

FINAL VERDICT

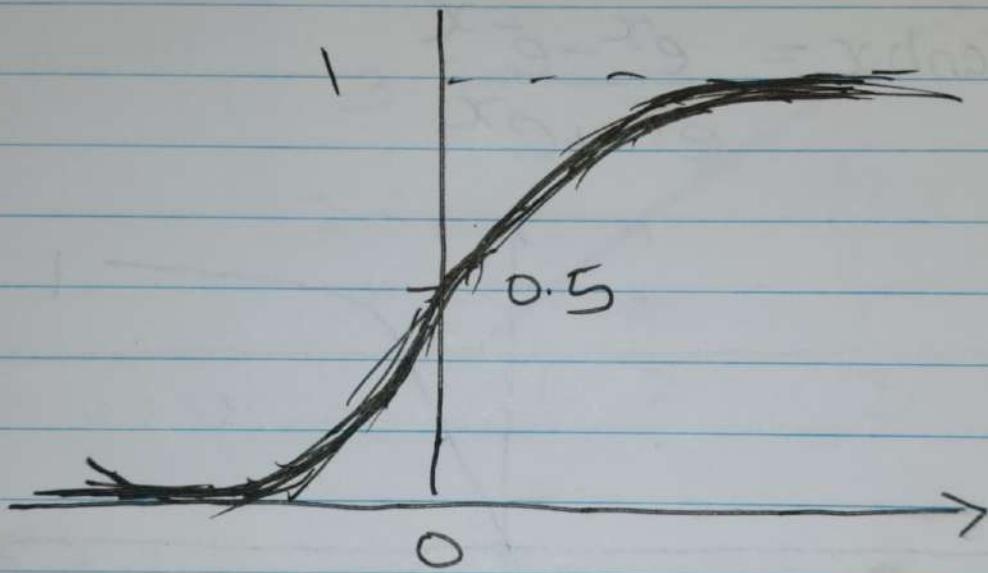
OF DEEP LEARNING

↳ Assume a Human Body

- NLP: Gave this body basic language & grammar knowledge.
- CNN → Gave eyes to this Body.
- RNN | LSTM: Gave Memory capability to this Body.
- GAN: Gave Imagination Power (High Quality) to this Body ↳ Painting
- Auto Encoder: Techniques to learn things in short cut (By making pattern)
→ GAVE
↳ VAE → Gave Body Imagination Power
↳ Variation Power like slight changes

ACTIVATION f(x)

$$\rightarrow \text{Sigmoid} \quad S(x) = \frac{1}{1+e^{-x}}$$

Merits

↳ ① Limits output $(0, 1)$

② Differentiable continuous
③ Used in Binary classification

Demerits

↳ ① Vanishing Gradient

↳ ② Computationally expensive

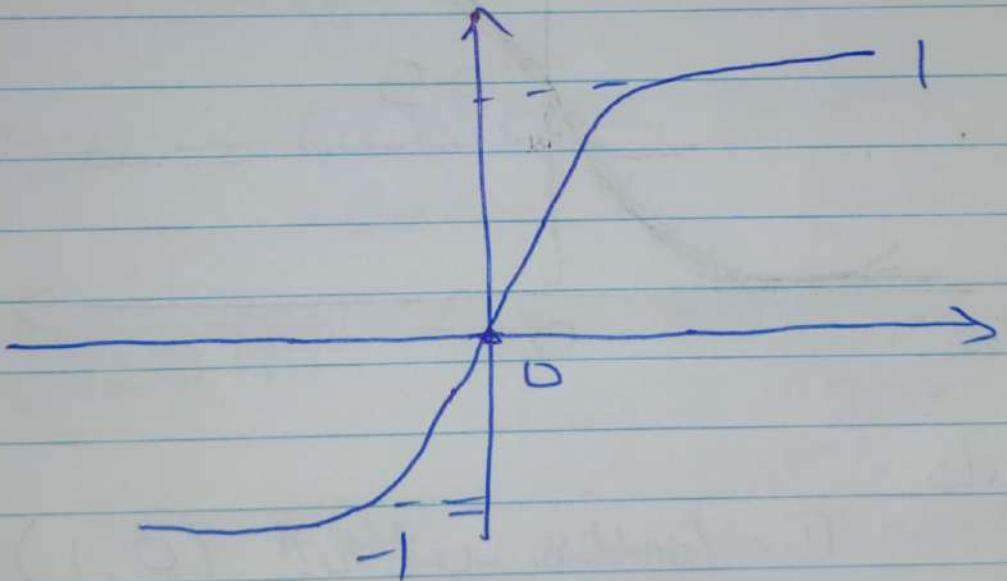
③ Not centred at zero

→ Tanh

$\text{Tanh} x =$

$$\text{Tanh } x = \frac{\sinh x}{\cosh x}$$

$$\text{Tanh } x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



Merits

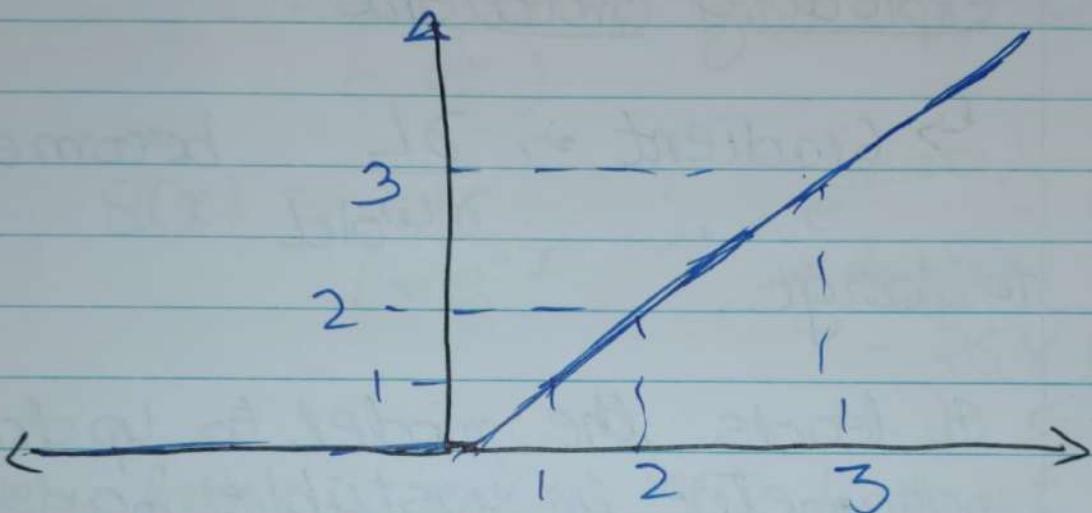
- ↳ (I) centered at zero
- ↳ (II) range $(-1, 1)$
- ↳ (III) continuous & differentiable

Cons

- ↳ (I) vanishing Gradient problem
- ↳ (II) computationally expensive

→ ReLU (Rectified Linear Unit)

$$f(x) = x^+ = \max(0, x)$$



Merits

- ① Easy to compute
- ② Don't cause Vanishing Gradient problem
- ③ Outperforms where the output should be non-ve.
- ④ fast & efficient.
↳ as neurons some are not active

Cons

- ↳ ① Cause Exploding Gradient problem
- ② Not centered

$$\frac{\partial L}{\partial \omega}$$

apsara

Date: _____

(III) Dying ReLU problem

Exploding Gradient

↳ Gradient $\Rightarrow \frac{\partial L}{\partial \omega_{old}}$ becomes
to large.

---> It leads the model to update
parameter in unstable manner

- Q How ReLU cause gradient exploding
problem?

↳ ReLU allows +ve value
to pass through w/o modification
↳ This can lead to this situation

- Q How Sigmoid can lead to Vanishing
gradient problem?

Ans

↳ In sigmoid the gradient is

$$\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$$

& For $x = \text{very large}$, $\sigma(x) = 1$
 $x = \text{very small}$, $\sigma(x) = 0$
 (-ve)

$$S(x) = \frac{1}{1 + e^{-x}} ; \quad \begin{aligned} e^{-(1000)} &= 0 \\ e^{\cancel{(-100)}} &= 26 \times 10^{-42} \end{aligned}$$

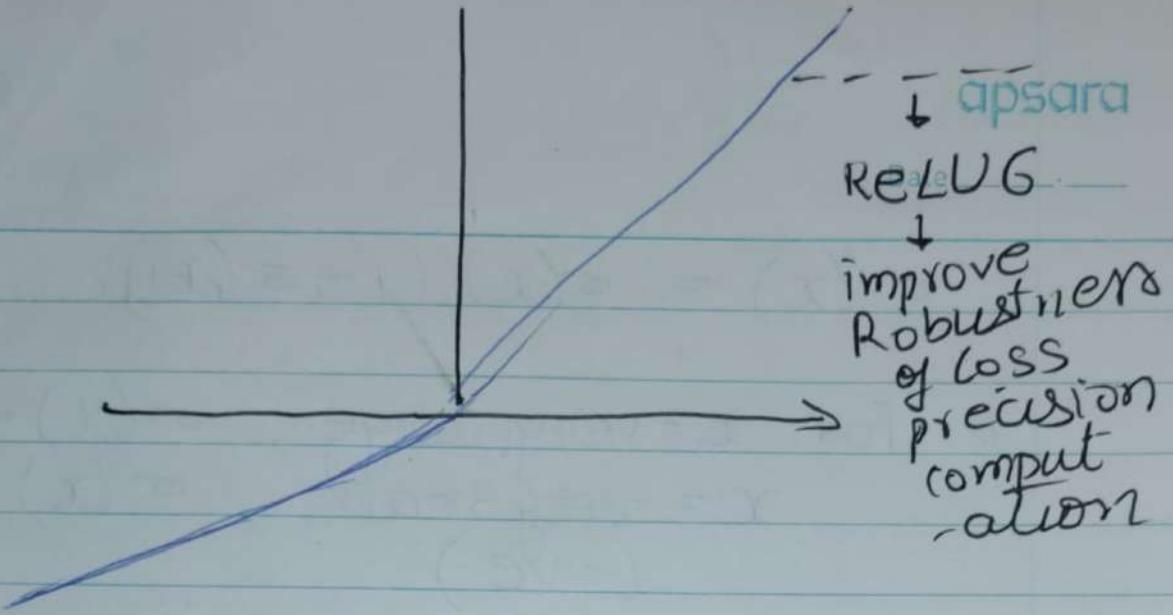
Very large : $\sigma'(x) = 0$ \Rightarrow No wt.
 very small : $\sigma'(x) = 0$ update.

Q Dying ReLU Problem?

↳ When many no. of neurons become inactive & output zero

LEAKY ReLU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise} \end{cases}$$



→ Solves the problem of back dying ReLU.

Pros : ① easy to compute
 ② doesn't cause V.G.P.
 ③ " " D.R.P.
 dying \Rightarrow ReLU.

Cons

- ① exploding gradient
- ② not centered at zero

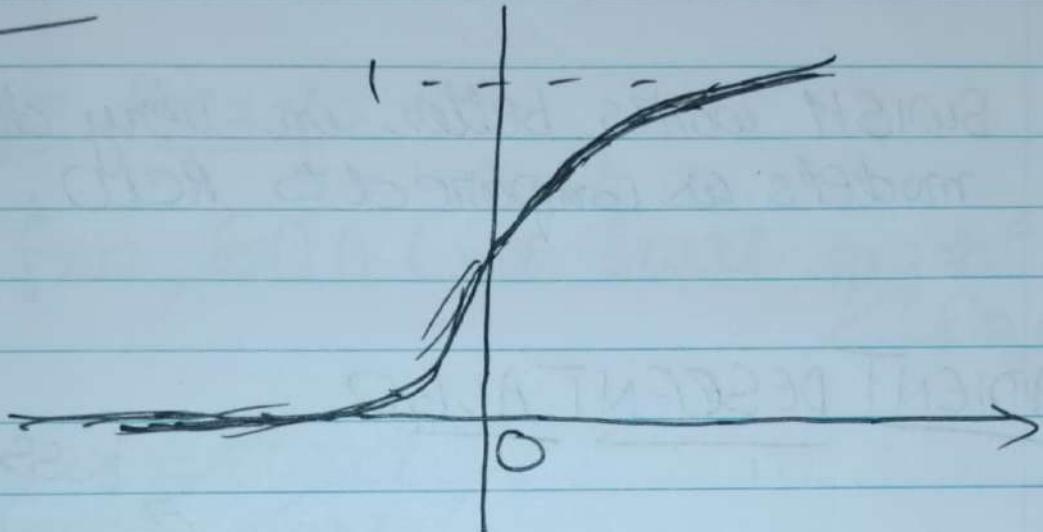
Parameterized ReLU

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

Trainable parameter

ReLU 6

$$\hookrightarrow f(x) = \min(\text{ReLU}(x), 6)$$

Softmax

$$\sigma(z) = \frac{e^z}{\sum e^x}$$

Pros
Cons (①)

① Used for Multiclass classification

① Computationally expensive

SWISH

$$\hookrightarrow \text{SWISH} = x \cdot \text{sigmoid}(x)$$

HARD SWISH

$$h\text{-swish}(x) = \frac{x \operatorname{ReLU6}(x+3)}{6}$$

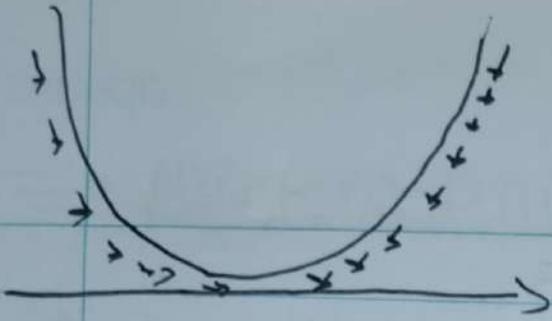
- ★ SWISH works better in very deep models as compared to ReLU.

GRADIENT DESCENT ALGO

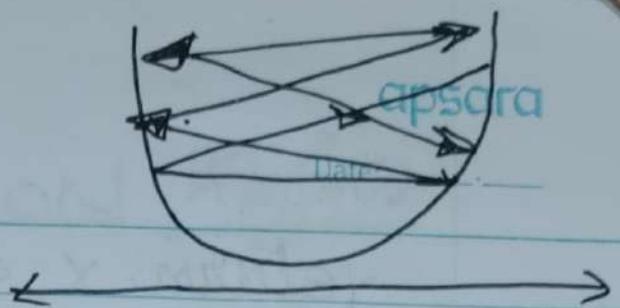
$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} - \eta \frac{\partial L}{\partial \mathbf{w}_{\text{old}}} \quad \begin{matrix} \curvearrowleft \\ \text{loss fcn} \end{matrix}$$

↳ learning rate

Loss function = error fcn = cost function



Small learning rate



Large learning rate

Algo Pseudo

~~fxn GDA (int start, gradient (float, float))~~

~~int / x = start;~~
~~int steps = start;~~

~~fxn x = start .~~

~~steps = [start]~~

~~for _ in range(max_iter):~~

~~diff = learn_rate * gradient(x)~~

~~if np.abs(diff) < tol:~~
~~break~~

~~x = x - diff~~

~~steps.append(x) → History Track~~

return x, steps.

7/0

IV

V

VI

→ BACK PROPAGATION ALGO

- ① Initialize weights randomly
- ② Forward propagate an input to n/o to get the predicted o/p.
- ③ Calculate error between predicted o/p & actual o/p.
- ④ Apply BPA to calculate gradient of loss fxn wrt. each weight
 - ↳ $\frac{\partial L}{\partial w_{jk}}$
- ⑤ Update wt. using Gradient Descent Algo.
- ⑥ Repeat 2-5 steps until weights converges

↓
Back propagation

BPA + GDA

FOR CALCULATION

apsara

Date: _____

IN BACKWARD PROPAGATION

$$\textcircled{1} \quad \frac{\partial E_{\text{total}}}{\partial y_{\text{out}}} = y_{\text{out}} - t_1$$

$$\textcircled{2} \quad \frac{\partial y_{\text{out}}}{\partial y_{\text{in}}} = y_{\text{out}}(1 - y_{\text{out}})$$

$$\textcircled{3} \quad \frac{\partial y_{\text{in}}}{\partial w_8} = h_{\text{out}}$$

↳ output of hidden layer.



will be used to calculate gradient for weight attach with outer layer.

$$\frac{\partial E_{\text{total}}}{\partial w_8} = \frac{\partial E_{\text{total}}}{\partial y_{\text{out}}} \frac{\partial y_{\text{out}}}{\partial y_{\text{in}}} \frac{\partial y_{\text{in}}}{\partial w_8}$$

For weight Attach with Input

$$\boxed{\frac{\partial E_{\text{total}}}{\partial \omega_1} = \frac{\partial E_{\text{total}}}{\partial h_{\text{out}}} \frac{\partial h_{\text{out}}}{\partial h_{\text{in}}} \frac{\partial h_{\text{in}}}{\partial \omega_1}}$$

(D) $\frac{\partial E_{\text{total}}}{\partial h_{\text{out}}} = h_{\text{out}}(1-h_{\text{out}})$

(I) ∂E

(II) $\frac{\partial h_{\text{out}}}{\partial h_{\text{in}}} = h_{\text{out}}(1-h_{\text{out}})$

(III) $\frac{\partial h_{\text{in}}}{\partial \omega_1} = \frac{\partial}{\partial \omega_1} (\omega_1 x_1 + \omega_2 x_2 + \dots + b)$

(1) $\frac{\partial E_{\text{total}}}{\partial h_{\text{out}}} = \frac{\partial E_{01}}{\partial \text{outh}_1} + \frac{\partial E_{02}}{\partial \text{outh}_1}$

(i) $\frac{\partial E_{01}}{\partial \text{outh}_1} = \frac{\partial E_{01}}{\partial y_{\text{out}}} \frac{\partial y_{\text{out}}}{\partial h_{\text{out}}} \frac{\partial h_{\text{out}}}{\partial h_{\text{in}}} \frac{\partial h_{\text{in}}}{\partial \omega_1}$

$$\textcircled{i} \quad \frac{\partial E_{01}}{\partial y_{1\text{out}}} = y_{1\text{out}} - t_1$$

apsara

Date: _____

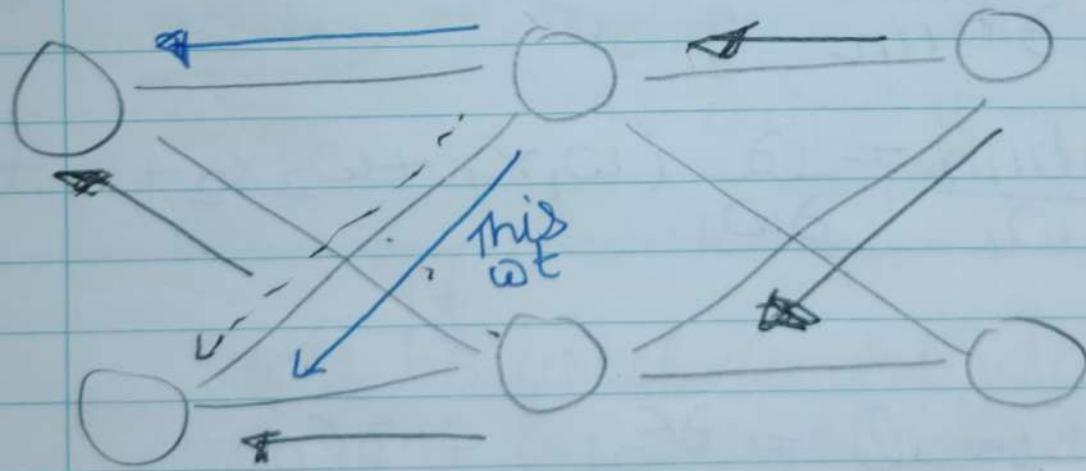
$$\textcircled{ii} \quad \frac{\partial y_{1\text{out}}}{\partial y_{1\text{in}}} = y_{1\text{out}}(1 - y_{1\text{out}})$$

