

100 days ML by Anup

Batch vs Online ML

Batch

- only stored data is known
- hardware limitation and less lot of data

Online

- Real time - New data is trained.
- performance improved dynamically as the new data is added.

Tensors

Container for storing numbers

Scalar → 0D Tensor

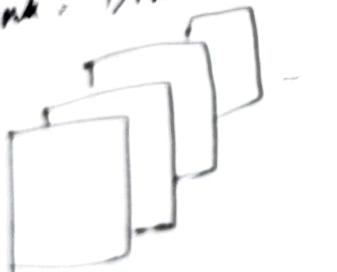
Vector → 1D Tensor

Cube → 2D Tensor

matrix → 2D Tensor

no of axis = Rank, Dimension

nD tensor =



One hot Encoding

Categorical Data

Ordinal

Nominal

one hot encoding
if there are many categories in the data how to convert that to numbers.

(1, 0, 0)	yellow
(0, 1, 0)	blue
(0, 0, 1)	Red
:	

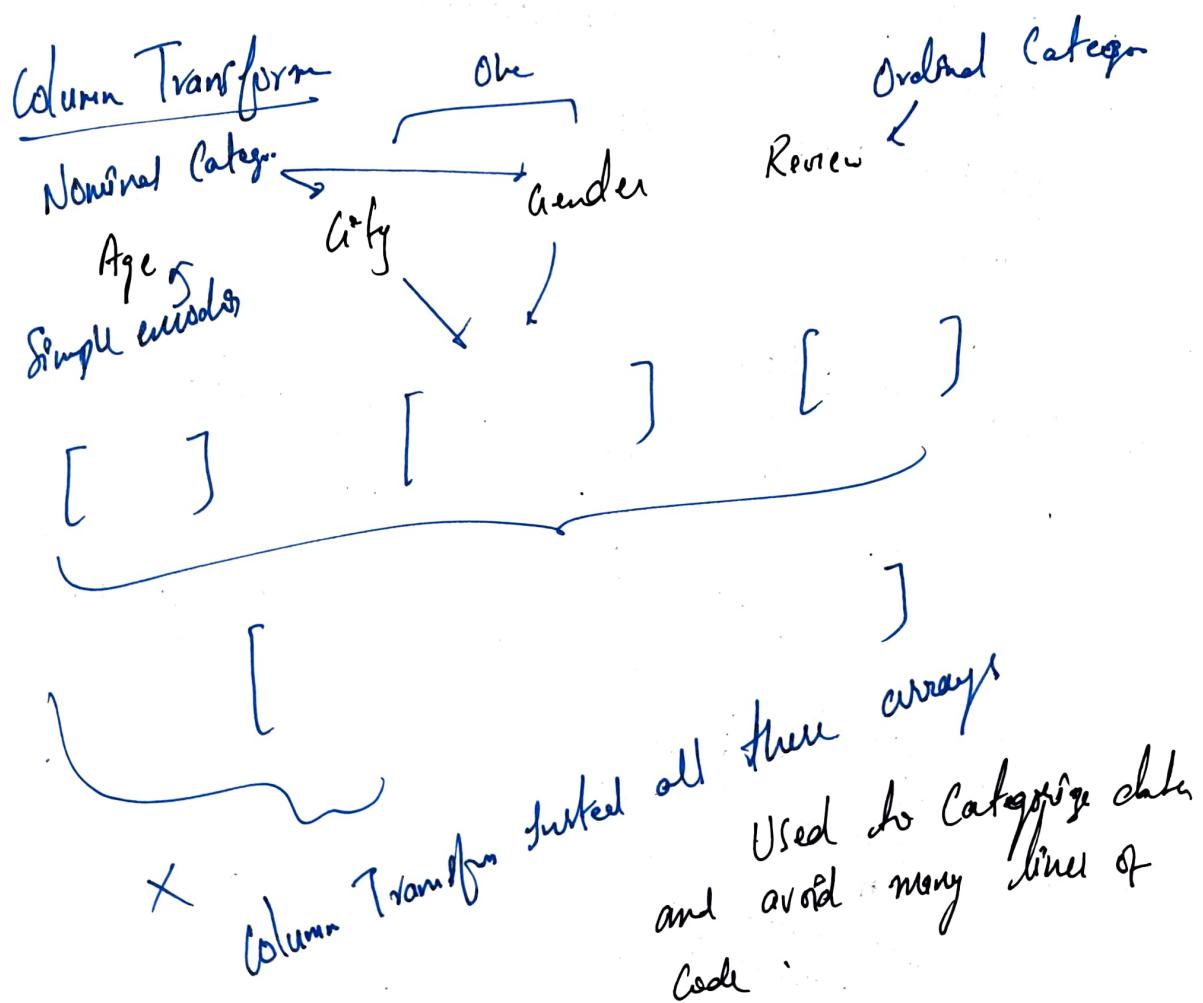
multicollinearity :-
 Can't let input vector to have
 multiplicative relation
 n categories ————— (n-1) columns.
 Columns.

pd.get_dummies(data, column)

One hot encode class in scikit learn → Used for ML

ohe = OneHotEncoder()

ohe.fit_transform(x-train["fuel", "owner"])

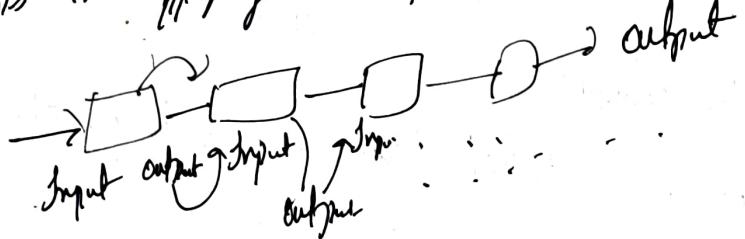


day 21

Pipelines

from scikit-learn

- Chains together multiple steps so that output of each step can be used as input of next step
- Helps in applying same preprocessing to train and test



Scoring function →

Decision tree

Pass list of tuples.

each tuple has list of transformers.

make-pipeline → formula carries

make-train

Binning and Binarization

converting Numerical data → Categorical Data

Discretisation (Binning)

Unsupervised	Supervised	Custom
→ Equal width		
→ Equal frequency	Decision tree ↳ binning	
→ k-means binning	Rare	

Binarization

Equal width binning / Uniform

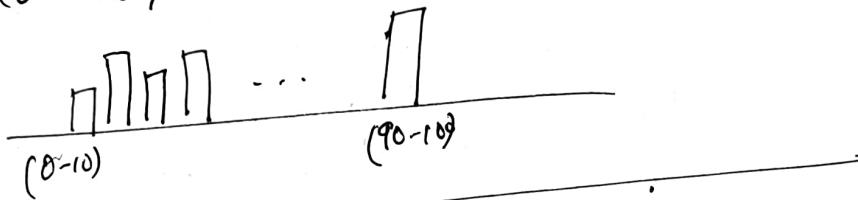
Bins = any number

let bins = 10

$$\frac{\max - \min}{\text{bins}} = \frac{100 - 0}{10} = 10$$

Range =

(0 - 10), (10 - 20), (20 - 30) ... (90 - 100)] 10 bins



Equal freq / Quantile Binning

Interval = any number

ext
0 - 10 1. 10 - 20 1. 20 1.
← 10 1.

70 - 100 Y.
100 Y.

Intervals → helps uniform value spread
→ widely used than equal width

K-means binning

- Used clustering algorithm. Forms clusters automatically
- Same algo is used.
- Used when clusters can be formed in the data
- here, Intervals = "Centroid".
- line that bisects clusters
- line that distance is measured from centroid
- every point distance is measured to that center to centroid, assigned to that cluster.

day 18

Data Gathering

CSV

JSON/SO

fetch API

Webscraping

CSV

comma separated value

a Pandas read_csv "html" → webpage documentation

dtypes parameter → overriding data types parameter in the file.
to convert float to int or anything like this
object = string

Dates → by default its strings →
should convert to date - type time
(datetime 64)

Keras functional Model

Non linear models

concatenating feature from different models (KNN, NN ... etc)

RNN

Recurrent Neural Network

used for sequential data and widely used in NLP

Type of sequential model.

example of sequential data → Text | Passage

→ Time series data

→ Audio | Speech

→ DNA Sequence

Why not NN?

- because input size cannot be varied in NN.

this can be solved but there still exists two problems
when the input size increase to thousands words,
million words ---
calculations will increase ---!

Need to do lot of unnecessary calculation

That's why a new architecture Only to predict Sequential data, they created RNN by Wilhelm Lenz and Ernst Ising

