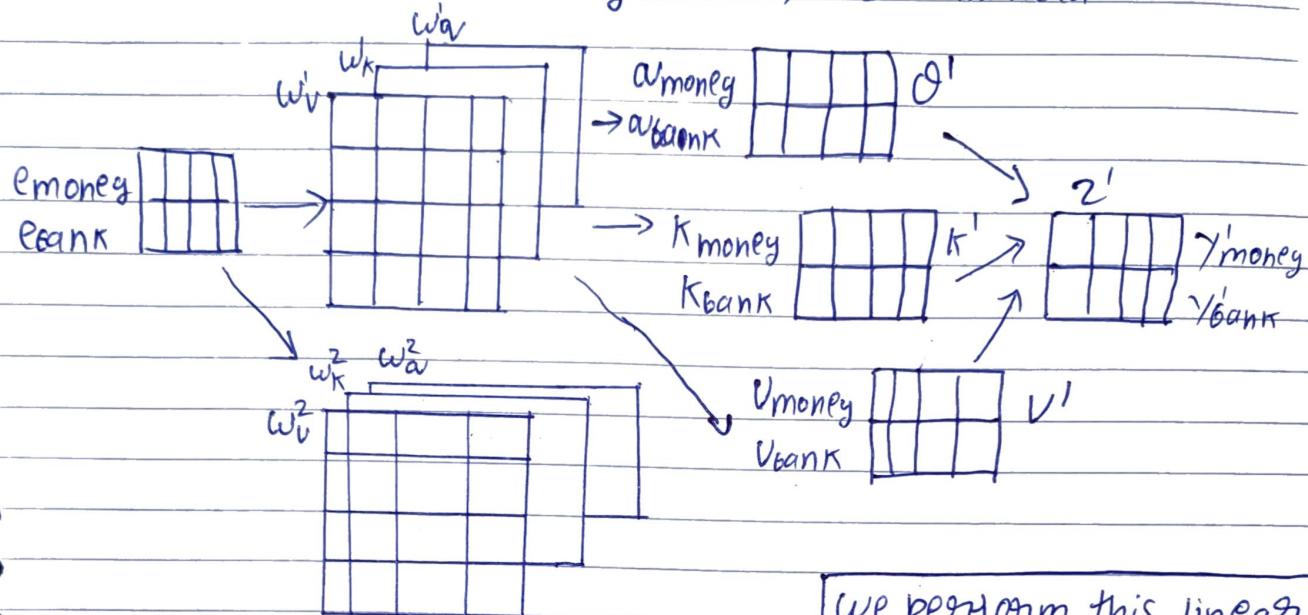
If a text has more than one perspective it will only gone capture the most prominent connection leading to miss out other important but obvious relation (single perspective)

### 14.6 Multi Head Attention

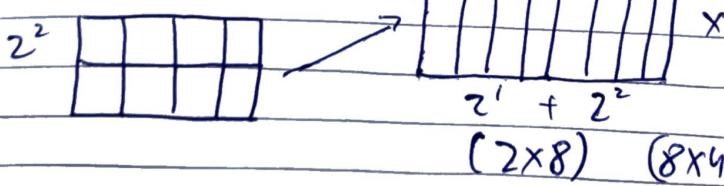
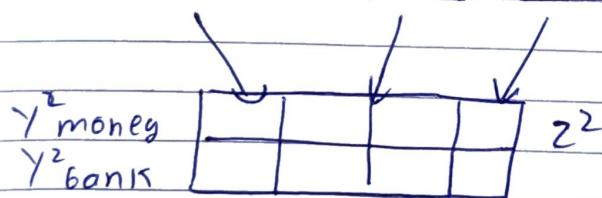
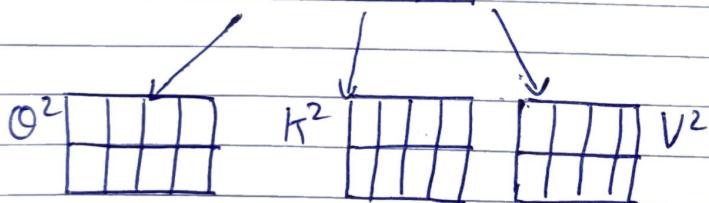
(i) Multi-head attention is a mechanism used in transformers to enhance the model's ability to focus on different parts of the input sequence simultaneously and capture various types of relationship and dependencies.