\textcircled{iii} $\frac{\partial y_{1\text{in}}}{\partial t_{1\text{out}}}$ \Rightarrow weight in between them

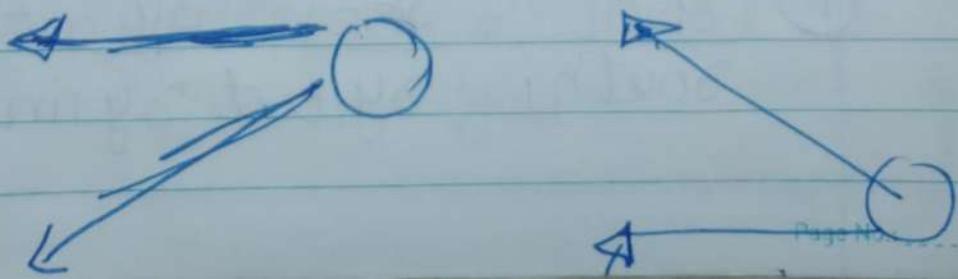
How calculate

$\rightarrow \tau E_{02}$
south,

this wt



Back Propagation always solve like this



GRADIENT DESCENT

apsara
Date: _____

ALGO

SGD

(random some)
datapoints) GRADIENT

Take one
datapoint

Calculate

Error

+ calculate
gradient

update the
weight

BATCH

GRADIENT

Take all
the sample

Calculate
Error

Update
the weight

MINI -
BATCH
GRADIENT

Take a Batch

Calculate
the Error

Update
the weight

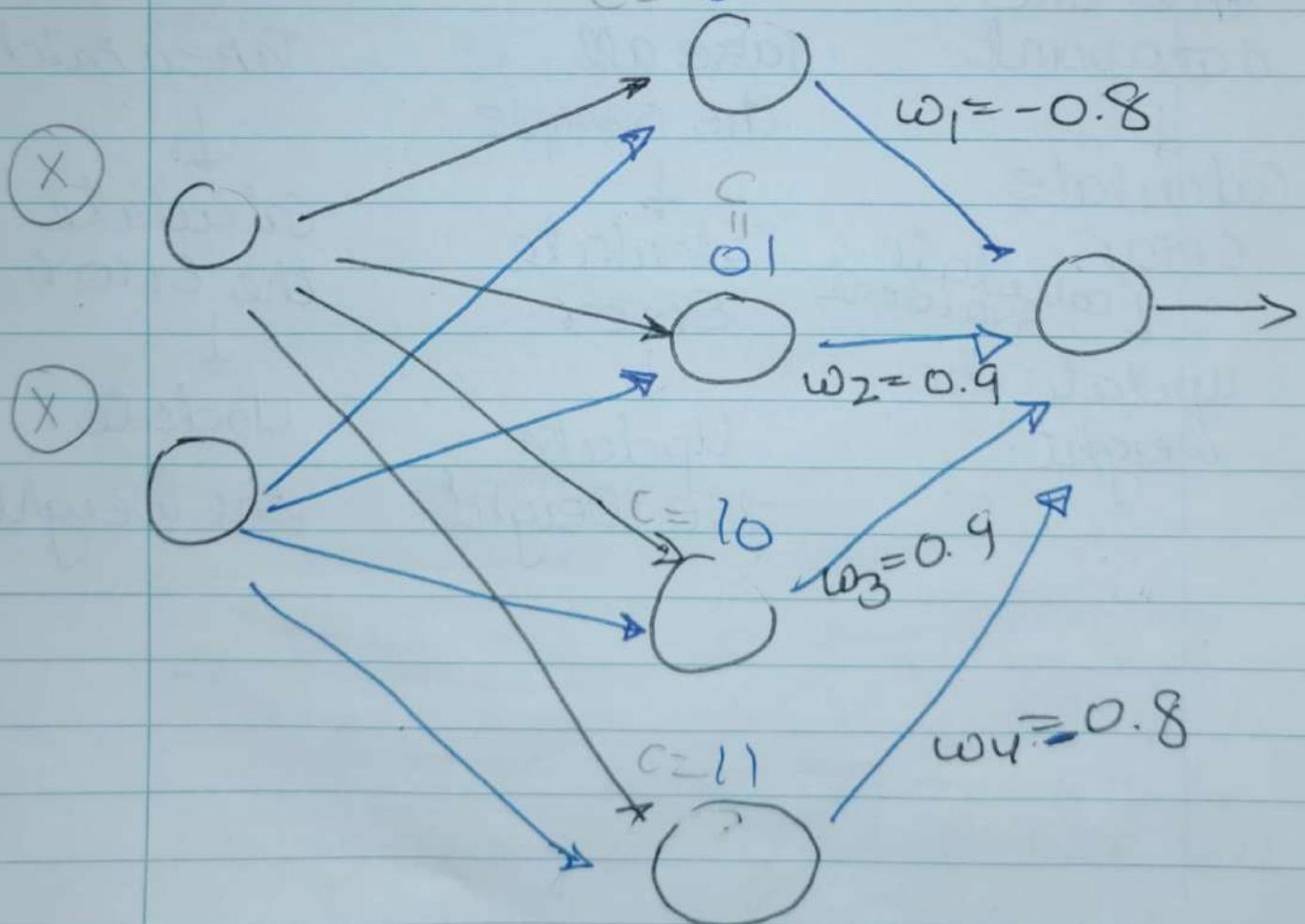
RBF

↳ Gaussian Radial Basis Fxn

$$\hookrightarrow h(x) = e^{-\frac{(x - c_1)^2}{\gamma^2}}$$

Construct XOR Gate using RBF Fxn

$$c = \begin{matrix} 00 \\ 01 \\ 10 \\ 11 \end{matrix}$$



Input 90

$$h_1(x) = e^{-x}$$

tion

$$h_1 = e^{-\frac{(x-c_1)^2}{\gamma^2}} ; c = (0,0)$$

$$h_2 = e^{-\frac{(x-c_2)^2}{\gamma^2}} ; c = (0,1)$$

$$h_3 = e^{-\frac{(x-c_3)^2}{\gamma^2}} ; c = (1,0)$$

$$h_4 = e^{-\frac{(x-c_4)^2}{\gamma^2}} ; c = (1,1)$$

Input (0,0)Distance of x from c ,

$$(x-c)^2 = (0-0)^2 + (0-0)^2 = 0$$

$$-(0)/2$$

$$h_1 = e^{-0} = 0$$

For input $(0,1)$

Distance of x from c_2

$$= (x - c_2)^2 = (0 - 0)^2 + (0 - 1)^2 \\ = 1$$

$$h_2 = e^{-\frac{1}{2}} = e^{-0.5} = 0.6$$

$$\boxed{h_2 = 0.6}$$

Distance of x from c_3

$$= (x - c_3)^2 = (0 - 1)^2 + (0 - 0)^2 \\ = 1$$

$$\boxed{h_3 = 0.6}$$

Distance of x from c_4

$$= (x - c_4)^2 = (0 - 1)^2 + (0 - 1)^2 \\ = 2$$

$$h_4 = e^{-1} = 0.4.$$

$$\text{output} \Rightarrow 1 \times -0.8 + 0.6 \times 0.9 + 0.6 \times 0.9 \\ + (-0.8) \times 0.4$$

$$\Rightarrow \cancel{0.3} - 0.04$$

$$0/p \Rightarrow \underline{\square}$$

For input (1,0)

Distance of x from $c_1(0,0)$

$$(x - c_1)^2 = (1-0)^2 + (0-0)^2 = 1$$

$$h_1 = e^{-1/2} = 0.6$$

Distance from x from $c_2(0,1)$

$$(x - c_2)^2 = (1-0)^2 + (0-1)^2 = 2.$$

$$h_2 = e^{-1} = 0.4.$$

Distance of x from $c_3(1,0)$

$$h_3 \quad (x - r)^2 = (1-1)^2 + (0-0)^2 = 0$$

$$h_3 = 1$$

$$\boxed{h_3 = 1}$$

Distance of $(1,0)$ from $c_4(1,1)$

$$(x - c_4)^2 = (1-1)^2 + (0-1)^2 \\ = 0+1=1$$

$$h_4 = e^{-1/2} = 0.6$$

$$0/p \Rightarrow h_1\omega_1 + h_2\omega_2 + h_3\omega_3 + h_4\omega_4$$

$$= 0.6x - 0.8 + 0.9 \times 0.4 + 0.9 \times 1$$

$$+ 0.6x - 0.8 = 0.3$$

$$\boxed{0/p = 0.3}$$

$$\hookrightarrow 0/p = 1$$

11) by for nearest two

Hypothesis Testing

↳ is a fundamental statistical method used to make inferences or draw conclusion about a population based on data sample data.

i) NULL Hypothesis ($=$)

↳ It is general given statement that there is no relationship b/w two measured cases difference

ii) Alternate Hypothesis (\neq)

↳ Two given state

↳ There is a effect or difference

iii) Level of Significance

↳ Degree of significance in which we accept or reject the Null Hypothesis.

--> Level of significance is selected 5%.
↳ α should be 95% confident

HYPOTHESIS TESTING

T-TEST

apsara

Date: _____

Q
=

A random sample of size 20 from a normal population gives a sample mean of 42 & sample standard deviation (S.D) of 6.

Test the hypothesis that population mean is 44.

i) Z-Test : sample size (n) > 30

ii) T-Test : sample size (n) ≤ 30

So)

$$n = 20$$

$$\bar{x} \text{ (sample mean)} = 42$$

$$s = 6$$

(standard deviation)

mean of population (μ_0) = 44

Step 1: Null hypothesis (H_0): $\mu = 44$

[i.e. population mean is 44]

Step 2: Alternative hypothesis (H_1):

* Technique of Identifying one Tailed & Two Tailed Test

① Use one Tailed Test

↳ less than, more than, greater than, smaller than, taller than, superior, inferior, increase, decrease, at least, at most, only, improved

② Use Two Tailed Test \Rightarrow otherwise

✓ For two Tailed Test

↳ $H_1: \mu \neq \mu_0 \Rightarrow \mu \neq 44$.

For one Tailed Test

↳ $H_1: \mu < \mu_0$ if $X < \mu_0$
else $H_2: \mu > \mu_0$ if $X > \mu_0$

So,

$\mu \neq 44$ [i.e. population mean is not 44]

Step 3:

Test statistics: under H_0 test

statistics is given by

T-Test Formula

if S.D
not given

$$\hookrightarrow t_{\text{cal}} = \frac{\bar{x} - u_0}{\frac{s}{\sqrt{n}}}$$

or

$$t_{\text{cal}} = \frac{\bar{x} - u_0}{\frac{s}{\sqrt{n-1}}}$$

Variance is given \Rightarrow S.D given
Then apply

$$S.D = \sqrt{\text{variance}}$$

$$\therefore t_{\text{cal}} = \frac{\bar{x} - u_0}{\frac{s}{\sqrt{n-1}}} = \frac{42 - 44}{\frac{6}{\sqrt{19}}}$$

$$t_{\text{cal}} = -1.45$$

$$\downarrow$$
$$|t_{\text{cal}}| = 1.45$$

Step 4: Level of Significance (α) = assume

$$\alpha = 5\% = \frac{5}{100} = 0.05$$

$$\therefore \text{Degree of Freedom (DOF)} = 20 - 1 = 19 \\ (n - 1)$$

$$|T_{\text{tab}}| = 2.093$$

from Table
require DOF
& α

Step 5: Comparison & Decision

$$|t_{tab}| > |t_{cal}|$$

Hence, it is not significant

So, we accept null hypothesis (H_0)
 & we reject alternative hypothesis (H_1)

→ one-Tail Test

Q A company claims that the mean life of the electric bulbs is 28 months. A random sample of 10 bulbs has the following life in months:

22, 24, 26, 32, 28, 20, 23, 34,
 30 & 43

Test the claim of the company at 5% level of significance

Sol

$$n = 10 \quad (\text{T-Test})$$

$$\bar{X}_0 = 28$$

Calculation of \bar{x} & s

$$\sum x = 282, \quad \sum x^2 = 8378$$

mean $\Rightarrow \bar{x} = \frac{\sum x}{n} = \frac{282}{10} = 28.2$
 sample

$\downarrow s = \sqrt{\frac{1}{n-1} \left[\sum x^2 - \frac{(\sum x)^2}{n} \right]}$

$$= \sqrt{\frac{1}{9} \left[8378 - \frac{(282)^2}{10} \right]}$$

$s = 6.87$

$n = 10, \quad \bar{x} = 28.2, \quad s = 6.87$

$u_0 = 28$

ep 1 : same

One

Step 2: Some Two Tailed Test

Step 3:

$$u > u_0$$

$$(\because \bar{x} > u_0)$$

$$u > 28$$

{ i.e mean life of electric bulb greater than 28 month }

Step 3: Test statistics: Under H_0 Test

statistics is given by,

$$t_{\text{cal}} = \frac{\bar{x} - u_0}{\frac{s}{\sqrt{n}}}$$

$$= \frac{28.2 - 28}{\frac{6.87}{\sqrt{10}}} = 0.09$$

$$|t_{cal}| = 0.09$$

Step 4: Level of Significance : 5%.

$$\begin{aligned} \text{Degree of Freedom (df)} &= n - 1 \\ &= 9 \end{aligned}$$

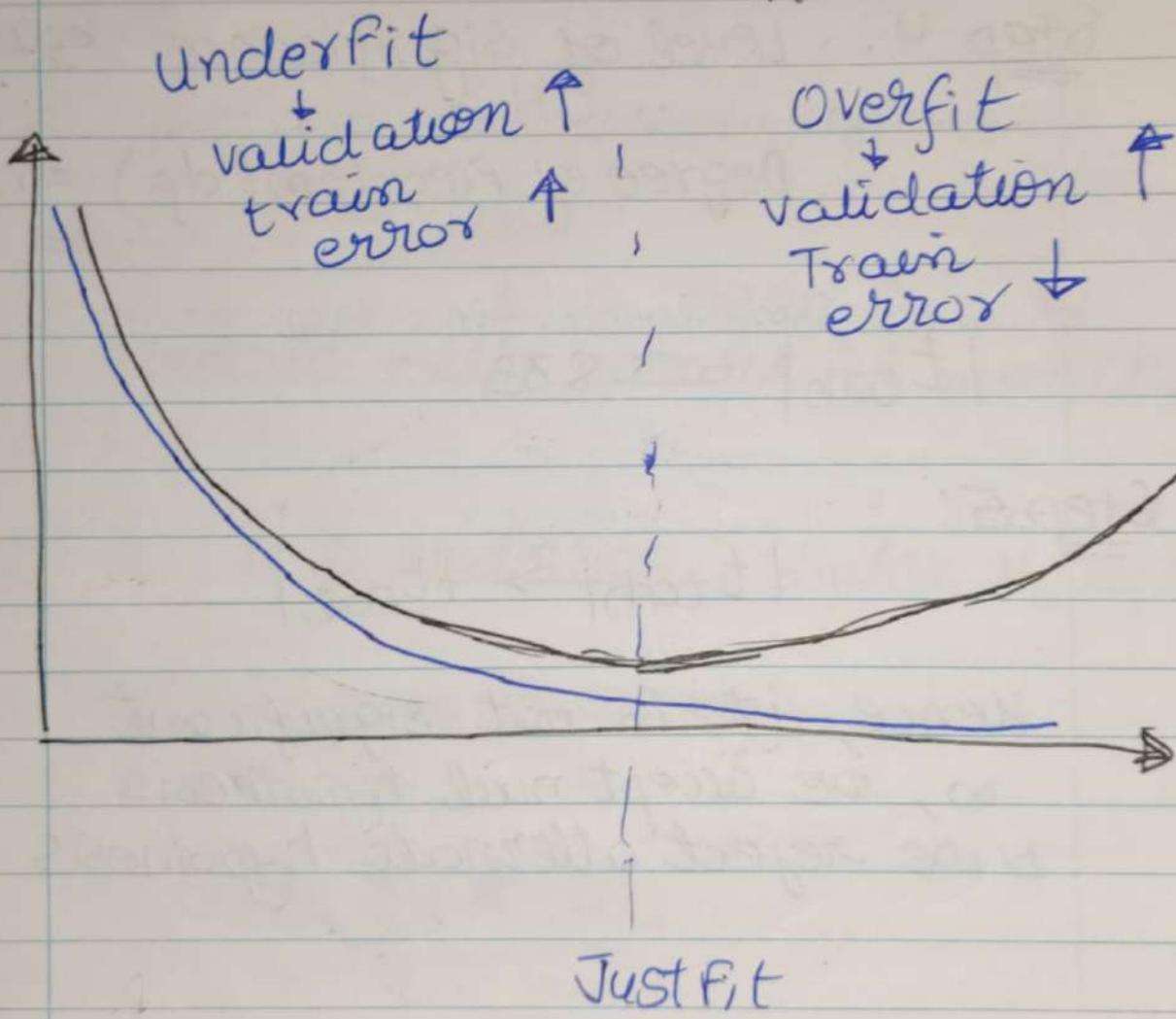
$$|t_{tab}| = 1.833.$$

Step 5: $|t_{tab}| > |t_{cal}|$

Hence, it is not significant
so, we accept null hypothesis.
we reject alternate hypothesis

ERROR GRAPH

FOR UNDERFIT, JUST FIT & OVERFIT



- In regularization, we reduce magnitude of features by keeping same no. of features

REGULARIZATION

apsara
Date: _____

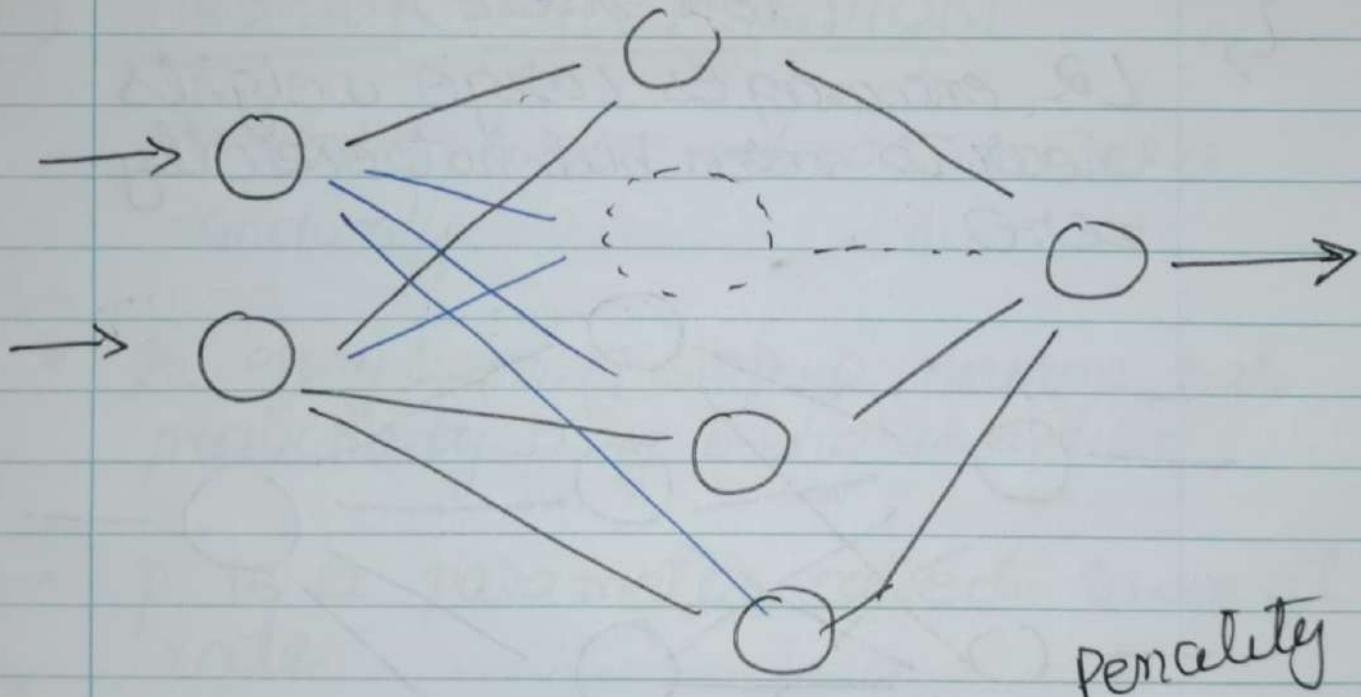
①

L1 REGULARIZATION (LASSO)

(Helps in Feature selection)

↳ In this we add absolute value of weights to the loss function

↳ L1 encourages weights to be as low as 0.0, resulting in more sparse n/w.