Road map →

Simple RNN - architecture

Backprop in RNN in time

Problems with RNN

LSTM

GRU

Types of RNN

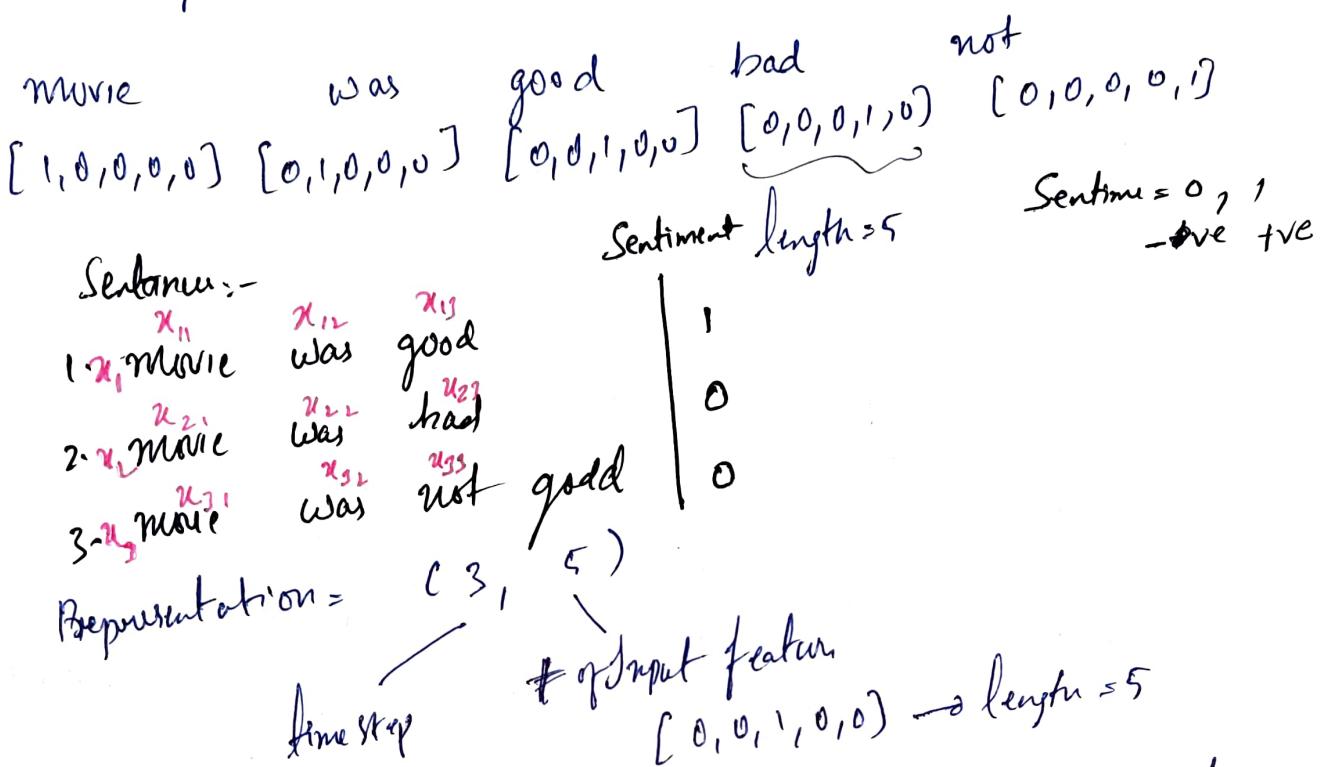
Deep RNN

Bi directional RNN
Sequence to Sequence model

Data in RNN :-

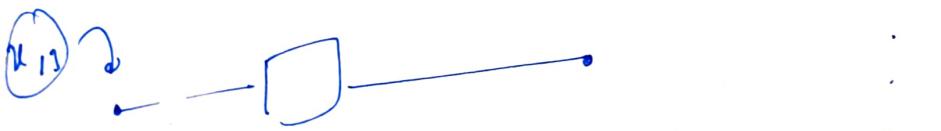
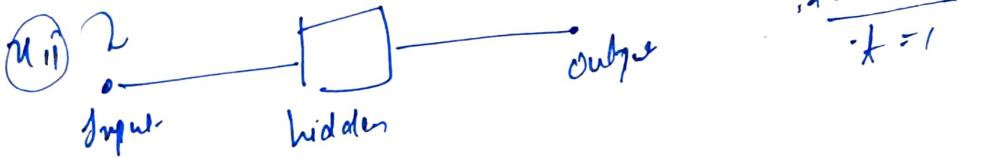
format = (timestep, input-features)

Can represent a sentence in One hot encoded format

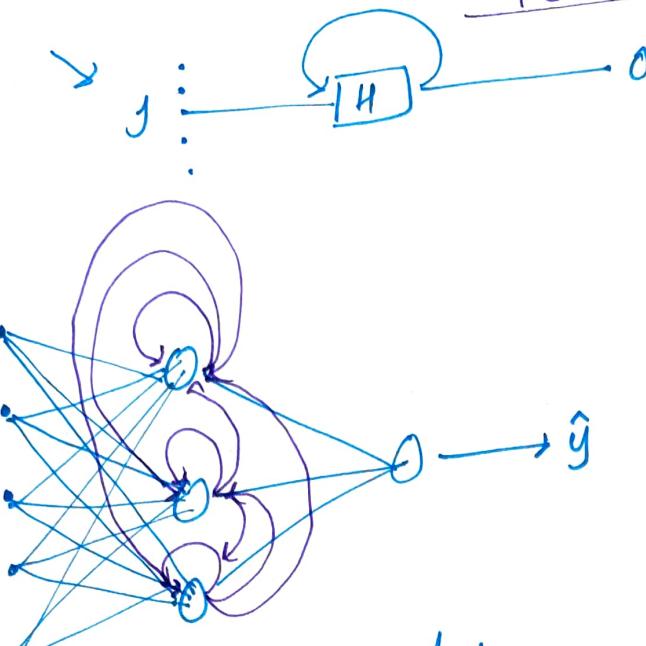


In Keras \rightarrow It has Simple RNN
format of Input data - Batch-size, Time step, Input-features
 $(3 \rightarrow 4, 5)$

Structure :-



there is a concept of State in RNN
 Feedback → Very Important

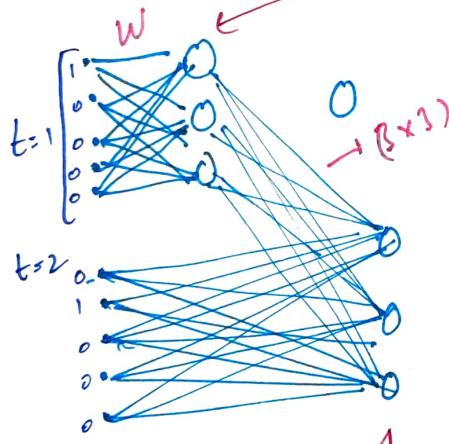


Input hidden

output
(Sigmoid)

better diagram :-

step 1 we send random values for feedback

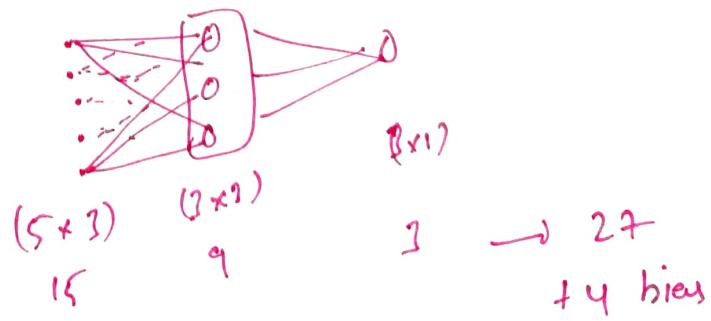


weight and Bias?

weights
27 + 4 Bias

$$W = 5 \times 3 = (5, 3)$$

$$\text{hidden} = (1, 3)$$



$1 \rightarrow 27$
+ 4 bias

Forward Pass \rightarrow Unfolding through time
Rows in Recurrent loop



$$t=1 \quad x_{ii} = (1, 0, 0, 0, 0)$$

dot product with weights.

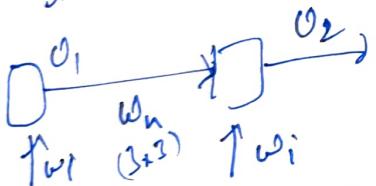
Inside Recurrent Network \rightarrow Simple NN function
Silent Activation - tanh

outputs

$$o_1 = \tanh(x_{ii} \cdot w_i) \quad \rightarrow \tanh(n) \text{ rule}$$

(1×3)
 (1×3)

$t=2$ Same weights are used

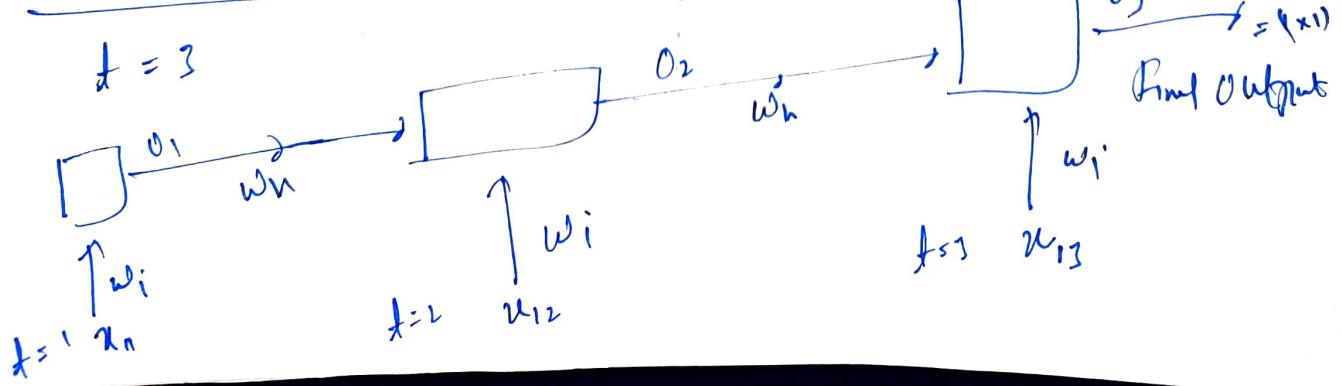


$$t_1 \quad x_{12} = (0, 1, 0, 0, 0)$$

$$x_{ii} \quad o_2 = \tanh[x_{12} w_i + o_1 w_n]$$

$(1 \times 5) \quad (5 \times 3) \quad (1 \times 1) \quad (1 \times 1)$
 $(1 \times 3) \quad (1 \times 3) \quad (1 \times 1)$
 (1×3)

$t=3$



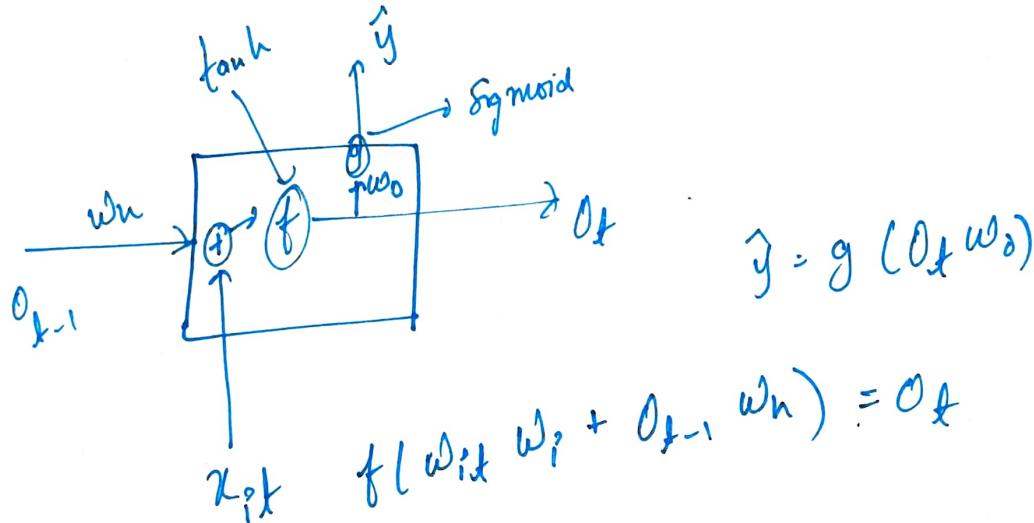
$$O_I = \tanh \underbrace{\left(\alpha_{13} w_i + \alpha_2 w_n \right)}_{(1 \times 3)}$$

except from $t=1$, rest rest all time step is having 2 inputs.
 [Words and output's from other steps] & Previous

$t=1$
 only 1 inputs
 So we give O_0 to keep consistency with random number or 0's

$$\text{so } f_i = \tanh (\alpha_{11} w_i + \alpha_0 w_n) \quad [\because O_0 = 0]$$

- Parameter and weights sharing is used here
- Simple RNN can predict sequence of last 10 time steps.