## (ii) Working

Sentence = "money bank", let multi head = 2



We perform this linear transformation so that we get the resultant embedding same size as input embedding.



The value of this grid is random initially.

## 14.7 Positional Encoding

(i) Transformers process all tokens in parallel. So without positional encodings the model would have no information about the order of the tokens in the sequence.

### (i) Formula

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

PE are positional encoding for the even and odd dimensions of the position pos

pos is the position in the sequence

i represents the index of position encoding vector

$d_{model}$  is the dimension of the model  
(the size of embedding vectors)

(ii) Example sentence - "Tariq Singh"  
embedding vector size = 6 (1x6)

Pos of Tariq = 1

Pos of Singh = 2

$d_{model} = 6$

Tariq embedding = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6]

Singh embedding = [0.6, 0.5, 0.4, 0.3, 0.2, 0.1]

now we need to calculate

$$\text{Tariun Positional encoding} = [?, ?, ?, ?, ?, ?, ?]$$

$$\text{Singh Positional encoding} = [?, ?, ?, ?, ?, ?, ?]$$

For Tariun

$$PE(0,0) = \sin\left(\frac{0}{10000^{0/16}}\right) = \sin(0) = 0$$

$$PE(0,1) = \cos\left(\frac{0}{10000^{1/16}}\right) = \cos(0) = 1$$

$$PE(0,2) = \sin\left(\frac{0}{10000^{2/16}}\right) = \sin(0) = 0$$

$$PE(0,3) = \cos\left(\frac{0}{10000^{3/16}}\right) = \cos(0) = 1$$

$$PE(0,4) = \sin\left(\frac{0}{10000^{4/16}}\right) = \sin(0) = 0$$

$$PE(0,5) = \cos\left(\frac{0}{10000^{5/16}}\right) = \cos(0) = 1 \rightarrow [0, 1, 0, 1, 0, 1]$$

Similarly we find = [0.84, 0.54, 0.0001, 1.0, 0.0001, 1.0]  
 Singh positional encoding.

Once we find the Positional encoding

Final embedding = Embedding + Positional Encoding

$$\text{Final embedding}_{\text{Tariun}} = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6] + [0, 1, 0, 1, 0, 1] = \underline{\text{Result}}$$

Final embedding Singh = Similarly we find this

(ii) without Positional encoding these 2 sentence will have same meaning  $\rightarrow$  Taran kill lion  
lion kill Tarun

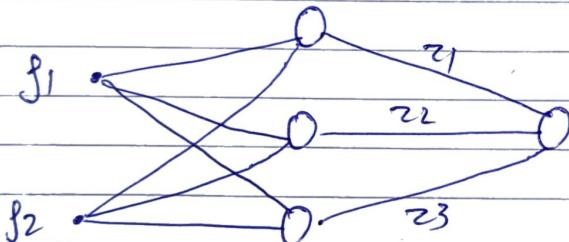
(iii) why we are using Sin and cosine

(i) Distinctive Pattern (ii) Smooth Variation

#### 14.8 Layer Normalization

(i) Unlike batch normalization which normalize across the batch dimension, layer normalization normalize across the features of the every individual sample.

(ii)



$z_1$	$z_2$	$z_3$	$g_1$	$g_2$	$g_3$
a	b	c	d	e	f

we normalize layer by layer.

(iii) why we don't use Batch normalization

- Dependency on Batch size
- Sequential Data Handling
- Does not support parallel processing.