Equation

$$\hookrightarrow \text{Cost} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m |w_i|$$

m = No. of features

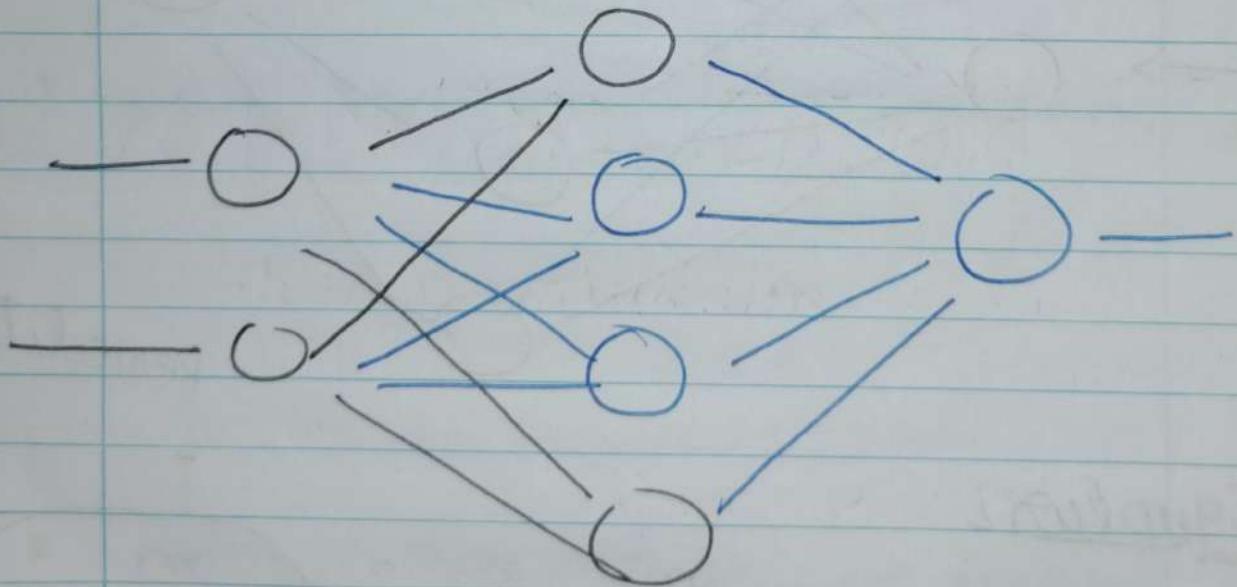
n = No. of examples

y_i = Act target value

⑪ L2 REGULARIZATION

(RIDGE REGRESSION or weight DECAY)

- Add the sum of square values of the weight to the loss function
- You will not have sparse w .
(less sparse)
- L2 encourages large weights near to zero but not exactly zero



Equation

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m w_i^2$$

Parameter ($0 \rightarrow 1$), fully penalty
no penalty \rightarrow value
apsara

↳ How much attention to pay to this penalty.

- Penalty too strong \rightarrow Underfit
- Penalty too weak \rightarrow Over fit.

(III) DROPOUT REGULARISATION

↳ Make some neurons inactive randomly.

- In every training step, a neuron has probability P of being inactive
- P is a parameter called dropout rate

(IV) EARLY STOPPING

↳ Stop the learning once the validation error is minimum

⑤ DATA AUGMENTATION

↳ Create variations of the data points to enrich the training data

* BIAS VARIANCE TRADEOFF

↳ Finding a Proper Balance b/w Bias & Variance is known as Bias variance Tradeoff.

(2015 introduce)

⑥ BATCH NORMALIZATION (STANDARDISATION)

NORMALIZATION

↳ collapse ip between 0 & 1

Standardisation

↳ making mean 0 & variance 1.

- Standardisation is part of Normalization

Date: _____

Types of Normalisation

- ① Min-Max Normalisation: scales the data b/w 0 & 1

$$x_{\text{Norm}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

- ② Z-Score Normalization (Standardization): scales the data to have mean of 0 & s.d of 1.

$$x_{\text{standardized}} = \frac{x - \mu}{\sigma}$$

- ③ Max Abs Normalization → Scales the data between 0-1 & 1.

$$x_{\text{max-abs}} = \frac{x}{|x_{\text{max}}|}$$

$-x -x -x -x -x -x$

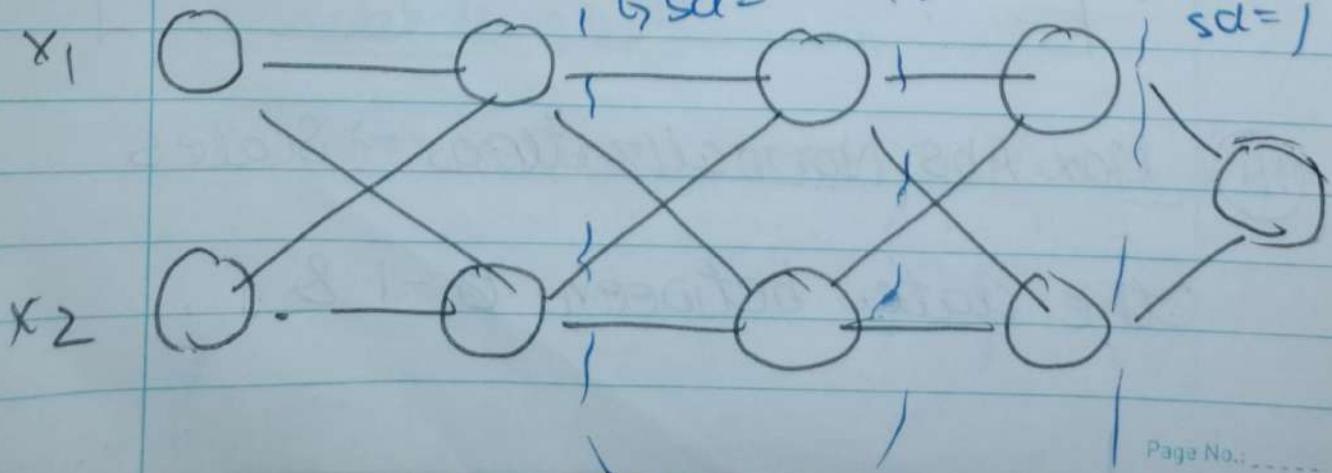
→ There is a unique type of Normalization that is used in Batch Normalization

i.e

$$\text{Batch Normalization} \Rightarrow \hat{x}(i) = \frac{x^{(i)} - \mu_B}{\sqrt{\sigma^2 + \epsilon}}$$

→ Batch Normalisation have mean 0 & standard deviation 1.

$$\begin{array}{l} \rightsquigarrow \text{O/P} \\ \text{mean} = 0 \\ \text{sd} = 1 \end{array} \quad \begin{array}{l} \rightsquigarrow \text{O/P} = \text{mean} = 0 \\ \text{sd} = 1 \end{array} \quad \begin{array}{l} \rightsquigarrow \text{mean} = 0 \\ \text{sd} = 1 \end{array}$$



How BATCH NORMALIZATION WORKS

→ Output of each hidden layer is standardized to have mean of 0 & variance 1 in order to achieve faster convergence

Steps

- ① Takes a mini-batch
 - ↳ (i) Calculate mean
 - (ii) Calculate variance

$$\text{④ } \bar{x}_B = \frac{x_1 + x_2 + x_3 + \dots}{n}$$

$$\text{variance} \Rightarrow \sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x}_B)^2$$

- ② Normalize the batch

- ③ Scale & Shift to ensure model can still learn effectively

$$\boxed{y_i = \gamma \hat{x}_i + \beta}$$

$\gamma \rightarrow \text{gamma}$ γ are trainable
 $\beta \rightarrow \text{Beta}$. β parameter

* weights need regularization where biases don't

L1 Regularization (Lasso)

- ① Penalty term is based on the absolute value of model parameter
- ② Produce Sparse Solution
- ③ Selects the subsets of the most imp. features
- ④ Sensitive to outliers

L2 Regularization (Ridge)

- ① The penalty term is based on the squares
- ② Produce non-sparse solution
- ③ All features are used by the model
- ④ Robust to outliers
→ that's why used significantly

Q Why L1 Regularisation is sensitive to outlier whereas L2 is not?

NORM PENALTIES AS CONSTRAINT

Date: _____

-ED OPTIMIZATION

- Constraint Optimization is about finding the best solution to problem while sticking to certain rules or constraint

Constraint Problem

↳ Want to minimize cost but also need to make sure you don't exceed budget

→ To solve these types of problem we use lagrange fcn

$$L(\theta, \lambda) = J(\theta, x, y) + \lambda \underbrace{S(\theta)}_{\substack{\text{only cost} \\ \text{fcn}}} - \lambda R$$

↑
const.
Penalty
constraint

Sol. to constraint problem

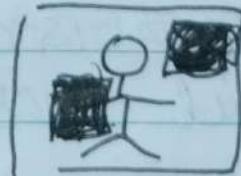
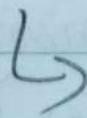
$$\theta^* = \arg \min_{\theta} \max_{\alpha, \alpha > 0} L(\theta, \alpha)$$

DATA AUGMENTATION

↳ Basically adding more synthetic data for training

• Main Data Augmentation Method

- ↳ I) flip
- II) rotate (1° to 3°) \Rightarrow many
- III) crop
- IV) color shift
- V) Noise Addition (Breast Cancer)
- VI) contrast change
- VII) information loss.



• NOISE ROBUSTNESS \Rightarrow (Noise Injection)

- i) Adding Noise to input
 - ii) Adding Noise to hidden layers
 - iii) Adding Noise to weights
 - iv) Injecting Noise
- ↓
used in
RNN

- Semi Supervised LEARNING

- ↳ Technique that combines a small amount of labelled with larger amount of unlabelled data during training
- ↳ This is a form of regularization as it improves the model generalization ability
- Unsupervised learning can provide useful clues how to group

EARLY STOPPING

- Maxout networks, being more flexible than ReLU, can learn more complex features but are still susceptible to overfitting

Q How EARLY STOPPING WORKS?

Epoch 1 : validation loss \downarrow : Save the parameter

Epoch 2 : val loss \downarrow : Update & Save Parameter

Epoch 3: val loss \downarrow : Update & Save Parameter

Epoch 4: val loss \uparrow : Don't Update

Epoch 5: Val loss \uparrow : " "

Epoch 6: val loss \uparrow : Early stopping.

* Mostly used in Deep learning.

→ Early stopping is a regularization technique used in ML to prevent overfitting by halting the training process once model validation loss starts rising.

It monitors the performance during training, & when no significant improvement is observed for certain no. of iterations or epochs, training stopped & the model is then restored to the state where it performed best on val. dataset.

→ During training, it saves parameter whenever val. loss ↑.

①

USE OF 2nd TRAINING

STEP IN EARLY STOPPING

↳ Early stopping requires validation set.

↳ Thus some training data is not fed into model

②

★ Better utilization of entire dataset apsara

Date: _____

To make use of this extra data, a second training step can be conducted after initial training with early stopping.

- 2nd training phase will include previous + valid. data (entire dataset)

Two BASIC STRATEGIES FOR 2nd

TRAINING PROCEDURE

- ① Fine Tuning
- ② Retraining from Scratch

① We use pretrained parameters obtained from initial training & then model is trained on few more epochs on full dataset for further improved.

② Fine Tuning allow model to leverage the knowledge gained from the initial training while adjusting to new datapoints

To make use of this extra data, a second training step can be conducted after initial training with early stopping.

- 2nd training phase will include previous + valid. data (entire dataset)

TWO BASIC STRATEGIES FOR 2nd

TRAINING PROCEDURE

- ① Fine Tuning
- ② Retraining from Scratch

- ① We use pretrained parameters obtained from initial training & then model is trained on few more epochs on full dataset for further improved.

(*) Fine Tuning allow model to leverage the knowledge gained from the initial training while adjusting to new data points

Not really used

⑪

Retrain the model from scratch means discarding the parameter update from first training & retrain the new model on whole dataset

↳ if validation dataset is too small (is one case, where we do this)

Early Stopping

- ↳ ① Early stopping effectively controls capacity of the model by limiting the number of training iterations
- ↳ ② It provides straight forward to regularize model w/o altering its str. or adding complexity

PARAMETER Tying REGULARISATION

- ↳ In this, instead of each layer in neural n/w having its own set of parameters (weights), certain

layers share same parameters.

PARAMETER SHARING

↳ refers to practice of using the same set of parameters (weights) across different layers of a network.

In CNN, the kernel parameters are shared

BAGGING REGULARISATION

↳ Bagging involve training multiple models on different subsets of the training data & then aggregating their predictions to produce final output.

Bagging \rightarrow Bootstrap aggregating

- Use concept of Model Averaging.

~~DROPOUT REGULARIZATION~~

PARAMETERS

- ① These are the entities learned via training from training data.
- ② They are set manually by designer
- ③ Example -> weight & Biases

HYPERPARAMETERS

- ① These are the entities set by the designer before training
- ④ Example ->
 - Learning rate
 - No. of layers
 - No. of neurons in each layer
 - Activation fxn

HYPERPARAMETER TUNING

- ↳ ① Optimizer
- ↳ ② Learning rate
- ↳ ③ Batch size

..... Tuning these manually makes you to try out 1000 more combination

Approaches for Hyperparameter Tuning

① Grid Search

- ↳ Try out every single option

- ↳ very slow

② Random Search (Used mostly)

- ↳ Try out a subset of all the options

- ↳ Work well for many problems

③ Zoom In

CODE → Little Bit Difficult.

Evaluation
of model

TRAINING SET	CROSS-VALIDATION SET	TEST SET
--------------	----------------------	----------

→ weight &
Biases

→ Model Selection
(Hyperparameter Tuning)

ENSEMBLE TECHNIQUE

Bagging
 (Feeding data
 fully to different
 model)

↓
 Random Forest

BOOSTING
 (Feeding data
 sequentially
 to a model)
 $\text{data} \Rightarrow O \rightarrow O \rightarrow O \rightarrow O$

- i) AdaBoost
- ii) Gradient Boost
- iii) XG Boost
- iv) CAT BOOST

Q Why we use the momentum
 in the overall update of weight
 -s in N.N?

↳ So, according to momentum is used
 in SGD or weight updation for
 speed up & faster convergence
 What momentum do is add previous
 fraction (\rightarrow update) to the current update
 of

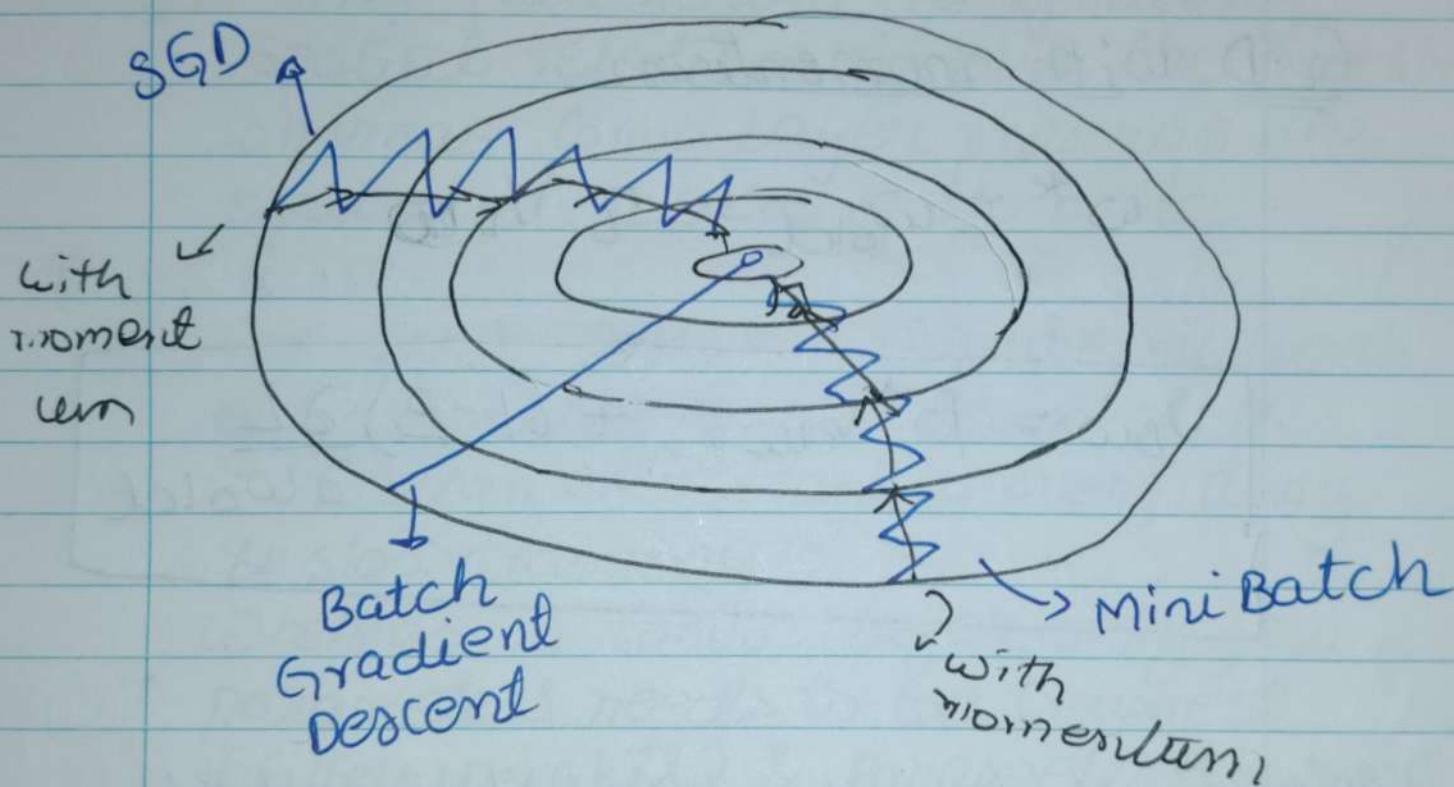
SGD

with Momentum

apsara

Date: _____

- ↳ aim is to remove vertical movement & improve the horizontal movement.