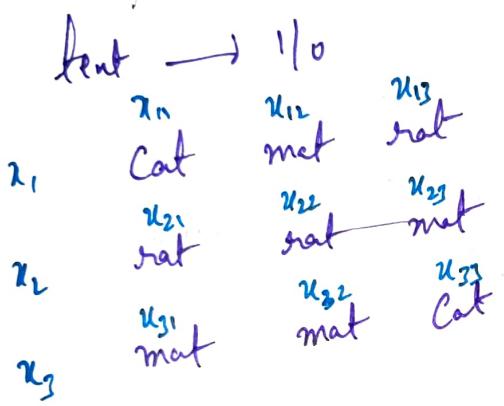
$$f_i \text{ with } w_i + O_{t-1} w_n = O_t$$

i : time #
 t : timestep

Embedding

does semantic meaning (words that are closer in vector space are expected to be closer in meaning)
 Reduces dimensions.
 (by avoiding Padding)

Backpropagation through time (BPTT)
 In Sentiment analysis model, many \rightarrow One row

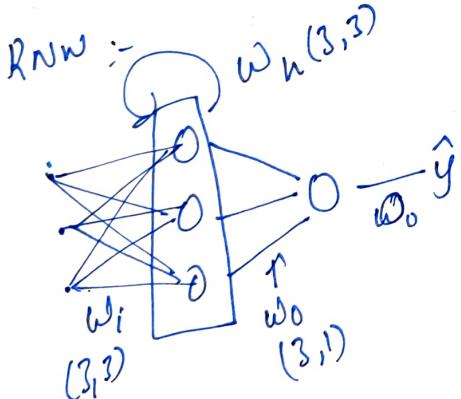


1	$[1, 0, 0]$	$[0, 1, 0]$	$[0, 0, 1]$
1	$[0, 0, 1]$	$[0, 0, 1]$	$[0, 1, 0]$
0	$[0, 1, 0]$	$[1, 0, 0]$	$[1, 0, 0]$

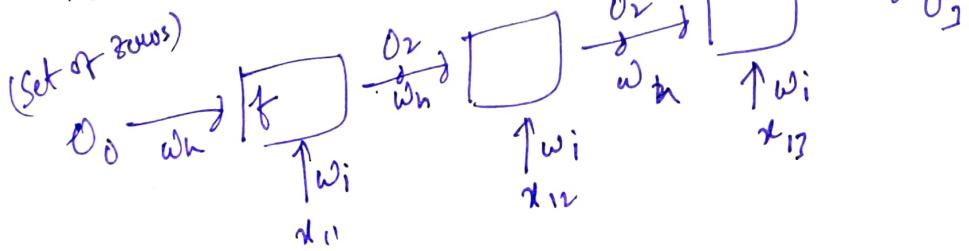
$$\text{vocab} = \begin{matrix} \text{cat} \\ [1, 0, 0] \end{matrix}$$

$$\begin{matrix} \text{mat} \\ [0, 1, 0] \end{matrix} \quad \begin{matrix} \text{Rat} \\ [0, 0, 1] \end{matrix}$$

$(3, 3)$ } dimension
 $(3, 1)$
 Ignoring bias here



Forward Prop:



$$q_1 \quad o_1 = f(x_i, w_i + \theta_0, w_n)$$

$$q_2 \quad o_2 = f(x_i, w_i + \theta_1, w_n)$$

$$q_3 \quad o_3 = f(x_i, w_i + \theta_2, w_n)$$

$$q_4 \quad \hat{y} = \alpha(o_3 \cdot w_0)$$

$$q_5 \quad L = -y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i) \quad \text{— loss fn}$$

Using gradient descent to find backward pass :-

$w_i \quad w_n \quad w_0$
 loss depends on \hat{y} , \hat{y} depends on
 o_3, w_0 . o_3 depends on previous term

$$\left| \begin{array}{l} w_i = w_i - n \frac{\partial L}{\partial w_i} \\ w_n = w_n - n \frac{\partial L}{\partial w_n} \\ w_0 = w_0 - n \frac{\partial L}{\partial w_0} \end{array} \right.$$

task is to find these 3 derivatives

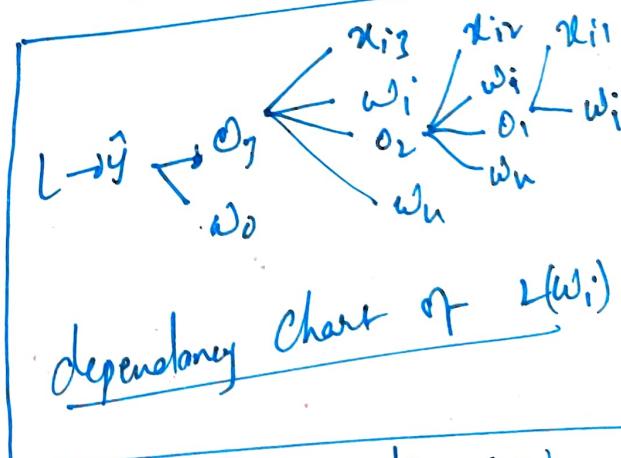
$$\boxed{\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_0}}$$

from q5 we can find $\frac{\partial \hat{y}}{\partial w_0}$

$$\boxed{\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_i} +}$$

$$+ \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_i}$$

$$+ \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_i}$$



dependency chart of $L(w_i)$

this is because the RNN unfolded 3 times (because there are 3 words in input)
 if there are n words then there will be ' n ' equations

n = time stamps

$$\frac{\partial L}{\partial w_i} \rightarrow \sum_{j=1}^3 \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_i}$$

$$\frac{\partial \hat{y}}{\partial o_1} = \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1}$$

$j=1$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_1} \quad \} \text{ for term 1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_1} \frac{\partial o_1}{\partial w_2} \quad \} \text{ for term 2}$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_3} \quad \} \text{ for term 3}$$