- ↳ lower weightage is given to previous datapoint & higher to current.

$$\bullet \boxed{v_t = \beta v_{t-1} + (1-\beta)\theta_t}$$

↑ current
data
point

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1-\beta)\theta_1$$

$$v_2 = \beta v_1 + (1-\beta)\theta_2$$

Gradient Descent w/o Momentum

apsara

$$\omega^* = \omega_{\text{old}} - \eta \frac{\partial L}{\partial \omega_{\text{old}}}$$

G.D with momentum

$$\omega^* = \omega_{\text{old}} - \eta V_D$$

$$V_D \omega = \beta^* V_D \omega_{t-1} + (1-\beta) \frac{\partial L}{\partial \omega_{\text{old}}}$$

Ans)

(I)

(II)

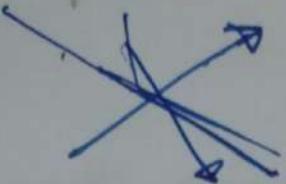
(III)

why Convo. layer preferred over fully connected NW (ANN) in CNN?

- ① ANN flattens the I/P ignoring spatial relationship b/w the pixels whereas Convo. Layer preserve the spatial structure of I/P data.
-) For CNN, ANN will require millions of parameter to be trained & that can also lead to overfitting & slow training. Whereas in Convo. layer only few parameters needs to be trained (filter weights) & parameter sharing also leads to massive reduction in no. of parameters.

CNN \Rightarrow computationally efficient
ANN \Rightarrow not

Why Batch Normalization is required for both input as well as hidden layers?



- (i) Improve Stability
- (ii) Performance
- (iii) Slow convergence
- (iv) Vanishing / Exploding gradient Problem
- (v) overfitting

$\stackrel{?}{=}$ How to get original data back
from PCA?

↳ Eigen vector * (Eigen value)

① zero^{*} = Eigen vector \otimes (Eigen value)^T
mean P. C value

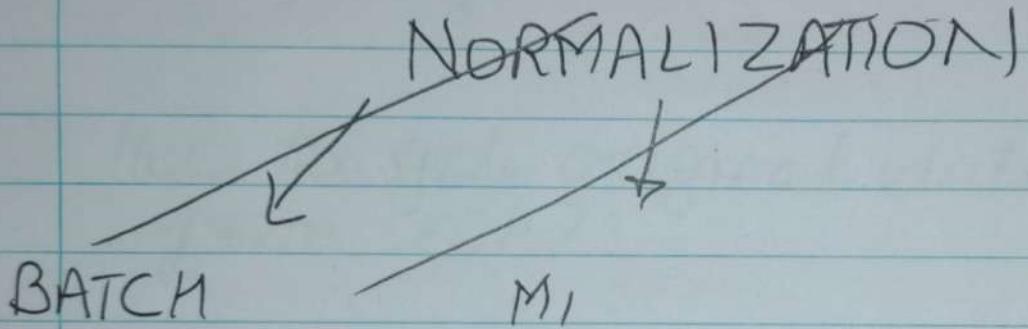
② zero mean = $\stackrel{*}{\text{zero mean}}$ ^T.

③ zero mean
original data = zero mean
+ mean
value

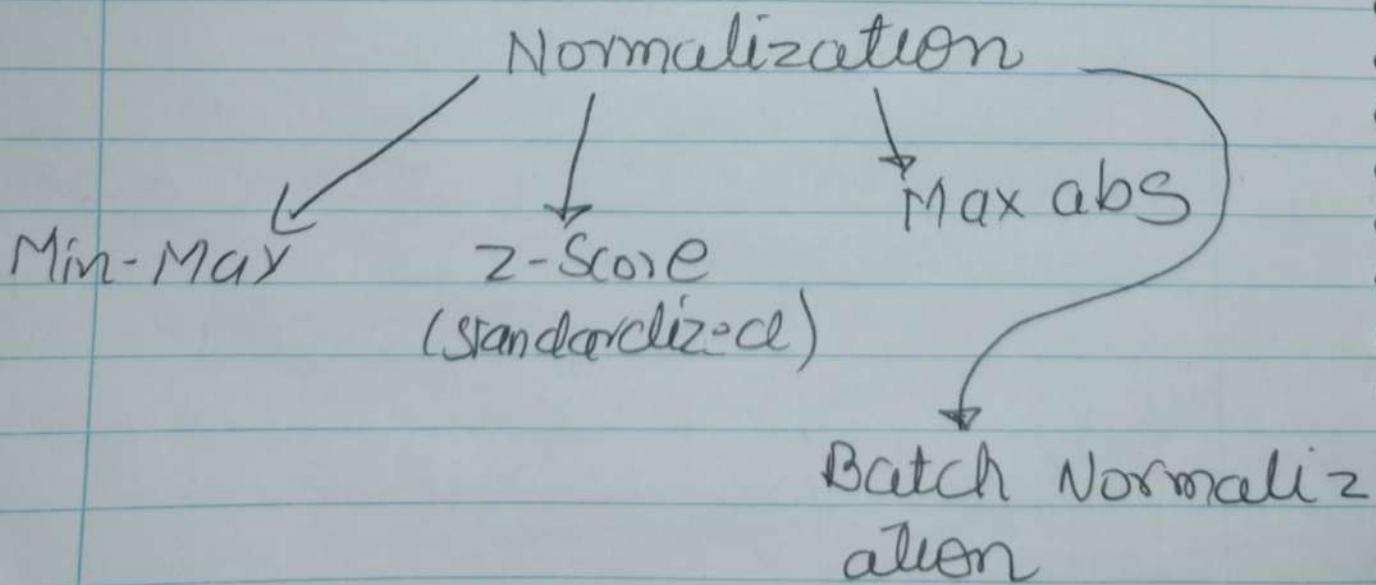
INTERVIEW

Zero Padding : layer used in CNN

to add zero values around the border of a 2D input.



- Normalization : Process that adjust data to a standard scale or distribution



---> Batch Normalization \Rightarrow Mean 0 & standard Deviation 1.

↳ o/p of each neuron is standardize to have mean 0 & variance 1 for faster convergence

★ Reduces the number of iterations required during training.

• Gaussian Blur: smoothens the

image by averaging pixel intensities with their neighbour.

• What is F1 score?

Ans It is a metric used to evaluate performance of a classification model

$$F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

CM

* FP is very dangerous
↳ Person don't have cancer
but ~~not~~ detected

Positive	TP	FP	GO	0
			01	0
			10	0
			11	1
Negative	FN	TN	Not have	cancer
			& Not detected	
			↳ Person don't have But selected	

①

$$\text{Accuracy} \Rightarrow \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FN} + \text{FP}}$$

②

Precision: Performance metric

used to evaluate the accuracy
of model's +ve prediction

$$\text{Precision: } \frac{\text{TP}}{\text{TP} + \text{FP}}$$

③

$$\text{Recall: } \frac{\text{TP}}{\text{TP} + \text{FN}}$$

↳ $\frac{\text{TP}}{\text{TP} + \text{FN}}$

* FN is dangerous.

Hyperparameter

apsara

Date: _____

Learning Rate: Learning rate defines how fast or slow model converges to global minima.

- > Step size \Rightarrow control the size of steps taken to update the model's parameter

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

↓
learningrate

Loss function: MSE

For Regression, Task generally one output neuron is required

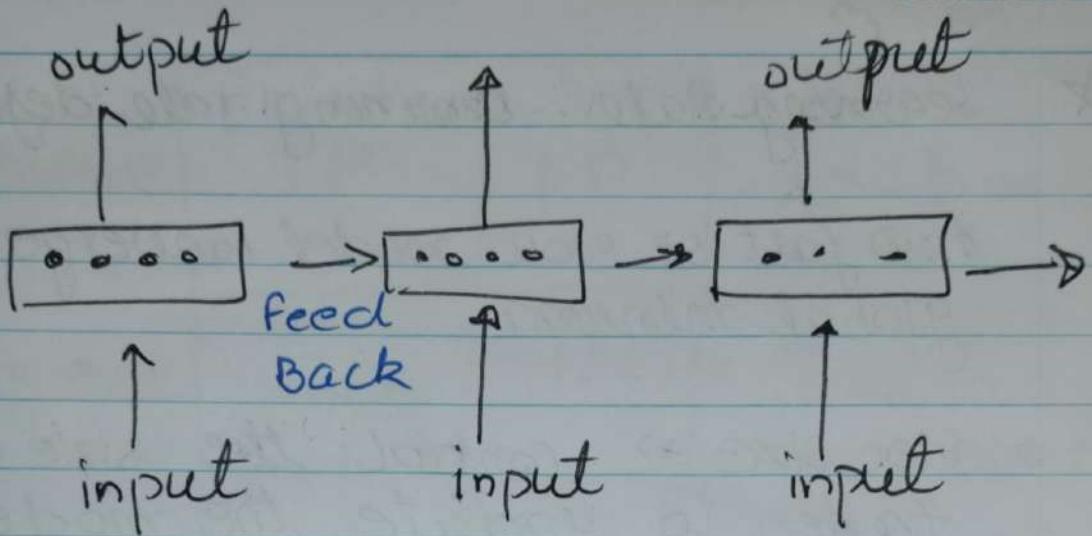
ANN \Rightarrow is a Fully connected NN

output of each previous layer is given input to every next neuron

RNN

RNN are also trained just like NNs

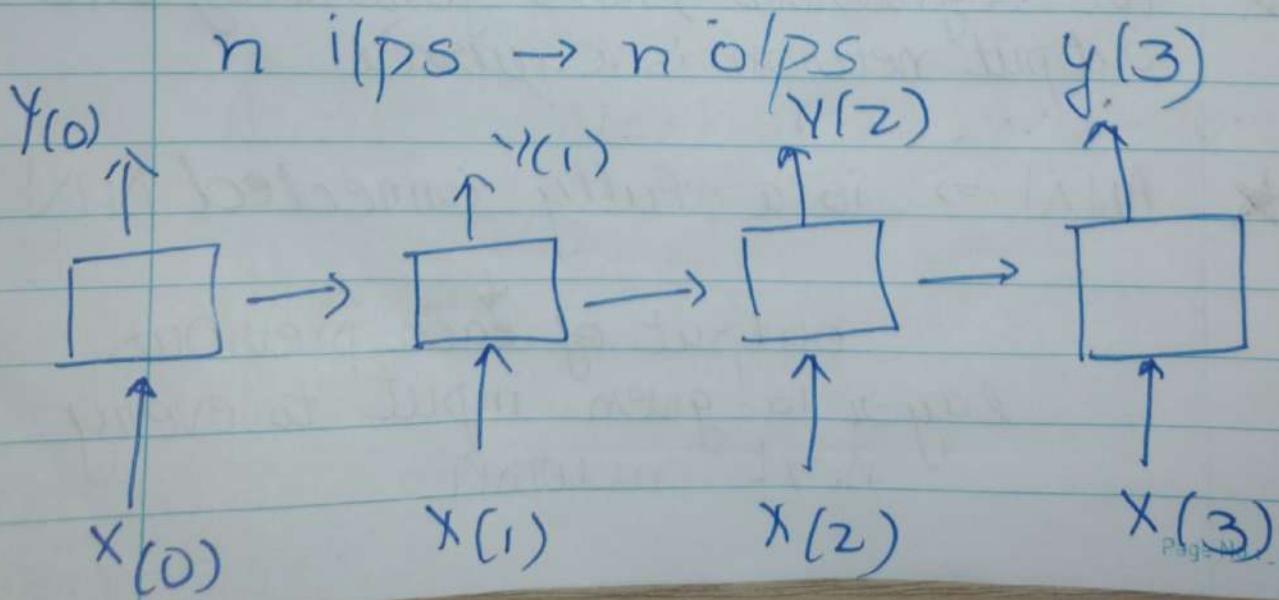
Date: _____



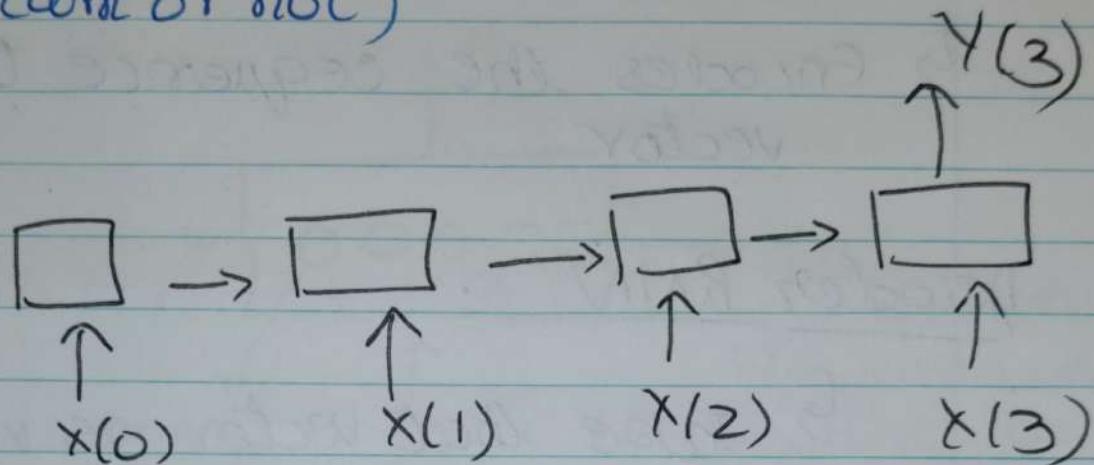
★ Output of previous layer is given as feedback to next layer.

Types of RNN

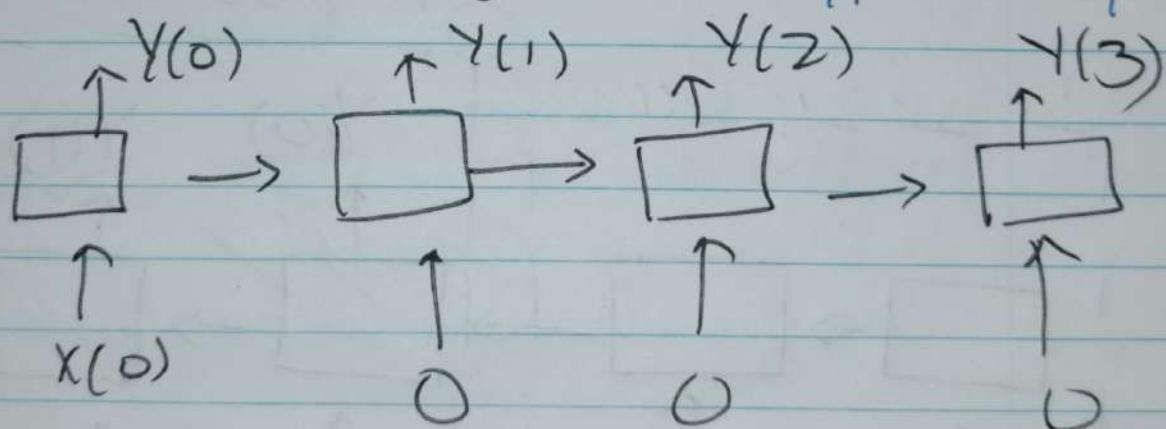
- ① Sequence to Sequence:
(Price forecasting)



- (II) Sequence to vector: $n \text{ i/p} \rightarrow 1 \text{ o/p}$.
(Scam or not)



- (III) Vector to Sequence: $1 \text{ i/p} \rightarrow n \text{ o/p}$



- (IV) Encoder-Decoder:
(Translation)

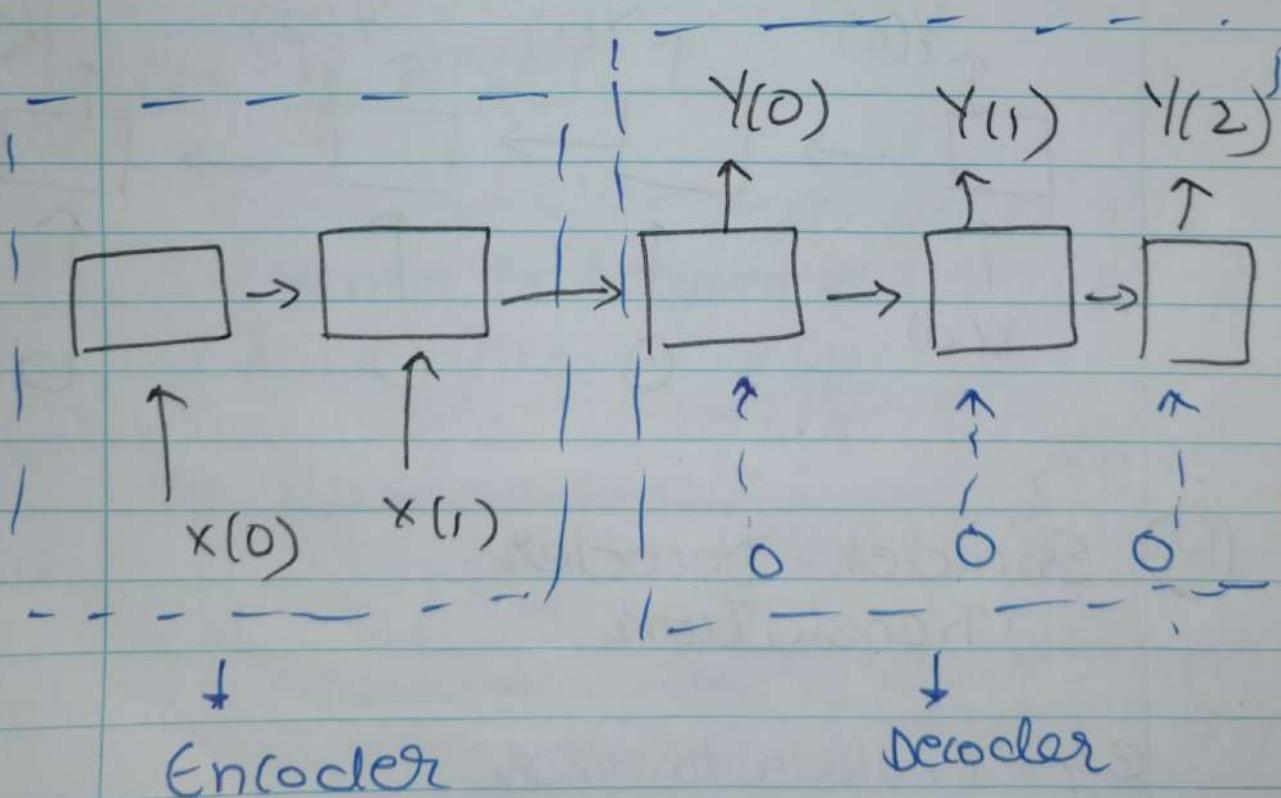
E.g. English to Hindi

Encoder RNN

- ↳ Encodes the sequence to a vector

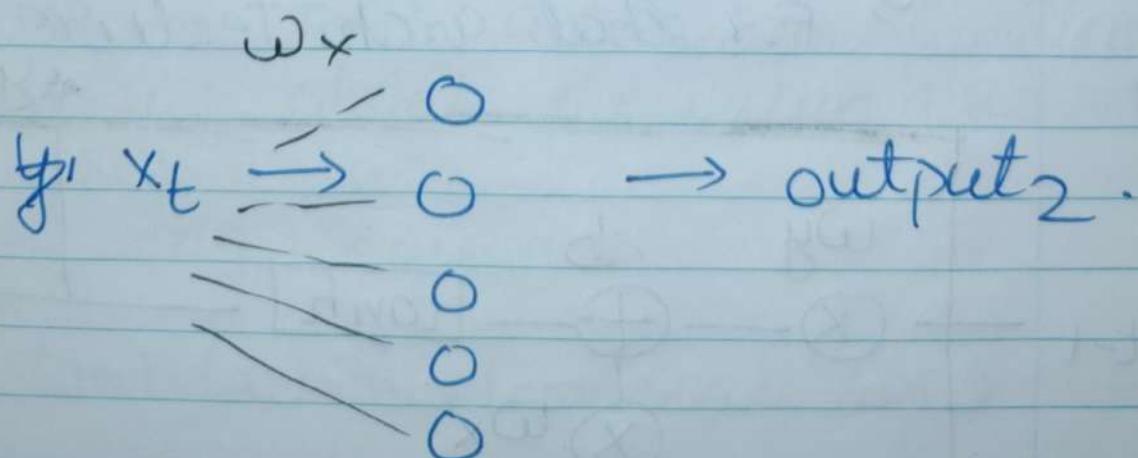
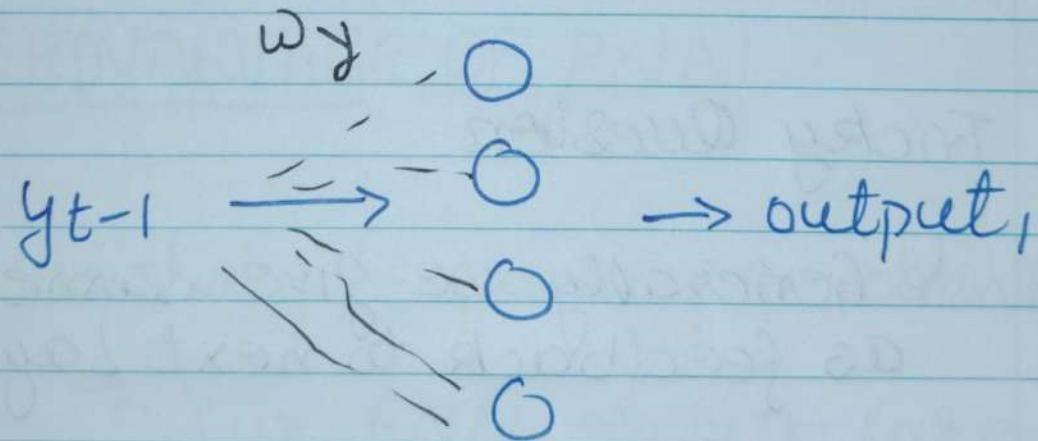
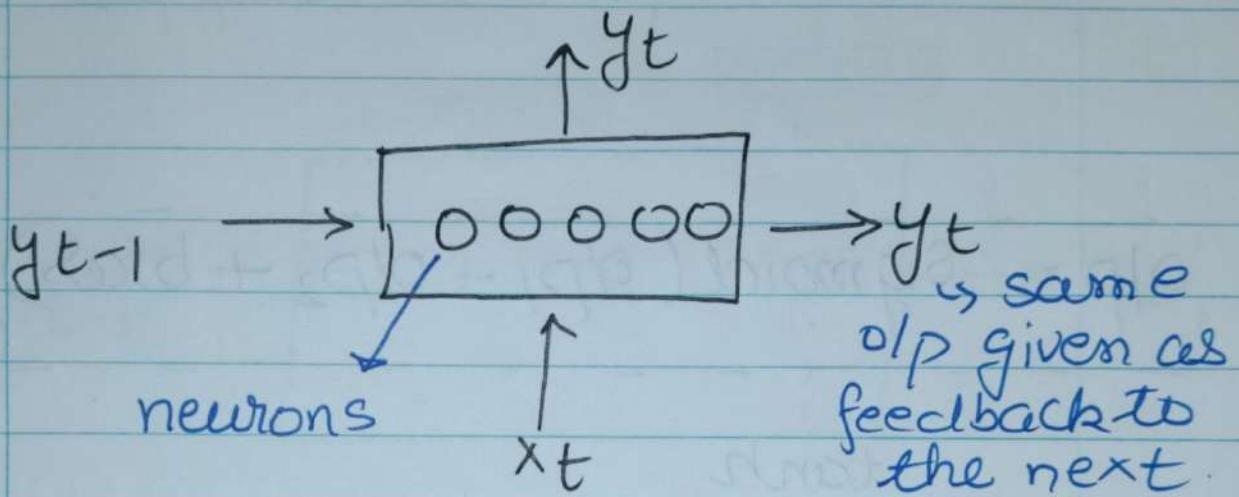
Decoder RNN

- ↳ Take that vector as input
- ↳ Reconstruct the original sequence or generate a related o/p sequence



Q

How output is calculated in RNN?



You can think it as double neuron connection one separate connection with IP & one with feedback

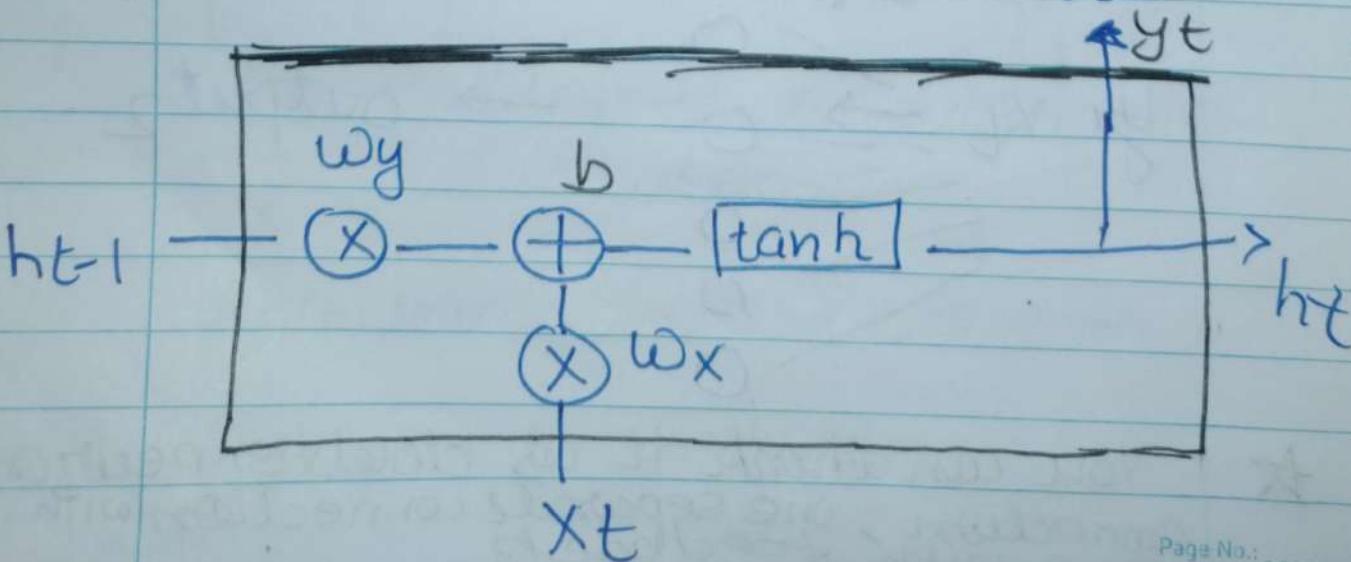
output = activation (output₁ + f_{xn} output₂ + bias)

$$\text{o/p} = \text{Sigmoid}(\text{o/p}_1 + \text{o/p}_2 + \text{bias})$$

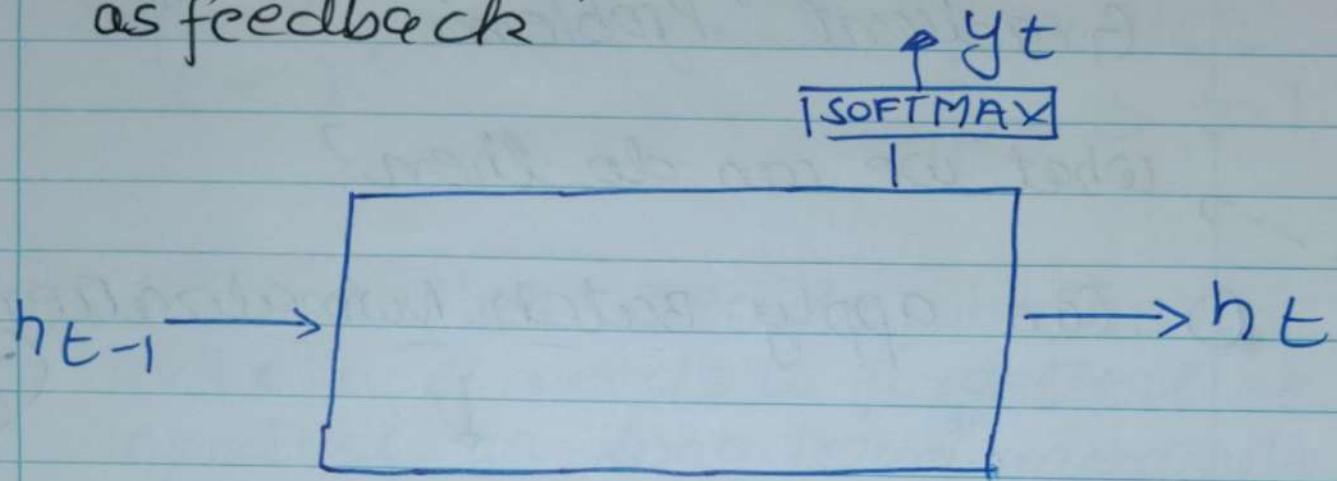
tanh
ReLU

Q Tricky Question

- ↳ Generally we give same o/p as feedback to next Layer
- ↳ For that archi teclure



→ In this case o/p will not be same as feedback



DISADVANTAGE OF RNN

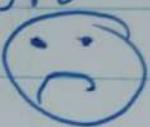
- ① Vanishing / Exploding Gradient:
 - ↳ cuz RNN architecture have Tanh activation function which shrink the value b/w -1 & 1.
 - ↳ Long Sequences.

$$\tanh'(x) = 1 - \tanh^2(x)$$

→ RNN also face exploding Gradient Problem

What we can do then?

① Can apply Batch Normalization



① normalize i/p across each layer



mean 0
variance 1

② Improve Training Speed

③ Long Term Memory does not work well:

They forget the beginning of long sequences:

→ What we can do?

↳ Use LSTM or GRU.

\rightarrow O/p of LSTM $\Rightarrow h_t$

LSTM

\rightarrow everything attach with \leftarrow is gate

apsara
Date: _____

\rightarrow Sigmoid fxn \rightarrow gate controller
 (σ)

O/p \Rightarrow 0 to 1

gate controller: decide whether something needs to be forgotten or inputted in long term memory.

- Forget Gate: which part of long term memory should be removed.
- Input gate: which part of the new info we got should be added to long term memory.
- Output gate: Which part of long term state should be added to the O/p at this timestamp?
- Sigmoid $\rightarrow 1 \rightarrow$ let everything pass
- Sigmoid $\rightarrow 0 \rightarrow$

→ Check on You.

LSTM can be used in:

① Natural language

② Time Series Data:

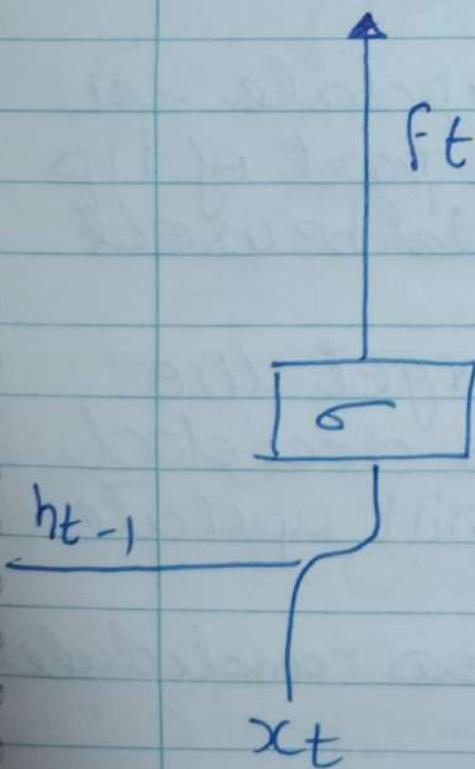
- Stock Price
- Energy usages

★ where the gap b/w the relevant info & predict position is small, RNN can work.

Q where RNN can't work but LSTM can?

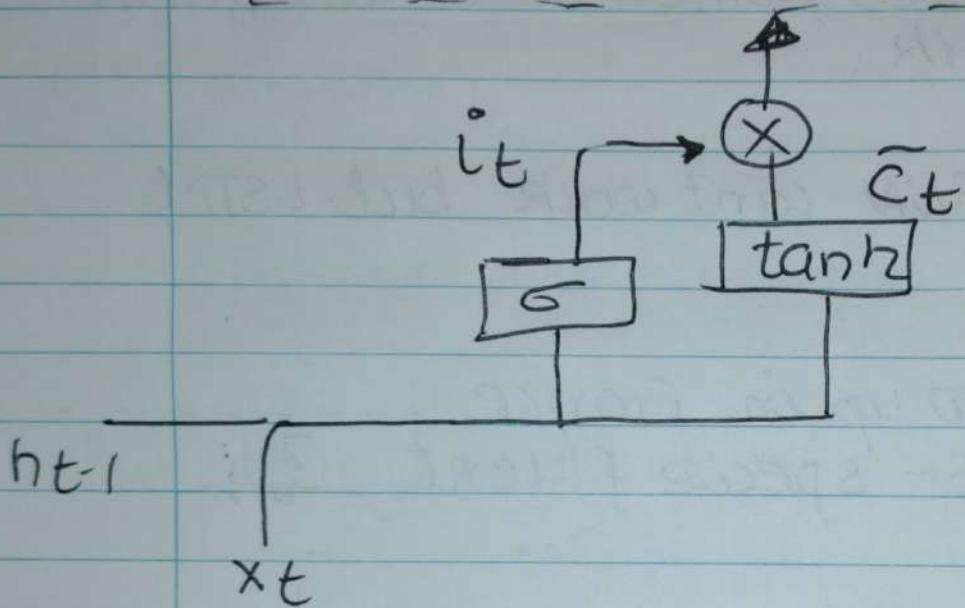
I grew up in France.....
I speak fluent ? -

Forget Layer \Rightarrow



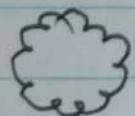
$$f_t = \sigma(w_f^x x_t + w_f^h h_{t-1} + b)$$

Input gate Layer



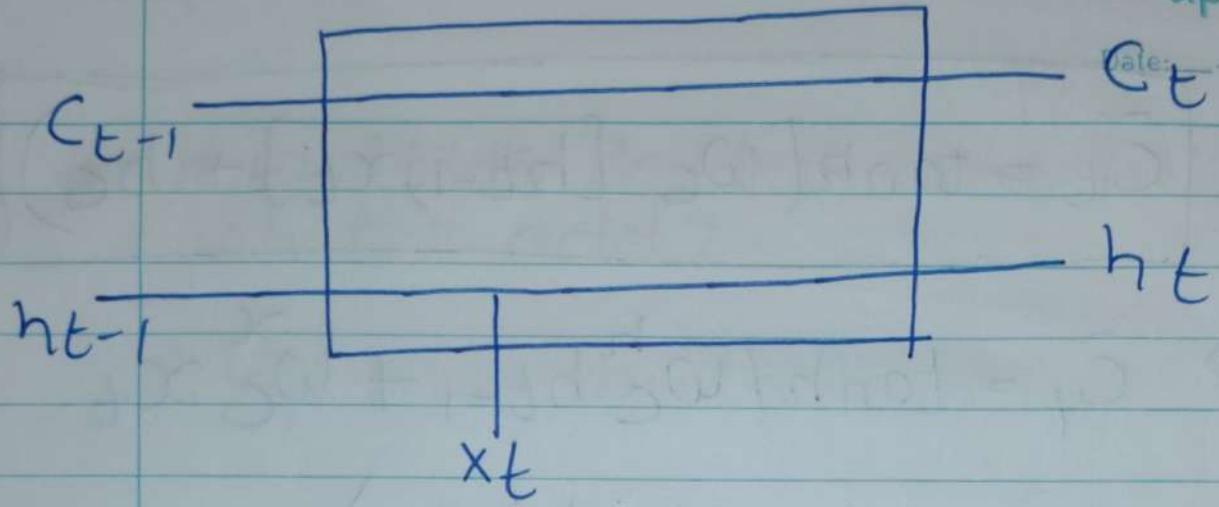
→ $i_t \rightarrow$ Input gate layer decides

which value we will update in
memory cell & which part of i/p
& previous hidden state should be used
to update.



Till Now, we have forgot the
things we want, we decided
which things we will update

→ $\tilde{c}_t \Rightarrow$ a vector of new candidates



$$C_t = f_t * \underline{C_{t-1}} + i_t * \tilde{C}_t$$

$$\begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & - \\ 0 & - & 1 & 0 \end{bmatrix} * \begin{bmatrix} \text{Previous} \\ \text{info.} \end{bmatrix}$$

- $i_t * \tilde{C}_t$

$$\begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} * \begin{bmatrix} \text{New} \\ \text{info.} \end{bmatrix}$$

New info. we can add.

apsara

Date:

$$\tilde{C}_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

$$\Rightarrow \tilde{C}_t = \tanh(w_c^h h_{t-1} + w_c^x x_t + b_c)$$

Forget Layer

$$f_t = \sigma(w_f^h h_{t-1} + w_f^x x_t + b_f)$$

(can
be
written
as)

$$\Rightarrow f_t = \sigma(w_f [h_{t-1}, x_t] + b_f)$$

Input layer

$$i_t = \sigma(w_i^h h_{t-1} + w_i^x x_t + b_i)$$

(can
also
be)

$$\Rightarrow i_t = \sigma(w_i [h_{t-1}, x_t] + b_i)$$

* Back Propagation in DL (Gate)

apsara

Date: _____

* How it is decided what to forget
what to add?