$$\frac{\partial L}{\partial w_i} = (j=1) + (j=2) + (j=3)$$

both are equal

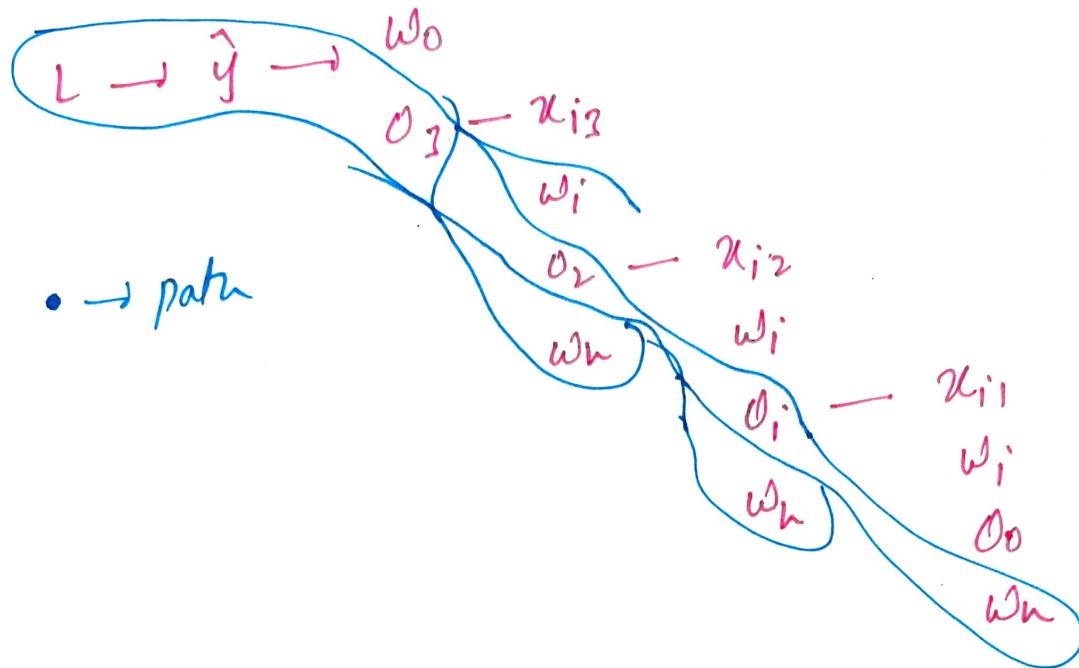
$$\frac{\partial L}{\partial w_n} = ?$$

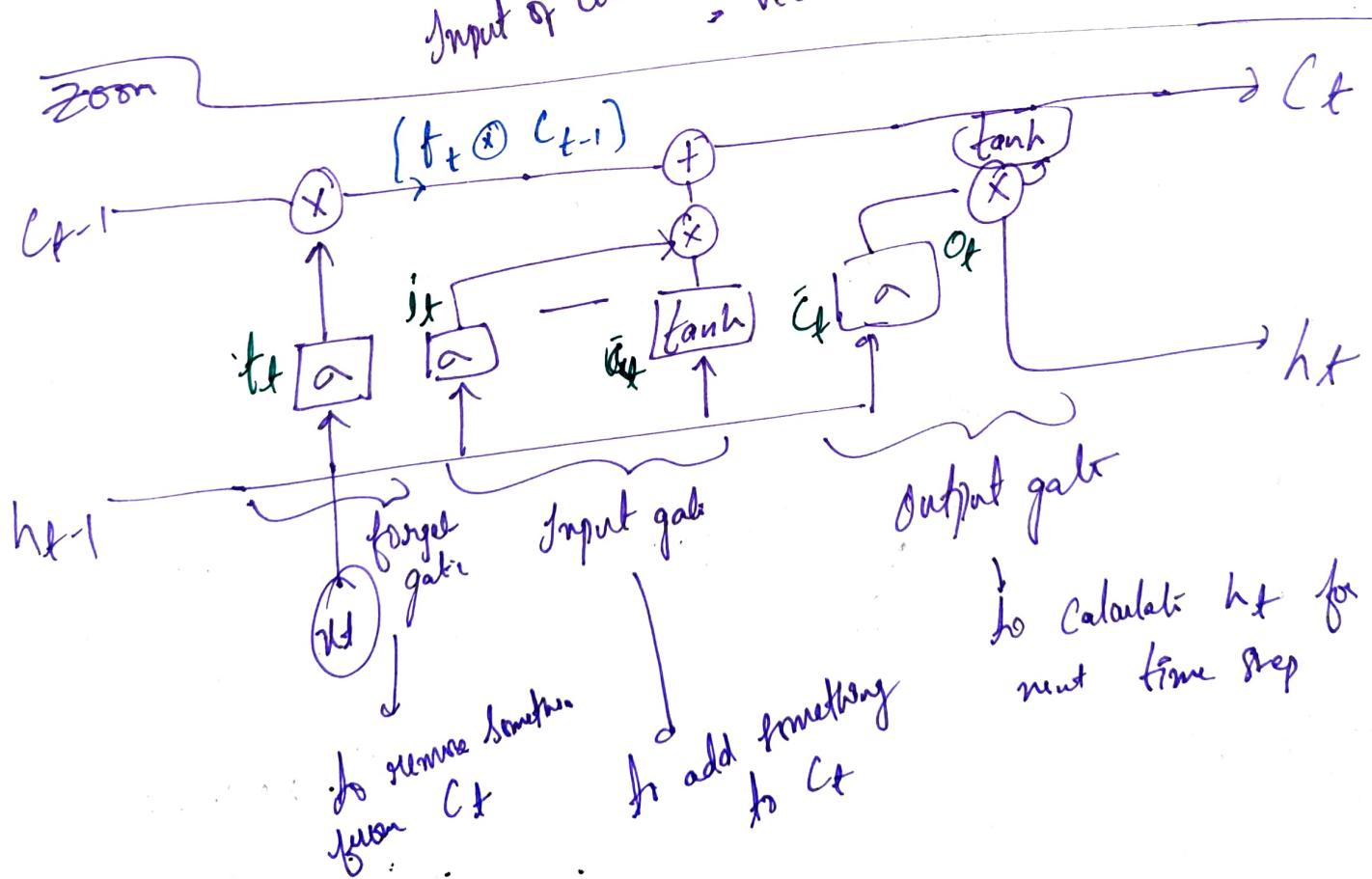
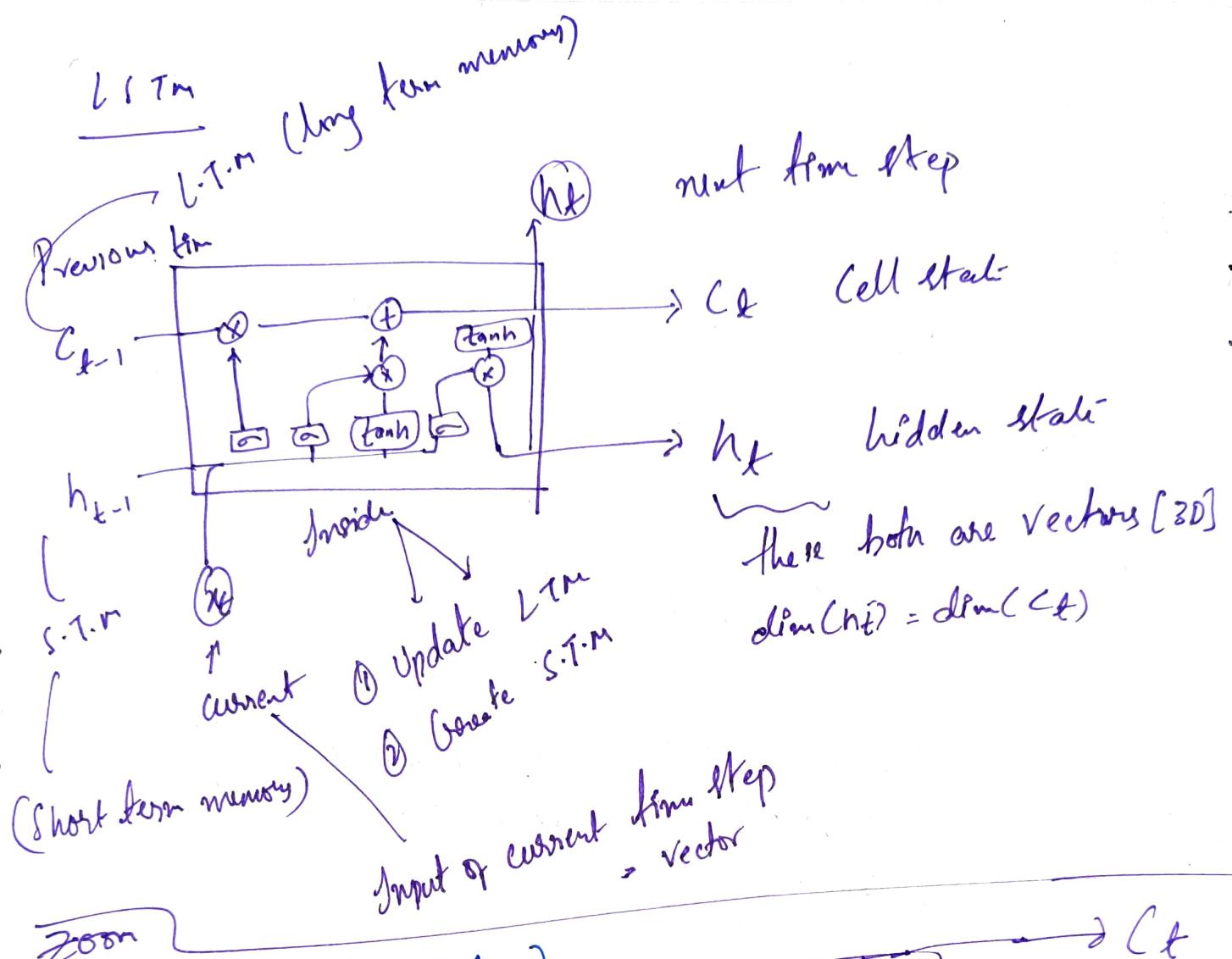
$$\frac{\partial L}{\partial w_n} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_n} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial w_n} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial w_n}$$

$$\frac{\partial L}{\partial w_n} = \sum_{j=1}^3 \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_n}$$

this is formed from
the dependency chart for
that w_n drawn in last paper.
in next paper all

let's do the dependency chart for w_n



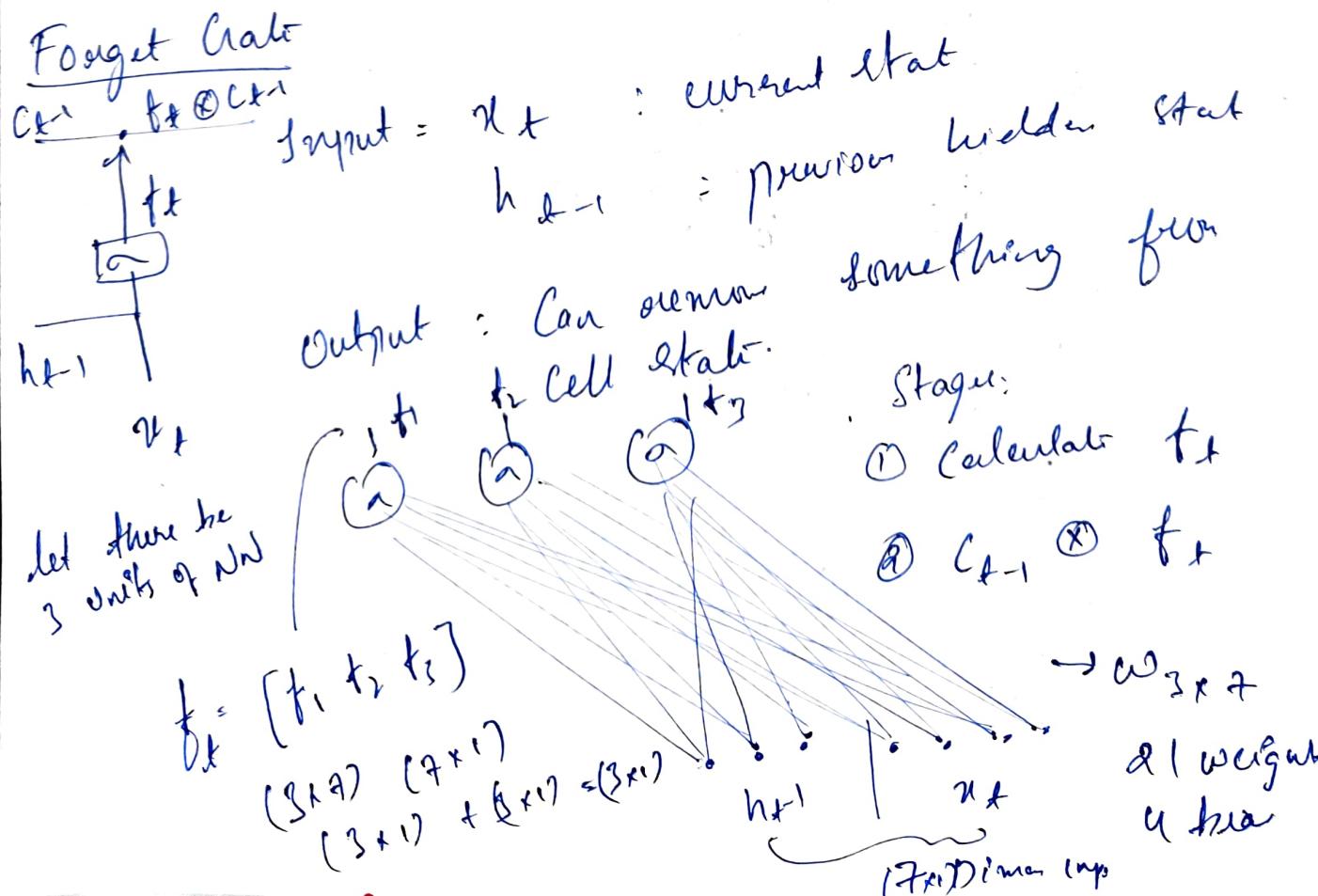


the gates:- } Forget Gate
 ② Input
 ① Forget
 ③ Output

multihot moded.
 forget it. \tilde{C}_t candidate cell state
 forget Input O_t output } vector
 $\dim(f_t) = \dim(\tilde{C}_t) = \dim(O_t) = \dim(C_t)$