↳ w_f & w_i are trained

Output Gate

$$o_t = \sigma(w_o[h_{t-1}, x_t] + b_o)$$

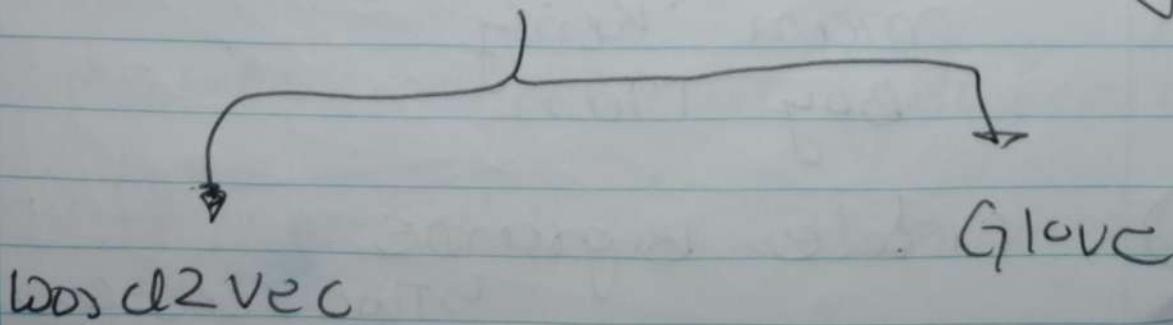
Final OLP

$$h_t = o_t * \tanh(c_t)$$

↳ Check in DL (Gate)
NLP

- x - x - x - x - x - x - x

There are Two Models to
Train Word Embedding



Word 2 Vec

apsara

Date: _____

→ To train Word Embedding via word2vec we use Neural Network

Dataset (Corpus)

The future King is the Prince
Daughter is the princess
Son is the prince
only a man can be King
only a woman can be queen
The - - - - -
-- -- -- -- --

① Remove Stop Words

future King Prince
Daughter princess
Son Prince
man King
Boy Man

② Create Bigrams

↳ Two consecutive words occurring in a sentence

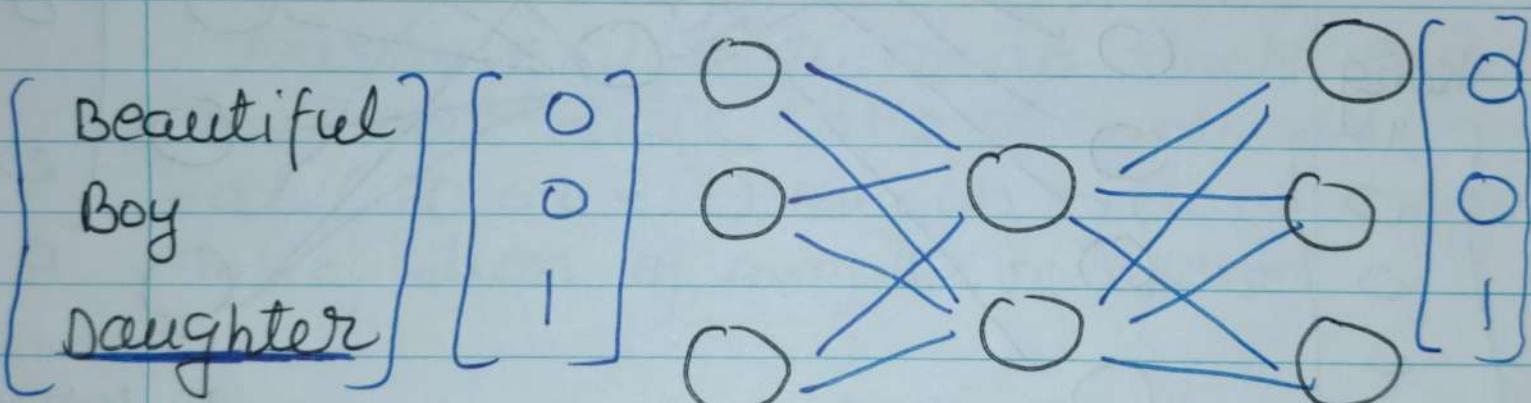
[future, King]

[King, Prince]

[Daughter, Princess]

(III)

Convert words to One hot.



* Split for training
test.

Beautiful
Boy
Princess

→ Models weights we get trained
↳ & we will use that weights
value

When two or more words
are used to predict a single word

Date: _____

CBOW

Future

[0 1
1 0
0 0]

O

O

O

O

O

O

O

O

King

[1 0
0 0
0 0]

O

O

O

O

O

O

O

O

O

O

O

O

O

O

O

O 7

O 1

O 0

O 0

O 0

O 0

Future

Skipgram

Prince

O

O

O

O

O

O

O

O

King

[1 0
0 0
0 0]

[0 0
0 0
0 0]

[0 0
0 0
0 0]

Page No.: _____

SOM

→ Dimensionality
Reduction
method

apsara

Date: _____

↳ Self organizing w/o no supervision is required

Learn on their own through unsupervised competitive learning

SOM : ANN that uses unsupervised learning approach & trained its nw through a competitive learning algo to map multidimensional data into low-dimensional (1D or 2D) which allows easy interpretation of complex problems

Map

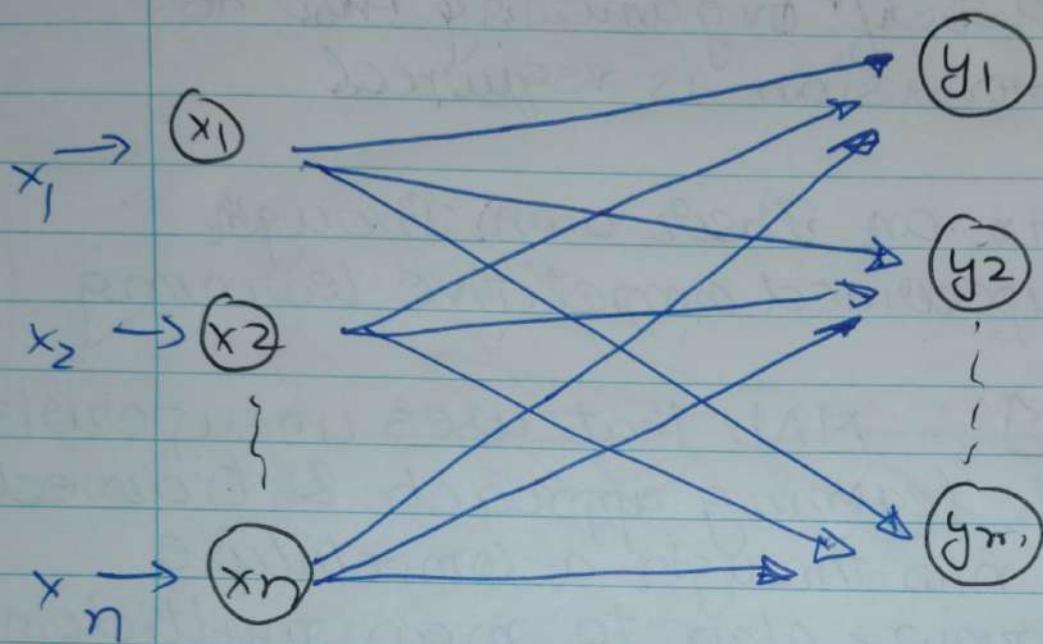
Multi Dimensional feature Map → Lower Dimensional feature Map

SOM layers \Rightarrow ① I/p layer
② O/p layer

If there n i/p's $\Rightarrow m$ o/p's

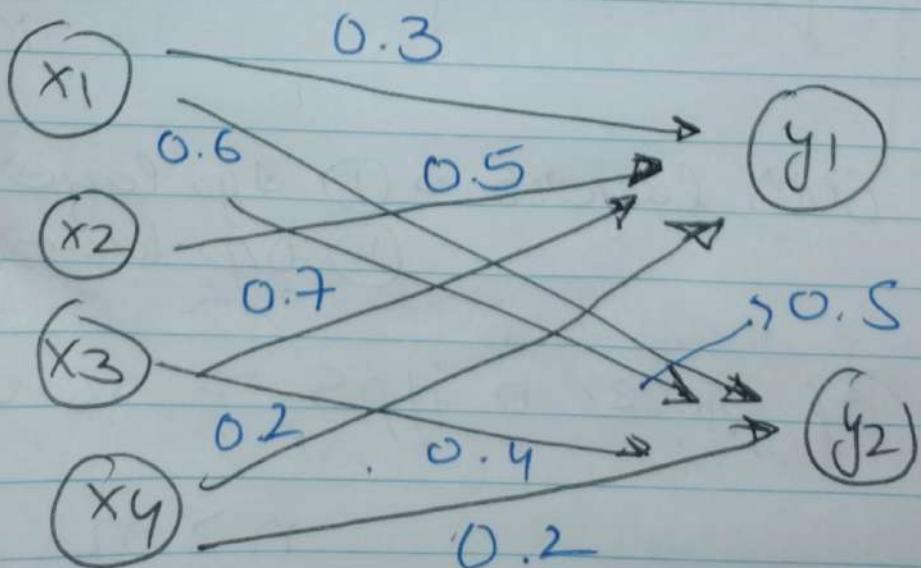
$$n > m$$

Architecture SOM (Fully Connected)



Q

Consider the network shown in fig which consider four training samples each vector of length 4 & two output units



→ Train SOM by determining the class membership of i/p data

$$x_1: (1, 0, 1, 0) \quad x_2: (1, 0, 0, 0)$$

$$x_3: (1, 1, 1, 1) \quad x_4: (0, 1, 1, 0)$$

Ans

→ Steps

↳ i) Calculate Euclidean Distance
one which have

less won ii) Update that particular weights which won

$$w_j(t+1) = w_j(t) + \eta(t)(x_s - w_j(t))$$

- Initial wt. Matrix

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.3 & 0.5 & 0.7 & 0.2 \\ 0.6 & 0.3 & 0.4 & 0.3 \end{bmatrix}$$

1st Iteration

Training Sample : $x_1 : (1, 0, 1, 0)$

→ Compute Euclidian dis b/w x_1 & unit 1 weights

$$\begin{aligned} d^2 &= (0.3 - 1)^2 + (0 - 0.5)^2 + (1 - 0.7)^2 \\ &\quad + (0 - 0.2)^2 \\ &= 0.87 \end{aligned}$$

→ lky for x_1 & unit 2 wts.

$$\begin{aligned} d^2 &= (0.6 - 1)^2 + (0.8 - 0)^2 \\ &\quad + (0.4 - 1)^2 + (0.3)^2 \end{aligned}$$

$$d^2 = 1.1$$

★ Unit 1 wins.

→ Update Unit 1 weights

$$\omega(t+1) = \omega(t) + \eta(t) (x_s - \omega(t))$$

$$\omega(t+1) = [0.3 \ 0.5 \ 0.7 \ 0.2]$$

$$+ 0.6 \begin{bmatrix} (1, 0, 1, 0) \\ (0.3, 0.5, 0.7, 0.2) \end{bmatrix}$$

$$\omega(t+1) = [0.72 \ 0.2 \ 0.88 \ 0.08]$$

$$\begin{bmatrix} \text{unit 1} \\ \text{unit 2} \end{bmatrix} = \begin{bmatrix} 0.72 & 0.2 & 0.88 & 0.08 \\ 0.6 & 0.7 & 0.4 & 0.3 \end{bmatrix}$$

Iteration 2

$$x_2 = [1, 0, 0, 0]$$

→ From Unit 1

$$d^2 = (1 - 0.72)^2 + (0 - 0.2)^2 + (0.88)^2 + (0.08)^2$$

$$d^2 = 0.0784 + 0.04 + 0.774 + 0.0064$$

$$= 0.8992$$

by x_2 : from Unit 2

$$d^2 = (0.6 - 1)^2 + 0.7^2 + 0.4^2 + 0.3^2$$

$$d^2 = 0.9$$

Unit 1 Wins

→ Update wt. of Unit 1

$$\omega_{t+1} = \omega_t - n(x - \omega_t)$$

$$\omega_{t+1} = [0.72 \ 0.2 \ 0.88 \ 0.08]$$

$$- 0.6 [1, 0, 0, 0]$$

$$- [0.72 \ 0.2 \ 0.88 \\ 0.08]$$

$$\omega_{t+1} = [0.72 \ 0.2 \ 0.88 \ 0.08]$$

$$- 0.6 [0.28 \ -0.2 \ -0.88 \\ -0.8]$$

$$\omega_t = [0.72 \ 0.2 \ 0.88 \ 0.8]$$

$$t \begin{bmatrix} -0.168 & 0.12 & 0.528 \\ 0.48 \end{bmatrix}$$

$$\omega_{t+1} = [-0.96 \ 0.32 \ 1.408 \ 1.28]$$

By calculate 3
by iteration

BOLTZMANN MACHINES

Date: _____

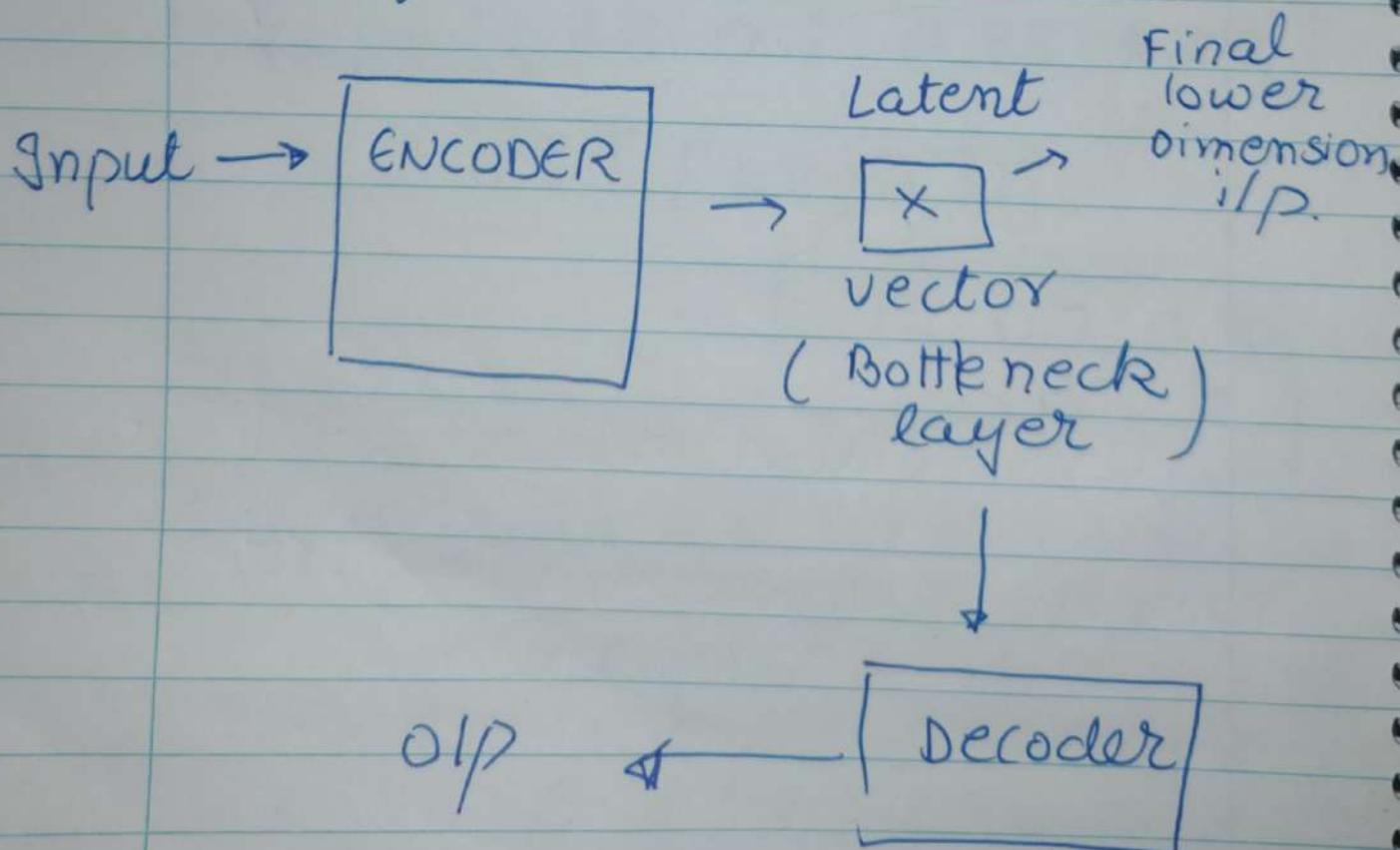
↳ Used to Solve
Two Problems

- ① Search Problem
- ② Learning Problem

- Restricted Boltzmann Machines
DON'T DETECT → THEY RECONSTRUCT

← AUTO ENCODER

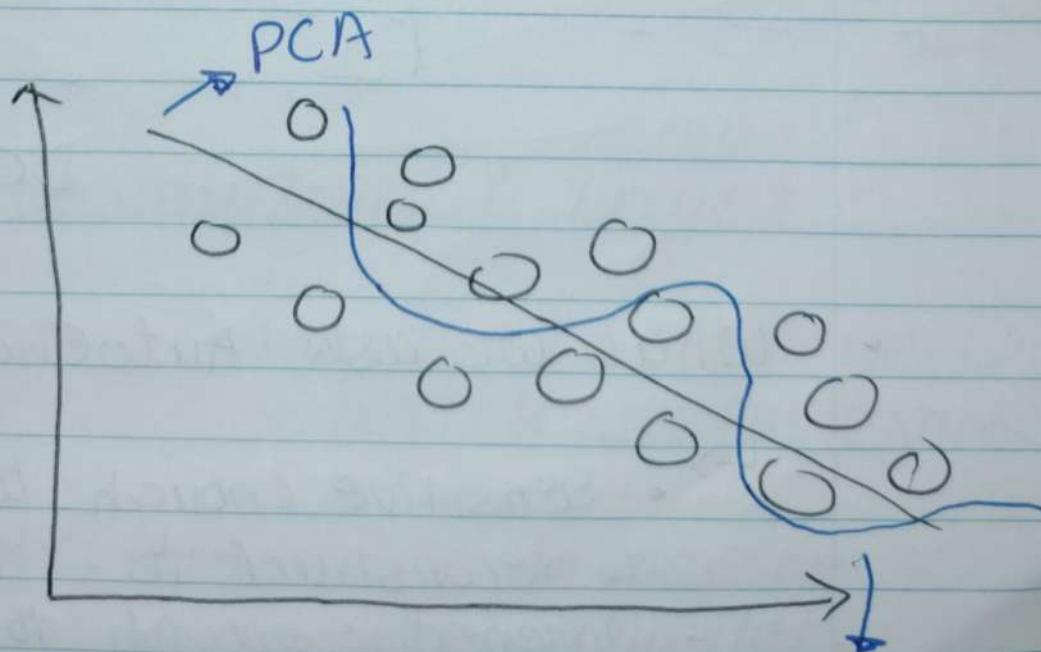
↳ is an unsupervised learning
technique



- Uses of AutoEncoder: Anomaly Detection
- Training Objective: To minimize Reconstruction error: error b/w o/p & i/p.

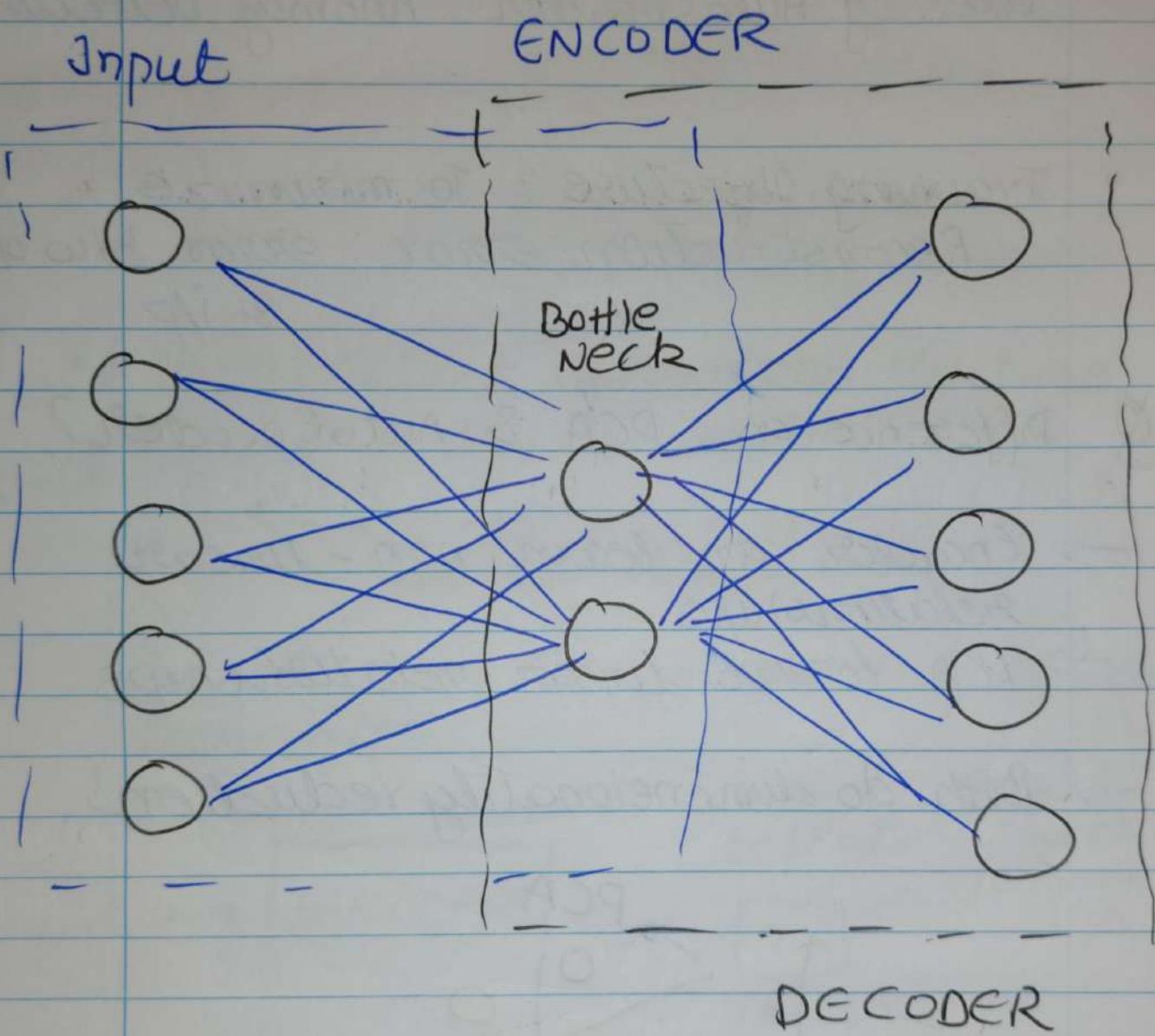
Q Difference in PCA & AutoEncoder?

- Encoder can learn non-linear relationships
- PCA learns linear relationships
- Both do dimensionality reduction



AutoEncoder

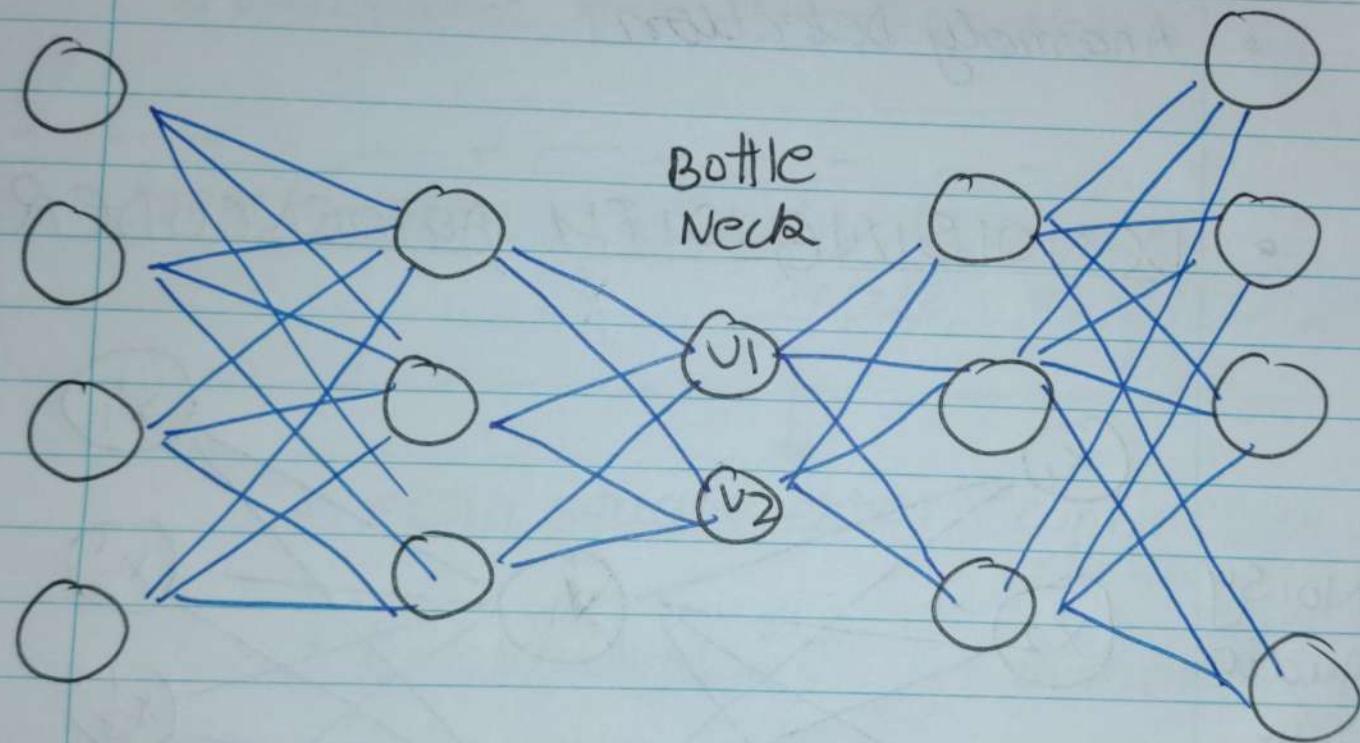
ARCHITECTURE



- What we ask AutoEncoder
 - ↳ • Sensitive enough to i/p data to reconstruct it.
 - Insensitive enough to i/p data not to overfit it

$$\rightarrow E(x, \hat{x}) + \text{Regularization}$$

→ HOW AUTOENCODER ACTUALLY LOOKS



Deep Convolutional Layer

- ENCODER : CONVOLUTION + LEAKY RELU + BATCH NORMALIZATION
- DECODER : Convolution transpose + Leaky Relu + Batch Normalization

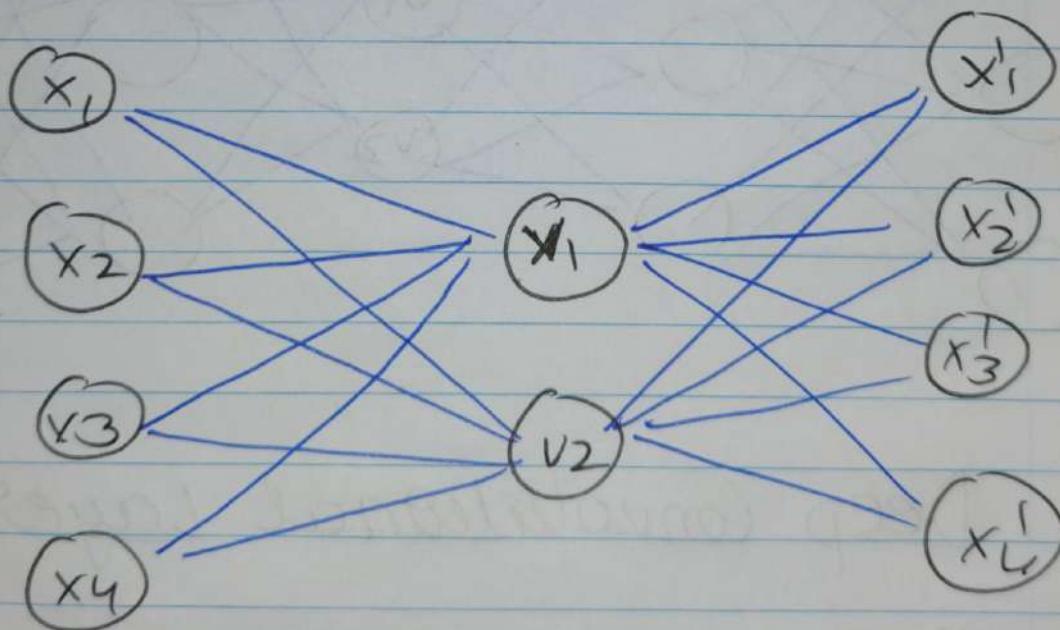
AutoEncoder Applications

Generation

Denoising \Rightarrow Images, Data

Anomaly detection

DENOISING WITH AUTOENCODER



You create pairs noisy & clean data

Noisy given as input to the autoencoder

- clean version is used to minimize error (reconstruction)
- Model's loss fn compares reconstructed audio to the clean audio & calculate reconstruction error.

• loss function : reconstruction loss :

input - reconstructed

↓
old

- mean square error : continuous data
- Binary Cross Entropy : Binary data

Q★

AutoEncoder is not a generative Model but Variational AutoEncoder is a generative model?

→

AutoEncoder lacks the ability to generate new, diverse sample.

→ Variational AutoEncoder can generate new sample.

AutoEncoder

- I) Latent Representation
 - Deterministic
- II) Not Capable of New Data Generation
- III) Function: Reconstruction
- IV) Loss fn: mse or bce

Variational AutoEncoder

- I) Latent Representation → Probabilistic
- II) Not Capable of New Data Generation
- III) Function: Generation + Reconstruction
- IV) Loss fn: mse / bce + KL Divergence

Tips for AutoEncoder

- ↳ i) keep no. of neurons less in hidden layer

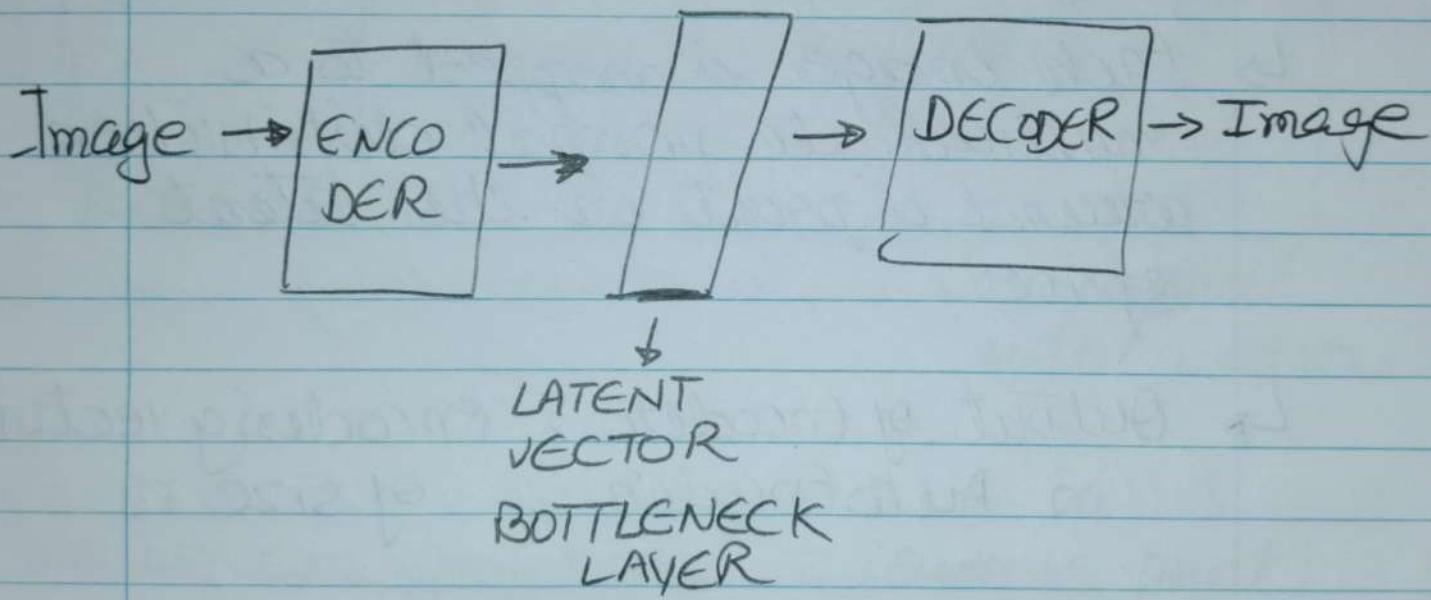
15

Variational ENCODER

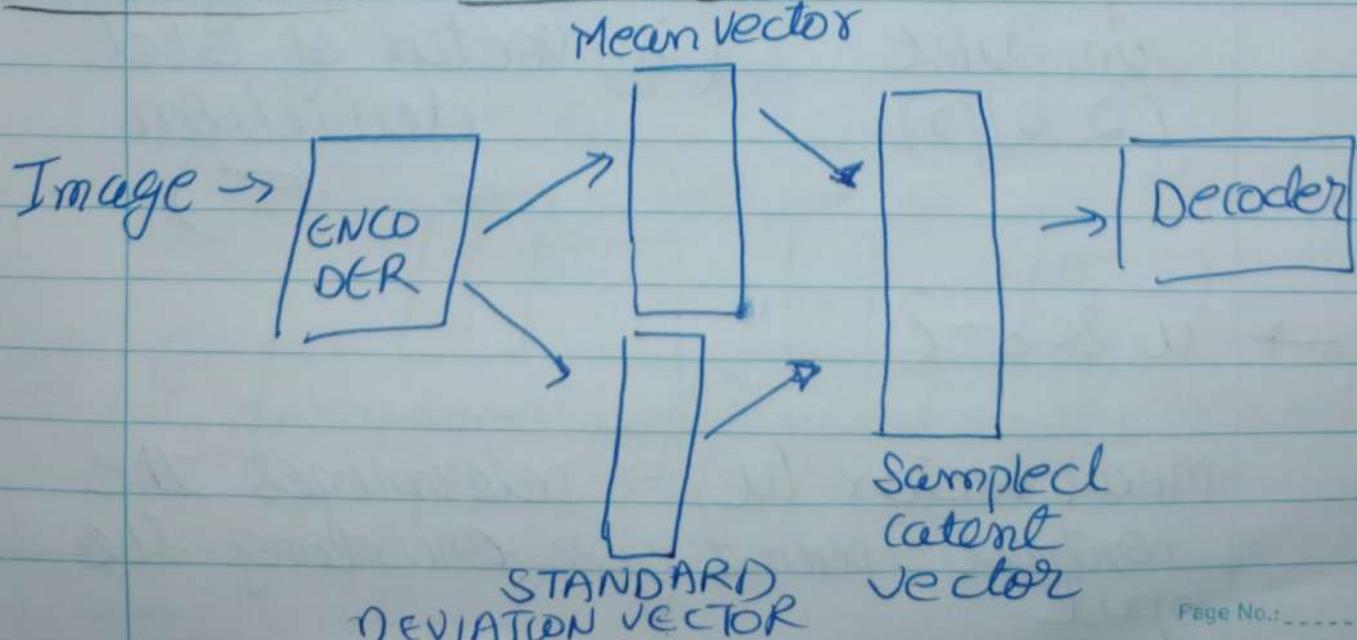
Date: _____

- * If you are working with images \Rightarrow You use convolutional layer.

AUTOENCODER



Variational Autoencoder



Problem in AutoEncoder

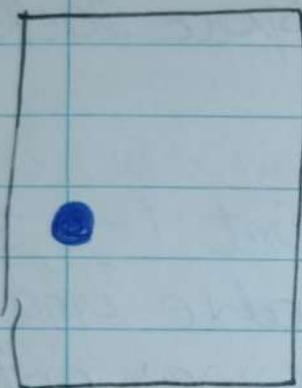
- ↳ Each instance is mapped directly to one point in latent space.

Variational AutoEncoder \Rightarrow Solver

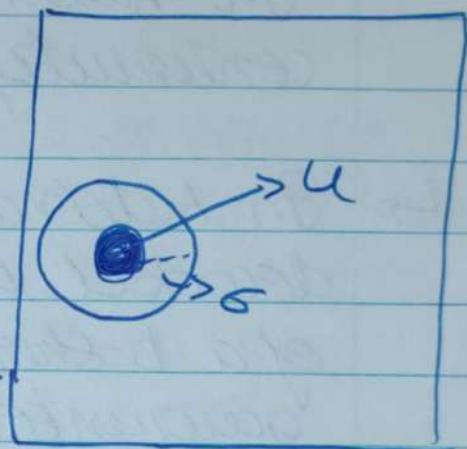
- ↳ Each image is mapped to a multivariate normal distribution around a point in the latent space.
- ↳ Output of encoder : encoding vector in AutoEncoder of size n
- ↳ output of Encoder. (i) vector of mean in VAE (2 O/P) (ii) vector of std deviation (σ)
 $\rightarrow u \& \sigma?$

Mean vector (u) \rightarrow determines the central point for encoding the input.

σ : How much the encoding can deviate from central point



AutoEncoder

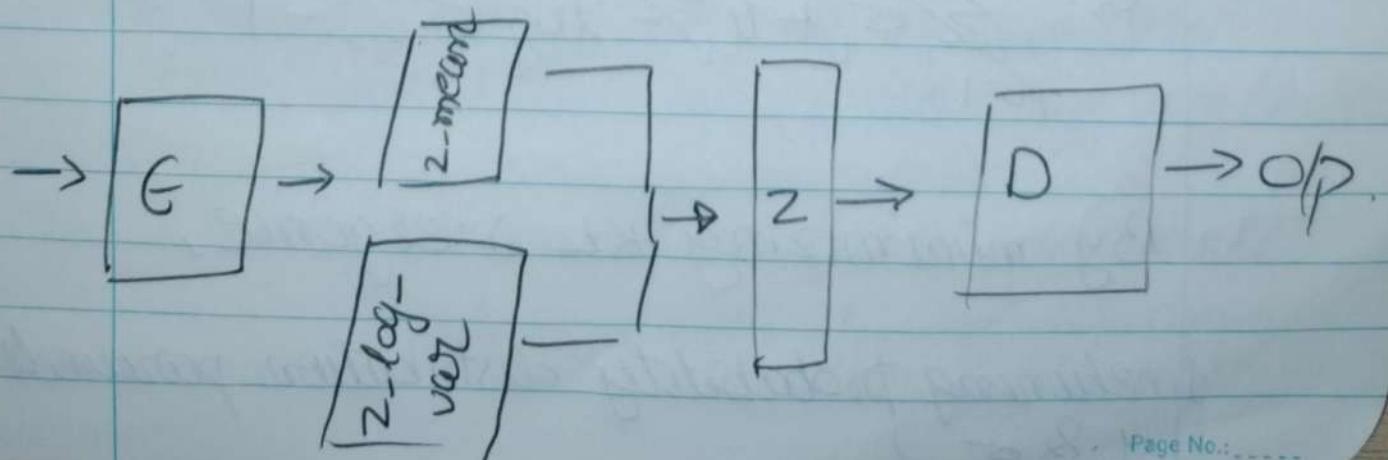


Variational
AutoEncoder

→ we can generate z (sample point)

$$\hookrightarrow z = z\text{-mean} + z\text{-sigma} * \text{epsilon}$$

$$z\text{-sigma} = e^{(z\text{-log-var} * 0.5)}$$



Q How this changes in encoder in VAE helps?