Pointwise operation $\rightarrow \oplus, \otimes, \text{tanh}$
 \hookrightarrow output = vector

a, tanh : neural networks collection of multiple nodes with sigmoid and tanh activations.
 Sigmoid:



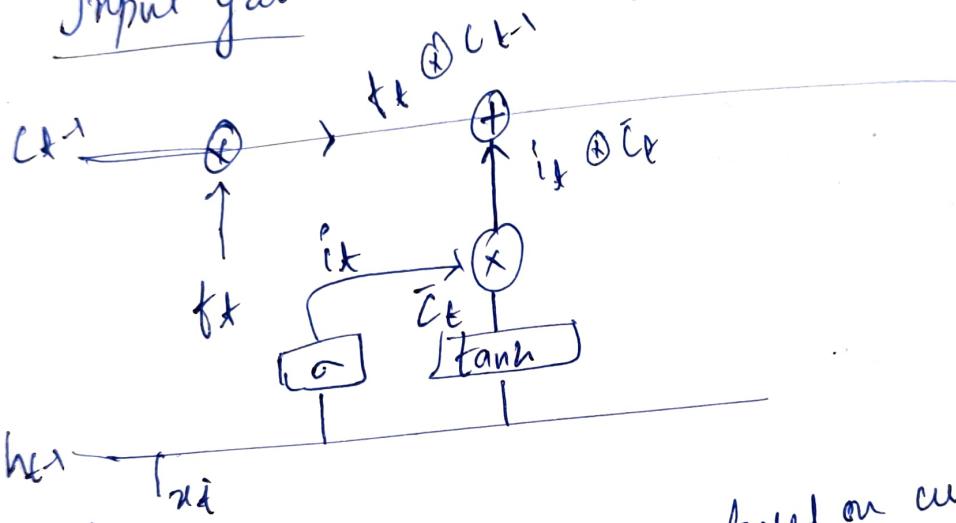
$$f_t = \sigma \left[w_f [h_{t-1}, x_t] + b_f \right]$$

$(3 \times 1) \rightarrow (3 \times 1)$

$[h_{t-1}, x_t] \xrightarrow{1 \times 2} \text{concatenation of 2 vector}$

$$f_t \odot c_{t-1} \xrightarrow{3 \times 1}$$

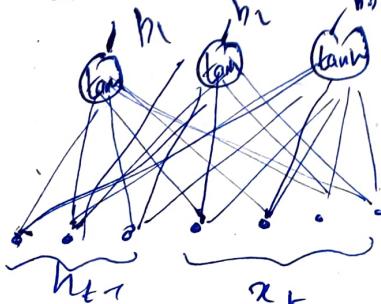
Input gate adds new info to the data



Stage 5

- ① \tilde{c}_t : Candidate cell state \rightarrow based on current input and previous cell state we add new info that's needed to be added
- ② i_t : decide what item needed to be added to final cell state.
- ③ c_t : calculate the final cell state.

Units from NN = 3



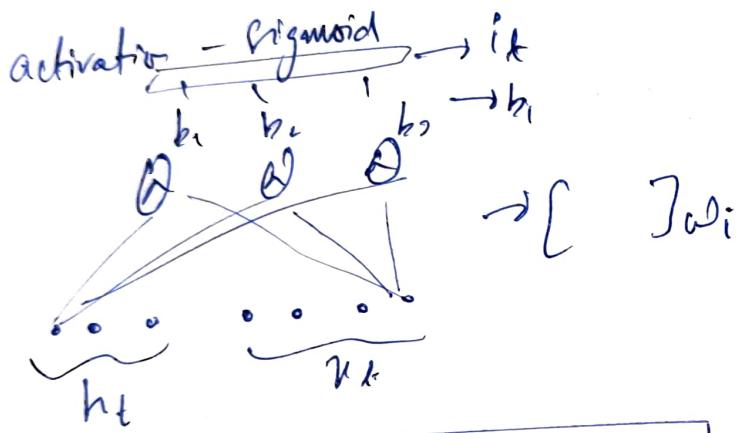
$$\rightarrow w_R (3 \times 3)$$

$$\tan(w_C(h_{t-1}, x_t) + b_C)$$

$$\bar{C}_t = \tanh \left[w_c \underbrace{(h_{t-1}, x_t)}_{(7 \times 1)} + b_c \right]$$

(3+9) (3+1)
 (3x1) (3+1)
 (3x1) (3x1)

Potential imports info = \bar{C}_t
 is follow from \bar{C}_t , next step, calculating it



$$i_t = \sigma \left[w_i \underbrace{[h_{t-1}, x_t]}_{(7 \times 1)} + b_i \right]$$

(3x7) (3+1)
 (3x1) (3+1)
 (3+1) 3+1

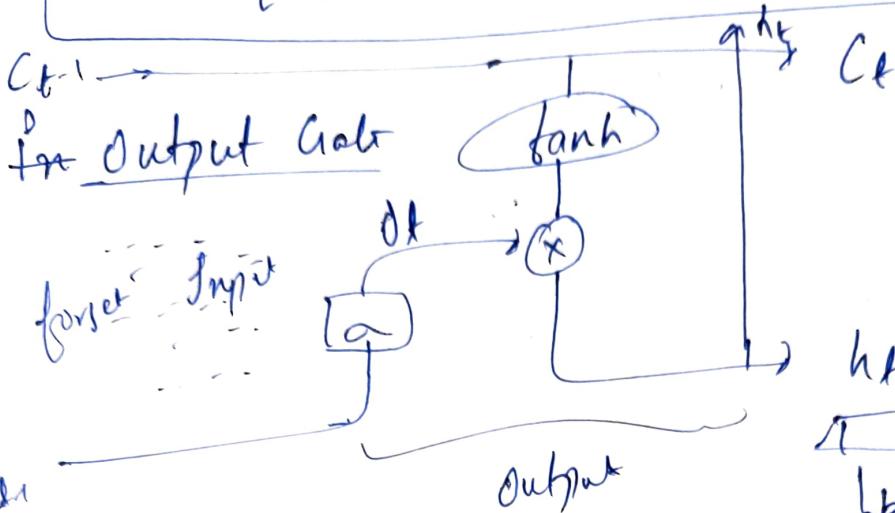
$$i_t \odot \bar{C}_t \rightarrow \bar{C}_t^+ \rightarrow \text{filtered Candidate Cell State}$$

(3x1) (3x1)

Now, finally

$$\bar{C}_t = \left[f_t \odot C_{t-1} \right] \oplus \left[i_t \odot \bar{C}_t^+ \right] = C_t$$

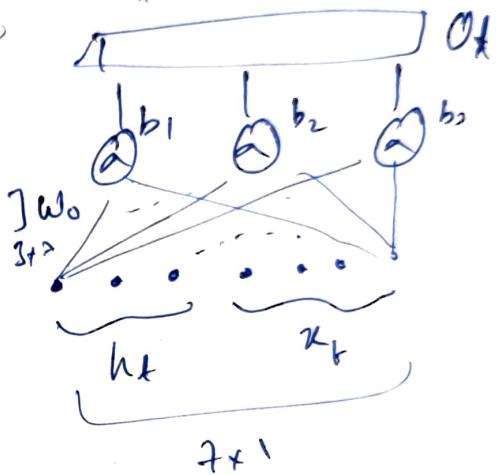
$$C_t = [f_t \odot C_{t-1}] \oplus [P_t \odot \bar{C}_t]$$



$$\odot_t \odot \tanh(C_t)$$

$$\odot_t = a [w_0 [h_{t-1}, u_t] + b_0]$$

3×7 7×1 2×1
 3×1
 (3×1)
 $a (3 \times 1) = (3 \times 1)$

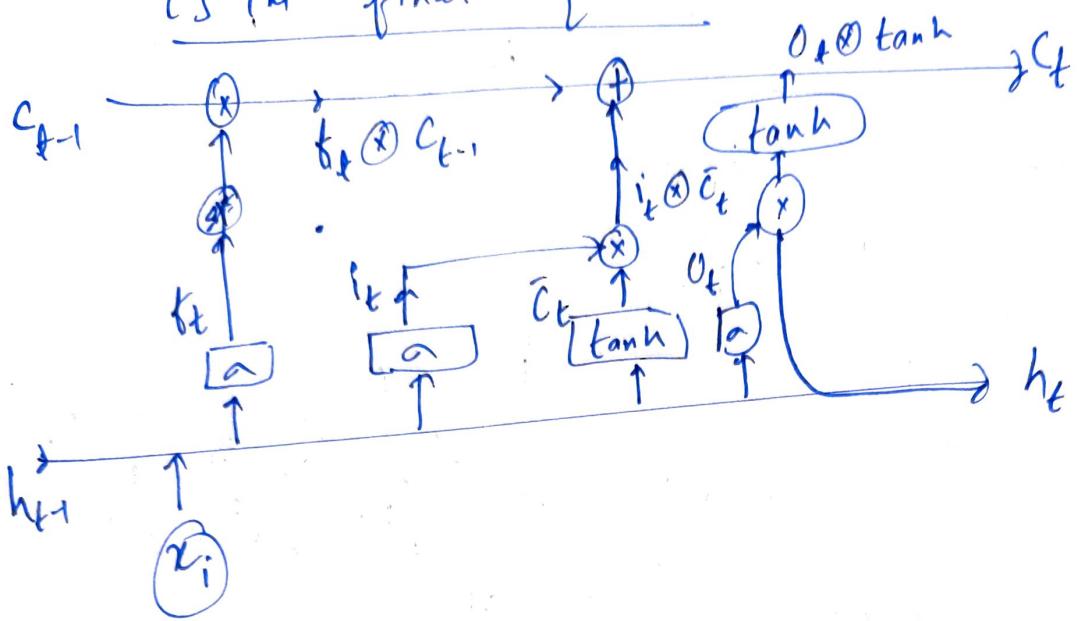


$$h_t = \odot_t \odot \tanh(C_t) \rightarrow 3 \times 1$$

3×1

3×1

LSTM final equations



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

$$② \tilde{c}_t = \tanh [w_c(h_{t-1}, x_t) + b_c]$$

$$③ i_t = \sigma [w_i(h_{t-1}, x_t) + b_i]$$

$$④ C_t = [f_t \otimes c_{t-1}] \oplus [i_t \otimes \tilde{c}_t] \#$$

$$⑤ o_t = \sigma [w_o(h_{t-1}, x_t) + b_o]$$

$$⑥ h_t = o_t \otimes \tanh(C_t) \#$$

here:-

$(h_{t-1}, x_t) \rightarrow$ concatenated Array

$\otimes, \oplus, \tanh \rightarrow$ Pointwise multiplication

Attention [is all you need!]
Used to remember longer sentences

$$\text{cos similarity} = \frac{\sum A_i \cdot B_i}{\sqrt{\sum A_i^2} \sqrt{\sum B_i^2}}$$

Self attention :

From each input vector [word embeddings]

- Create 3 vector

- 1. Query vector $\rightarrow x \cdot w_{Qk}$
- 2. Key vector $\rightarrow x \cdot w_k$
- 3. Value vector $\rightarrow x \cdot w_v$

} self attention
using dot product

$$\text{then } \rightarrow \text{self attention score} = \frac{Q_1 \cdot K_1}{\sqrt{d_k}} \quad | \quad \text{Query Vector}$$

↓
do this for each word

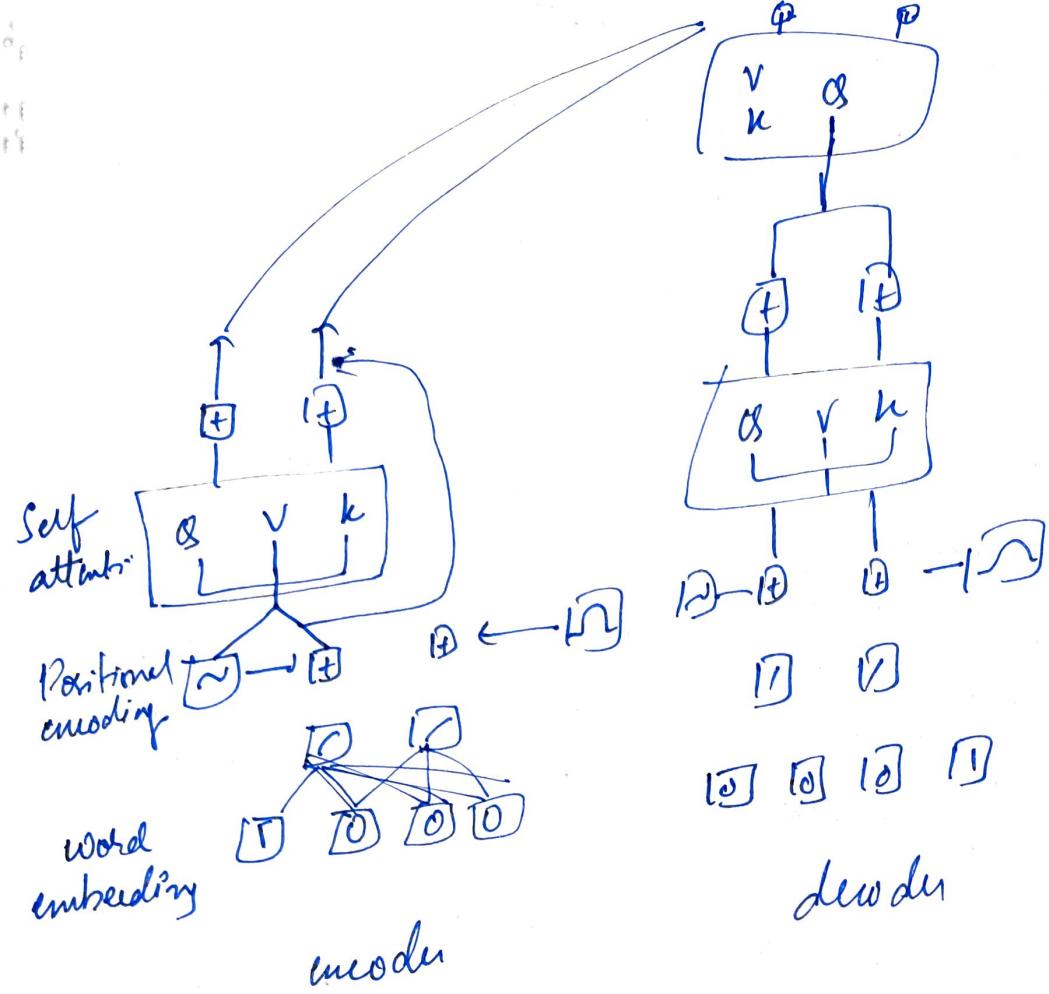
$$\Sigma \left[\text{Softmax} \left(\frac{Q_1 \cdot K_1}{\sqrt{d_k}} \right) \right] \times V_1 = Z$$

$$\text{layer Norm} = X + Z$$

Positional Encoding

↓
self attention

↓
Add and Normalize
& feed forward Network \rightarrow add and normalize.



GS 6910 / CS 7015 Attention is all you need!

they introduced new attention function.

$$\text{Score}(s_{t-1}, h_i) = V_a^T \tanh(W_{\text{att}} s_{t-1} + W_{\text{att}} h_i)$$

Key Vector

$$k_i = W_k \cdot h_i$$

Value Vector

$$v_i = W_v \cdot h_i$$

Query Vector

$$q_i = W_q \cdot h_i$$

$$q_i = W_q \cdot h_i$$

$$x = \frac{x^1}{\sqrt{m}}$$

$$\hat{x}_{ij} = \frac{1}{\sqrt{m}} \left[x_{ij} - \frac{1}{m} \sum_{k=1}^m x_{kj} \right]$$

$$\text{Covariance matrix} = x^T x = \frac{1}{m} (x')^T x'$$

$$\frac{1}{m} \sum \sum (x_{ij} - \hat{x}_{ij})^2$$

doubt = linear

low = sq. err

enough = linear

$$\min \Sigma \sum (x_{ij} - \hat{x}_{ij})^2$$

$$\min (||x - Hw^*||_F)^2$$

$$H = U_{-1 \leq k} \Sigma_{k, k}$$

$$= (x u^T) (x u^T)^T U_{-k} \Sigma_{kk}$$

$$= (x V \Sigma^T U^T) (U \Sigma V^T V \Sigma^T U^T)^{-1} U_{-k} \Sigma_{kk}$$

$$= x V \Sigma^T U^T U (\Sigma \Sigma^T)^{-1} U^T U_{-k} \Sigma_{kk}$$

$$= x V \Sigma^T (\Sigma \Sigma^T)^{-1} U^T U_{-k} \Sigma_{kk}$$

$$= x V \Sigma^T \Sigma^{-1} \Sigma^{-1} \Sigma_{kk}^{-1} U^T U_{-k} \Sigma_{kk}$$

$$= x V \Sigma^{-1} I_{-k} \Sigma_{kk}$$

$$= x V I_{-k}$$

$$W = X V I_{-k}$$

Convolutional Neural Networks

ANN - matrix multiplication

CNN - convolutional operation

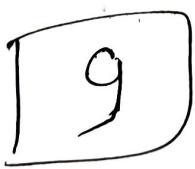
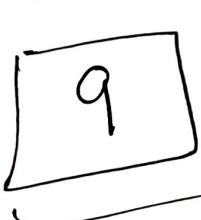
has grid-like topology (1D data / 2D data)
[]
Time series Image

Inspired by humans visual cortex

History - found by Yann LeCun (bank cheques)
then microsoft did OCR.

Why not ANN? ?
High computation cost } Image
Overfitting
less of data (spatial arrangement of pixels)
[distance between images will
be changed]

CNN Intuition
Should detect any kind of image that the model
is trained for



all of them are 9.

Convolution layer → extract features by filters

Applications

- Object localization → where is the particular obj in a Image
- Face recognition
- object detection
- Image segmentation → divide a image to segment
- Improve resolution
- Convert b/w to color Image
- Body posture
- ⋮

Road map :-

- Biological connection
- Convolution operation
Padding + stride
- Pooling layer
- CNN vs ANN

Biological Connection Human visual cortex

Visual cortex have 2 Cells -

detect particular format edge
orientation cell
or feature detector
↑
Simple cell (Edge detector)
Complex cell
|
bigger
receptive field.

Convolution layer

Conv operation

Start layer: detect edges
 next layers 2. combine edges and detect complex features
 : 3. completely identify image

types
 Vertical edge detection } done by kernels / filters
 Horizontal " "
 Slant "
 :
 Disk " "
 edges are intensity changes.

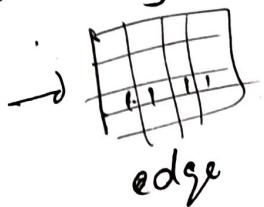
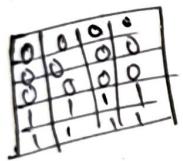


Image
 Filter
 feature map

Element wise multiplication w/ filter kernel

Output dimension after applying kernel:

Image	kernel	Output
$m \times n$	$m \times m$	$(m-m+1) \times (n-n+1)$
$m \times m \times c$	$n \times n \times c$	$(m-n+1) \times (n-n+1)$ Single channel

Padding and Strides

After many layers the size of image keeps reducing which is not a natural process so we use padding and stride and boundary pixels will be very less part of the convolution operation rather than the middle pixels

after padding, image size = $(n + 2p - f + 1)$

5×5 3×3
image kernel
↓
↓

$$\text{Padded} = s + 2(1) - 3 + 1$$

Padding Part

Padding "Valid" \rightarrow this means padding is not used
"same" \rightarrow will keep padding (0)

Stride
when applying filter on the image, it moves the filter over one pixel if stride is kept 1.
If stride is (2×2) then the filter will move a pixel
[stride ↑, dimension ↓] $\left[\frac{n+2p-f}{2} + 1 \right]$

Strided convolution
• Stride value other than 1

Stride = 1 \rightarrow default stride

stride ≠ 1 \rightarrow strided convolution

$\lceil \frac{n-f+1}{s} + 1 \rceil$ or floor operation

n = image
 f = filter
 p = padding

S = stride

Why stride / Strided convolution?