↳ In Autoencoder, latent space lacks continuity

↳ In Autoencoder, if the point (-5, 3) decoded into a recognizable image of a t-shirt, then there was not guarantee that (-5.1, 3.1) would produce something similar

→ But in VAE, decoder ensure nearby yield similar decoded images

KL Divergence

$$\hookrightarrow \sum_{i=1}^n (\sigma_i^2 + u_i^2 - \log(\sigma_i) - 1)$$

↳ By minimizing KL divergence,

finetuning probability distribution parameter (u & σ)

Uses Variational AutoEncoder

- ↳ i) Generating Purely synthetic music
- ii) Generating fake human faces.

→ DL-Gate → Applications of encoder

↳ Imp to learn concept

REPARAMETERIZATION

TRICK

- ↳ In VAE, if we try to take sample from a probability distribution to feed it to decoder it is a random process & the process of generating sample from distribution is non-differentiable
- ↳ So, we can't compute gradients
- ↳ So, can't do Backprop.

Q Now, Why this process of random picking sample is non-differentiable?

↳ ∵ it don't involve depend on neural n/w's learned weights.

Q ONE MORE QUESTION \Rightarrow HOW TO USE VAE FOR NEW IMAGE GENERATION?

↳ So, for Generation of new images from VAE, only involve decoder part no role of encoder X

★ Reparameterization trick is used to make sampling process differentiable.

Q How Reparameterization trick make sampling process differentiable

By Expressing Random sample as a deterministic fxn of model's parameter and fixed noise term

$$z = u + \sigma(\epsilon) \rightarrow \begin{array}{l} \text{sample point} \\ \text{from standard} \\ \text{normal distribution} \end{array}$$

\hookrightarrow from our distribution

$u, \sigma \Rightarrow$ o/p/s of the encoder

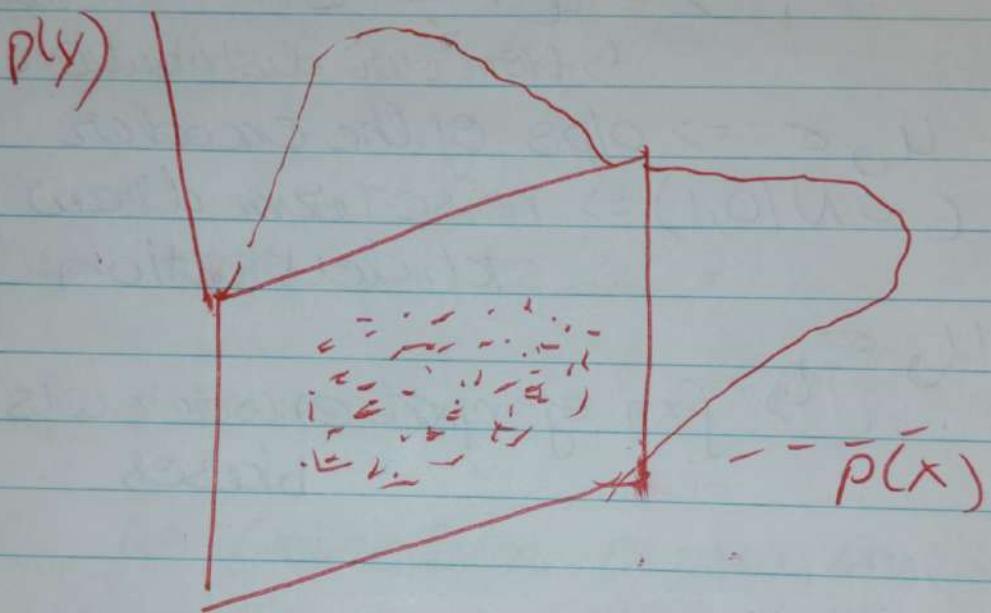
$\epsilon \sim N(0,1) \Rightarrow$ noise term draw from std. distribution

$u, \sigma \hookrightarrow$ fxn of input encoder wts. & biases

Q HOW TO GO FROM VANILLA AUTO ENCODER TO VARIATIONAL AUTO ENCODER

- ↳ • Modify Encoder Component
- ↳ • Modify loss fxn

Visualise Multivariate Normal Distribution



$$z = u + \sigma \cdot \epsilon$$

$$\sigma = e^{\frac{\log(\sigma^2)}{2}}$$

- Modify Loss Function

$$\text{LOSS} = \text{RMSE} + KL$$

↓
Root mse

KL \Rightarrow difference b/w normal distribution (mean, log variance vector) from standard normal distribution

$$D_{KL} = \frac{1}{2} \sum (1 + \log(\sigma^2) - \mu^2 - \sigma^2)$$

* What does this KL loss term do?

\rightarrow Penalizes observation where mean & log variance vectors differ significantly from parameters of Standard Normal Distribution

* KL Divergence Loss: Promotes Symmetric Term

around origin

- Decrease chances of large gaps b/w clusters of points

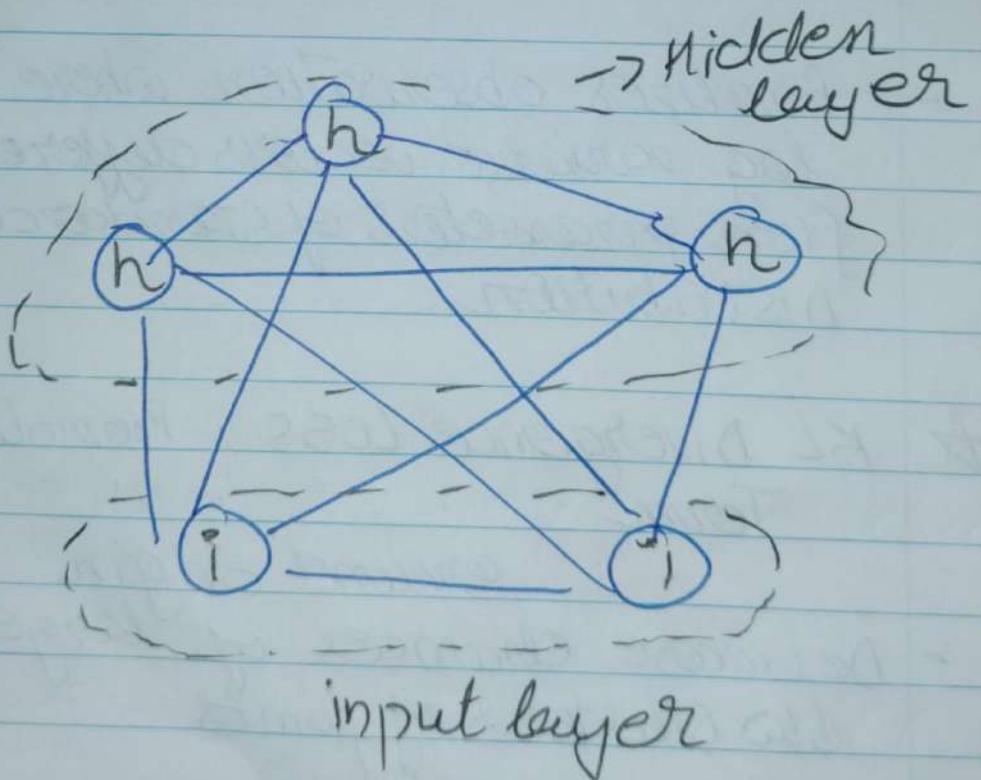
BOLTZMANN MACHINE

Date: _____

↳ Unsupervised

Network of symmetrically connected neuron-like units that make stochastic decision about whether to be on or off

- Fully connected
- No data flow direction



- No output layer
- input units are also interconnected

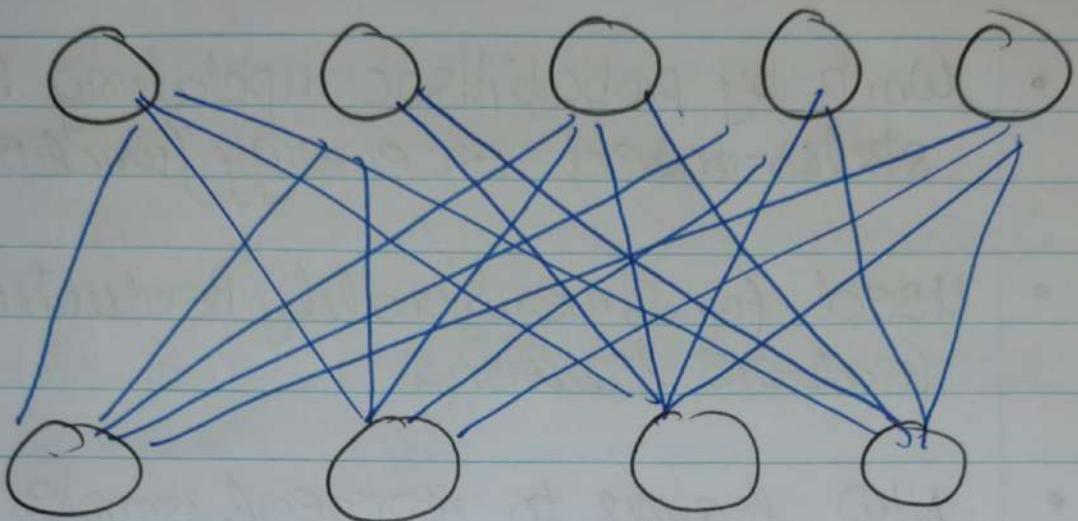
- Work by probabilistic updating their states based on energy function
- Used for Dimensionality Reduction & feature learning
- Now learns to represent complex prob. distribution
- During Training, optimize weight by minimizing loss function \Rightarrow Energy of the system

RESTRICTED BOLTZMANN

↳ i/p - i/p & i/p - hidden connection was removed. \rightarrow hidden - hidden

Q Why shifted to Restricted Boltzmann?

- ↳ Reduced complexity of model \Rightarrow R.B
- ii) Slow convergence & high computational overhead \Rightarrow B

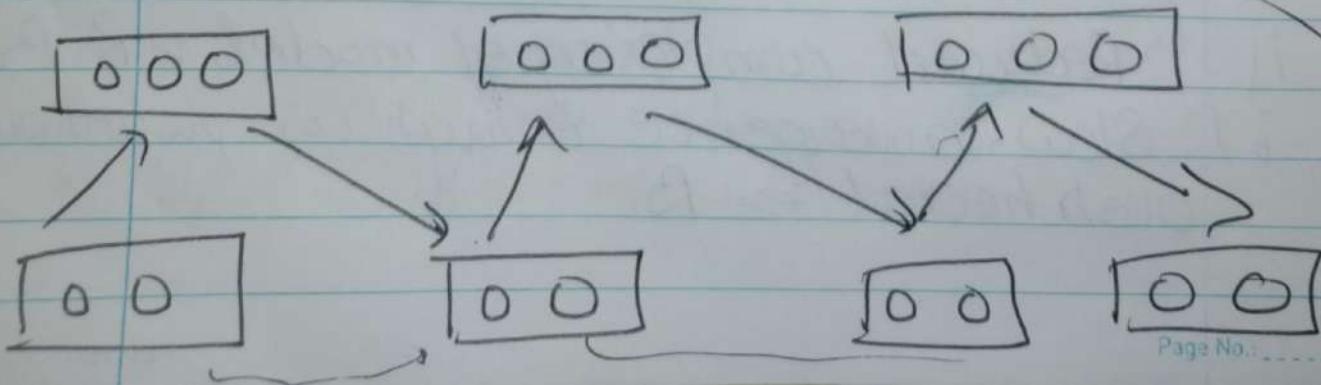


i/p
(visible)
node

RB → Generative Stochastic ANN that can learn probability distribution over its set of inputs

* Can be used in Movie Recommendation System

GIBBS SAMPLING & CONTRASTIVE DIVERGENCE



WOW
Just Amazing

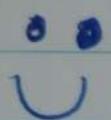
2014 (Introduced)

apsara

Date: _____

GAN

↳ Generative Adversarial Network



Text to Image Generation

WaterMark Remover

Low Resolution to High Resolution

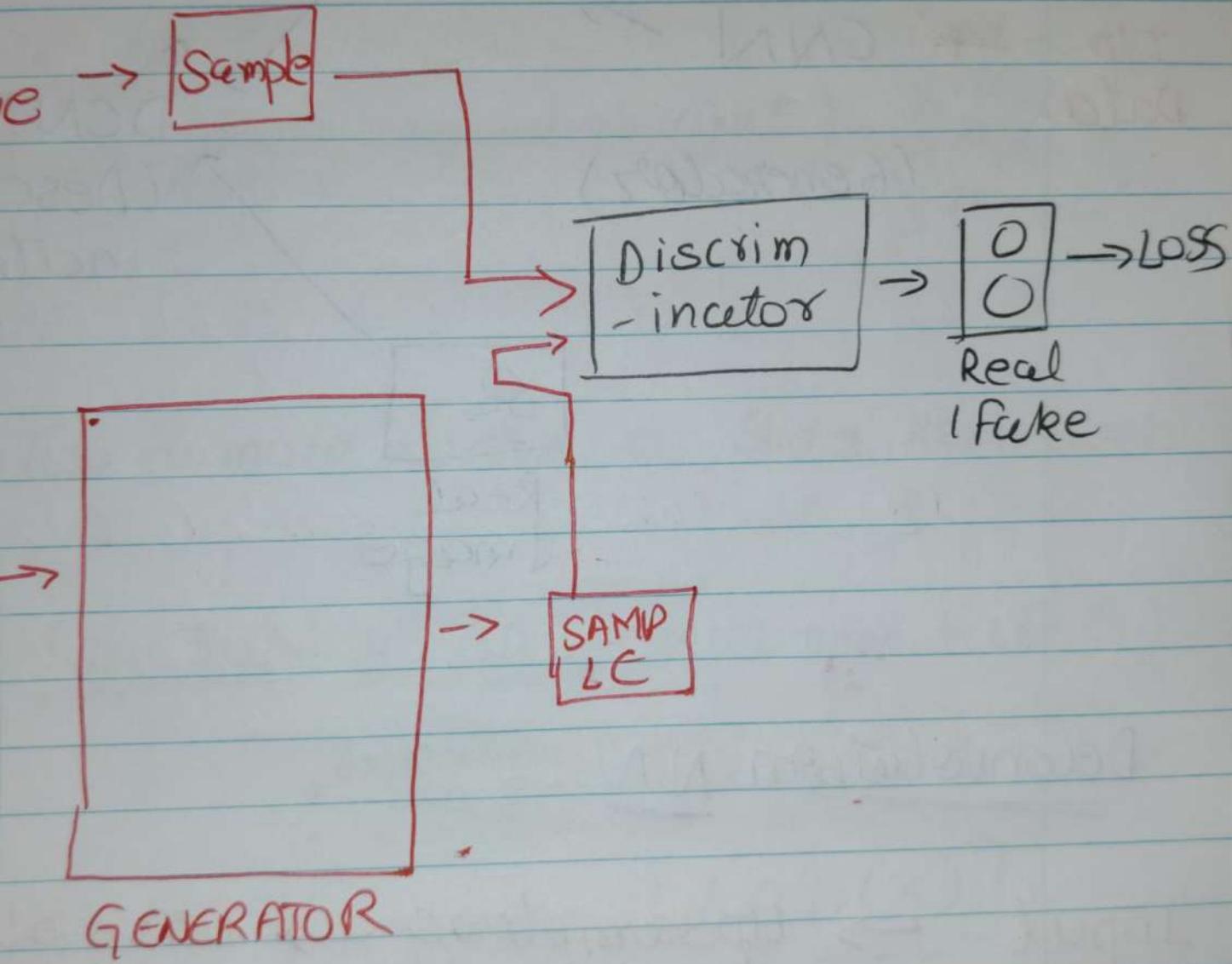
Lip Lip-Syncing



AFTER MODEL IS TRAINED, we
use only Generator to Generate
Images

• Generator: ConvO + Batch Norm
Layer +
ReLU

• DISCRIMINATOR: Takes TWO i/P
i) Generator o/P.
ii) Real image

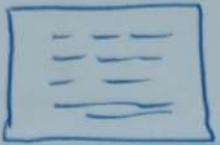


GAN

I/P
Data

→ CNN →

(Generator)



apsara
deconvolutional

↑
DCNN

(Discrim-
inator)



Real
Image

DeConvolution NN

Input → upsampling → DeConvolutional layer
(Flatten)
(Opposite of
Down Sampling)
Max Pooling

(Feature
Generation
(Image
Reconstruction))

O/P ← Batch
Normalization
(optional) ↓
Activation
Function

Date: _____

1st \Rightarrow 02nd \Rightarrow 03rd \Rightarrow 04th \Rightarrow 15th \Rightarrow 0

!

$\rightarrow 0/1$ (For many we put as input)
 images

- Two formula written in Slides Remember it.

LOSS FUNCTION OF GAN: $\min_G \max_D V(D, G)$

\uparrow expectation of discriminator when
 we give real images

$$V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)]$$

$$+ \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(h(z)))]$$

\hookrightarrow expectation of discriminator
 when we give fake images

G = Generator

D = Discriminator \rightarrow probability

$p_{\text{data}}(x)$ = distribution of real data

$p(z)$ = Distribution of generator.

$D(x)$ = Discriminator $N(x)$

$G(z) \Rightarrow$ Generator $N(\omega)$

g_t can also be written as

$$L_D = -[\log(D(x)) + \log(1 - D(G(z)))]$$

o/p
of
Discriminators

$D(x) \Rightarrow$ o/p of real data by Discriminator
 $G(z) \Rightarrow$ fake data generated by generator
 $D(G(z)) \Rightarrow$ Discriminator o/p for fake data

Real Data \rightarrow Discriminator $\rightarrow D(x)$

Fake Data \rightarrow " $\rightarrow D(G(z))$

O/p $\Rightarrow 0$ or 1

\hookrightarrow we want to train our model in such a way that it should

(from Discriminator) $D(G(z)) \Rightarrow$ o/p for fake data $\Rightarrow 0$
 $D(x) \Rightarrow$ o/p for Real data $\Rightarrow 1$

→ It will train the model perfectly by making contest

DC-GAN → Deep Convolutional GAN

- ~ It uses Convolutional DC
- ~ Simple GAN: Use convolutional layer (Deconvolution layer) to generate images
- ~ Simple GAN: Use fully connected Dense Layer

GAN

DCNN

① Generic Arch with fully connected layers

① Arch. specially designed for image generation

② used for generating various type of data

② Generating Images

(iii) Can take any form of i/p data including vector, matrix

(iv) Fully Connected Layer

(v) Image Data

(vi) Stabilize Training by Batch Normalization, avoiding fully connected layer

Other Types of GAN

i) Conditional GAN :

ii) Super Resolution GAN : low Resolution
→ High Resolution

↳ DANN + Adversarial NN

iv) Info GAN :