- so that the model don't miss the info and
- need high level features, no need of low level features.

then can use stride > 1

- for computing Power.

$$\left[\frac{n+2p+f}{s} + 1 \right]$$

Pooling

convolution has issues

1. Memory
2. Translation Variance
3. Location dependency

way to down sample feature map is given
Applied after non linear feature map is given

Max, min etc

Type - max, min etc

Advantages :- ① feature map size reduces
② Translation Invariance.

Disadvantage :- ① Can't use in image segmentation (75%)
② lose lot of information
(depends on application)

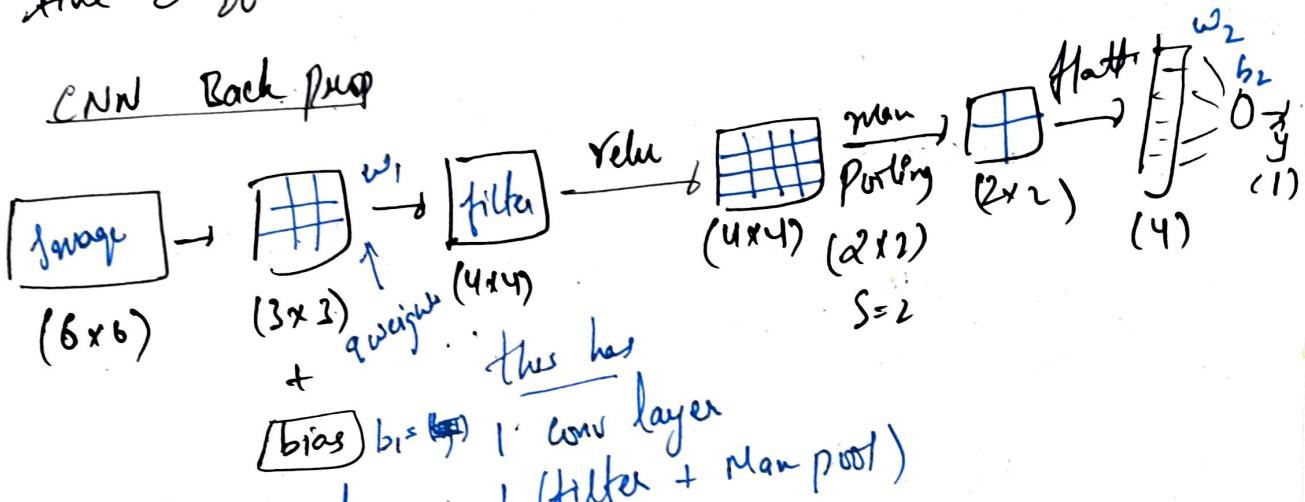
Architecture

convolution \rightarrow padding/Stride \rightarrow Pooling

CNN vs ANN

both are kinda same but small multiblock
size difference

CNN Back Prop



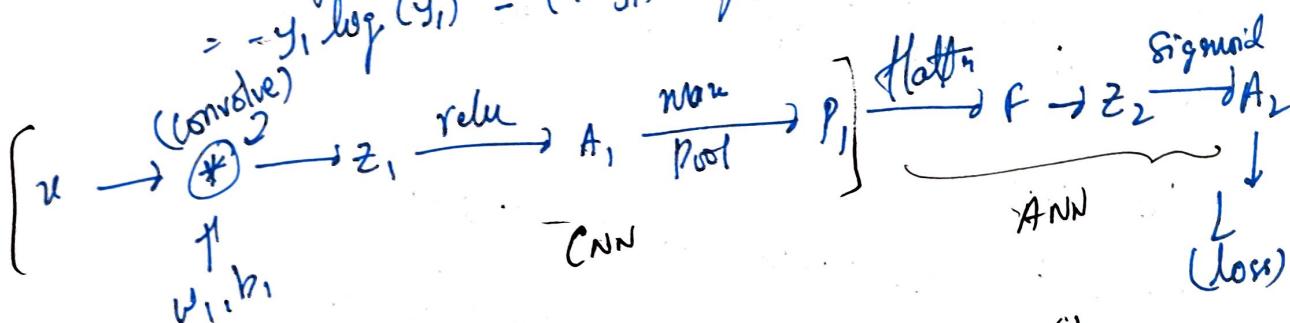
Trainable parameters

$$w_1 = (3, 3) \quad w_2 = (1, 4) \quad \left. \begin{array}{l} \text{total: 15 parameters} \\ \text{bias } b_1 = (1, 1) \end{array} \right\}$$

$$b_1 = (1, 1) \quad b_2 = (1, 1)$$

loss = binary classification (like cat/dog)

$$= -y_1 \log(\hat{y}_1) - (1-y_1) \log(1-\hat{y}_1)$$



Forward propagation

$$z_1 = \text{conv}(x, w_1) + b_1$$

$$A_1 = \text{relu}(z_1)$$

$$P_1 = \text{maxpool}(A_1)$$

$$F = \text{flatten}(P_1)$$

$$Z_2 = F w_2 + b_2$$

$$A_2 = \text{sigmoid}(Z_2)$$

Shapes

$$F = (4, 1)$$

$$w_2 = (1, 4)$$

$$Z_2 = (1, 1)$$

$$A_2 = (1, 1)$$

$$(1, 4) \cdot (4, 1) = (1, 1)$$

end goal : Apply Gradient Descent on the final output and find out loss and minimize it

Gradient descent :-

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

$$\text{gradient} \quad b_1 = b_1 - n \frac{\partial L}{\partial b_1} ; \quad b_2 = b_2 - n \frac{\partial L}{\partial b_2}$$

divide backprop into 2 parts ① CNN ② ANN

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} \underbrace{,}_{\substack{\therefore w_2 \text{ depend on } A_2, z_2 \\ \text{①}}} \quad \text{[because } w_2 \text{ depend on } A_2, z_2]$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2} \quad \text{②}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial P} \times \frac{\partial P}{\partial P_i} \cdot \frac{\partial P_i}{\partial A_i} \cdot \frac{\partial A_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_1} \quad \text{③}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial F} \cdot \frac{\partial F}{\partial P_i} \cdot \frac{\partial P_i}{\partial A_i} \cdot \frac{\partial A_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial b_1} \quad \text{④}$$

lets solve now :-

$$\frac{\partial L}{\partial a_2} = \frac{1}{\partial a_2} [-y_1 \log(a_2) - (1-y_1) \log(1-a_2)]$$

$$= \frac{-y_1}{a_2} + \frac{1-y_1}{1-a_2} = \frac{-y_1(1-a_2) + a_2(1-y_1)}{a_2(1-a_2)}$$

$$\frac{\partial L}{\partial a_2} = \frac{a_2 - y_1}{a_2(1-a_2)}$$

$$\frac{\partial A_2}{\partial z_2} = \alpha(z_2) [1 - \alpha(z_2)] = \alpha_2 (1 - \alpha_2)$$

$$\boxed{\frac{\partial A_2}{\partial z_2} = \alpha_2 (1 - \alpha_2)}$$

$$\boxed{\frac{\partial z_2}{\partial w_2} = f}$$

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

$$\frac{\partial L}{\partial w_2} = \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} \cdot \frac{\partial z_2}{\partial w_2}$$

$$\begin{aligned} \frac{\partial L}{\partial w_2} &= \frac{\alpha_2 - y_1}{\alpha_2 (1 - \alpha_2)} \times \underbrace{\alpha_2 (1 - \alpha_2) \times f}_{\substack{(1,1) \\ (1,1) \\ (1,1)}} \underbrace{f}_{(1,4)} \\ &= (\alpha_2 - y_1) f = (A_2 - y) f \end{aligned}$$

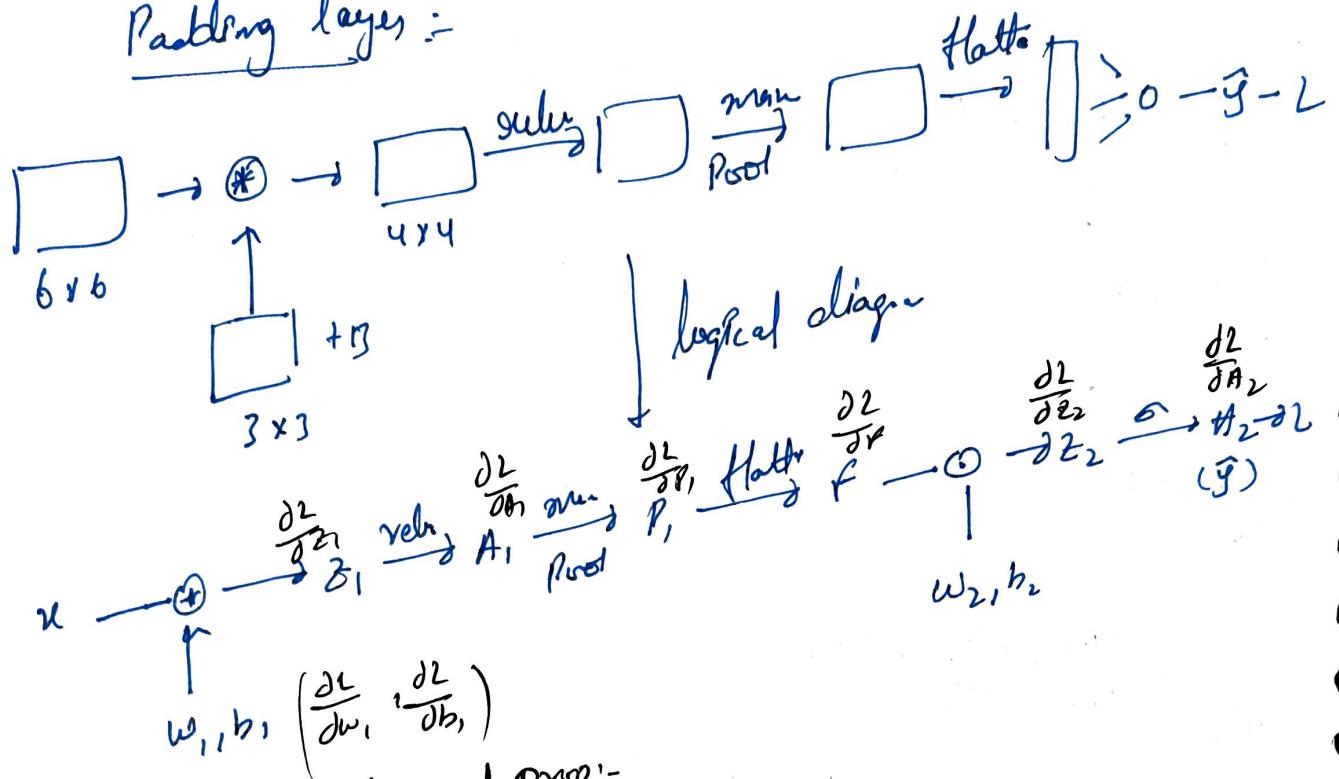
$$\boxed{\frac{\partial L}{\partial w_2} = (A_2 - y) f^T} \quad \text{Shape} = (1, 4)$$

$$\frac{\partial L}{\partial b_2} = \frac{\alpha_2 - y_1}{\alpha_2 (1 - \alpha_2)} \alpha_2 (1 - \alpha_2) \times 1$$

$$\boxed{\frac{\partial L}{\partial b_2} = (A_2 - y)}$$

Now back prop for convolution layer, max pool layer and

Padding layers :-



Equations for forward prop:-

$$\left. \begin{array}{l} z_1 = \text{conv}(x_1, w_1) + b_1 \\ A_1 = \text{relu}(z_1) \\ P_1 = \text{max pool}(A_1) \\ f = \text{Flatten}(P_1) \end{array} \right\} \begin{array}{l} z_2 = w_2 f + b_2 \\ A_2 = \sigma(z_2) \\ L = \frac{1}{m} \sum_{i=1}^m (-y_i \log(A_2) - (1-y_i) \log(1-A_2)) \end{array}$$

Should calculate ③ and ④

$$\frac{\partial z_2}{\partial f} = \frac{\partial}{\partial f} (w_2 f + b_2) = w_2$$

Should move to shape all the term
 $\frac{\partial z_2}{\partial v} = w_2$

$$\frac{\partial f}{\partial P_1} = \frac{\partial}{\partial P_1} \text{Flatten}(P_1) \rightarrow \text{Should reshape this to its original matrix size}$$

$$\frac{\partial F}{\partial P_1} = \text{reshape}(P_1, \text{Shape})$$

$$\frac{\partial L}{\partial w_1} = (A_2 - y) \omega_2 \cdot \text{reshape}(P_1, \text{Shape}) \frac{\partial P_1}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

So $\frac{\partial P_1}{\partial A_1} = \frac{\partial}{\partial A_1} \text{maxpool}(A_1) \rightarrow$ Should do inverse of maxpooling
Only more elements will have a role in loss prediction

$$\begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix}_{2 \times 2} \rightarrow \begin{bmatrix} 0 & x_1 & x_2 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & x_3 & x_4 & 0 \end{bmatrix}_{4 \times 4}$$

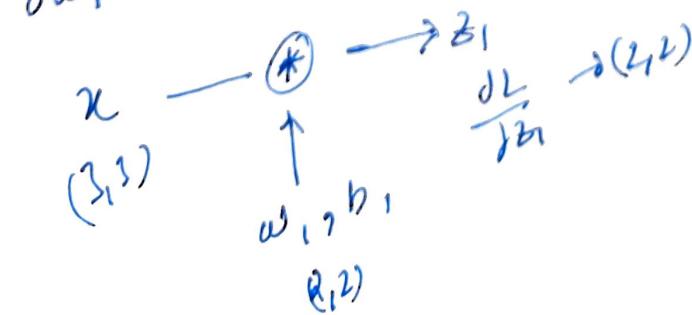
(4×4) matrix ~~(2×2)~~

Replacing x_1, x_2, x_3, x_4 to their original place ~~but~~ before max pooling is done if A_{mn} is the max element

$$\frac{\partial L}{\partial A_1} = \begin{cases} \frac{\partial L}{\partial P_1, xy} & \text{if } A_{mn} \text{ is the max element} \\ 0 & \text{if, otherwise} \end{cases}$$

$$\frac{\partial A_1}{\partial z_1} = \begin{cases} 1 & \text{if } z_1, xy > 0 \\ 0 & \text{if } z_1, xy \leq 0 \end{cases}$$

$\frac{\partial z_1}{\partial w_1} \rightarrow$ convolution backprop, this has training parameter



$$\chi = \begin{bmatrix} \chi_{11} & \chi_{12} & \chi_{13} \\ \chi_{21} & \chi_{22} & \chi_{23} \\ \chi_{31} & \chi_{32} & \chi_{33} \end{bmatrix} \quad \omega_1 = \begin{bmatrix} \omega_{11} & \omega_{12} \\ \omega_{21} & \omega_{22} \end{bmatrix} \quad z_1 = \begin{bmatrix} z_{11} & z_{12} \\ z_{21} & z_{22} \end{bmatrix}$$

$$z_{11} = \chi_{11} w_{11} + \chi_{12} w_{12} + \chi_{21} w_{21} + \chi_{22} w_{22} + b_1$$

$$z_{12} = \chi_{12} w_{11} + \chi_{13} w_{12} + \chi_{22} w_{21} + \chi_{23} w_{22} + b_1$$

$$\frac{\partial L}{\partial z_{12}} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \end{bmatrix}$$

$$z_{21} = \chi_{21} w_{11} + \chi_{22} w_{12} + \chi_{31} w_{21} + \chi_{32} w_{22} + b_1$$

$$z_{22} = \chi_{22} w_{11} + \chi_{23} w_{12} + \chi_{32} w_{21} + \chi_{33} w_{22} + b_1$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1} = \frac{\partial L}{\partial z_{11}} \cdot \frac{\partial z_{11}}{\partial b_1} + \frac{\partial L}{\partial z_{12}} \cdot \frac{\partial z_{12}}{\partial b_1} + \frac{\partial L}{\partial z_{21}} \cdot \frac{\partial z_{21}}{\partial b_1}$$

$$+ \frac{\partial L}{\partial z_{22}} \cdot \frac{\partial z_{22}}{\partial b_1}$$

$$= \frac{\partial L}{\partial z_{11}} + \frac{\partial L}{\partial z_{12}} + \frac{\partial L}{\partial z_{21}} + \frac{\partial L}{\partial z_{22}} = \text{sum} \left(\frac{\partial L}{\partial z_i} \right)$$

$$\boxed{\frac{\partial L}{\partial b_1} = \text{sum} \left(\frac{\partial L}{\partial z_i} \right)}$$

Scalar



Practical (for own)

Variational Auto Encoders

normal
model } Used for compressing
Used for segmentation in self driving cars
Reconstructing parts of a image
VAE is encoded to distribution
mean vector - Standard deviation

$$\text{loss} = \text{Reconstruction loss} + \text{KL divergence}$$

Sampling = latent Vector.

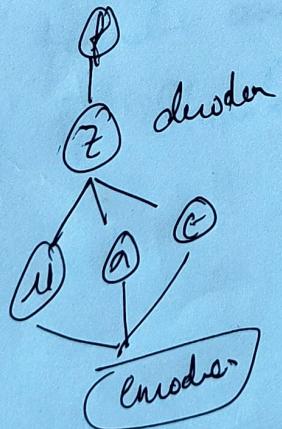
Can't push gradients to this sample b/c close to $(0,1)$
hence backprop is hard.

Reparameterization trick.

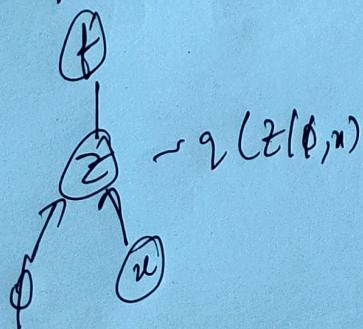
$$z = \mu + \sigma \theta$$

$\theta \sim \text{Normal}(0, 1)$

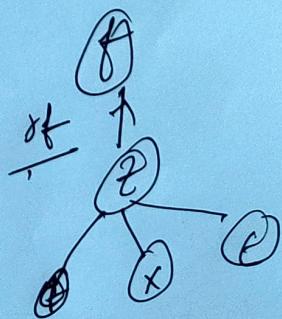
/
mean
std



original

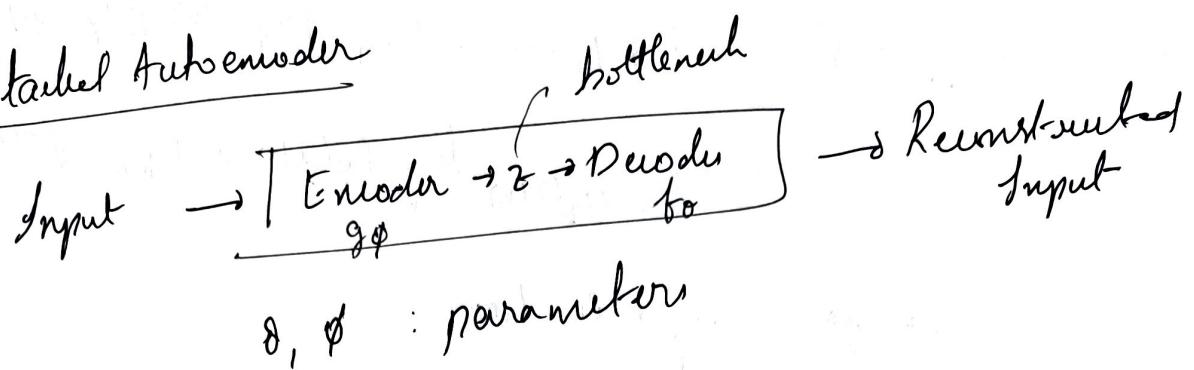


Reparameterization form



VAE by Abhishek Kumar on youtube.

Stacked Autoencoder

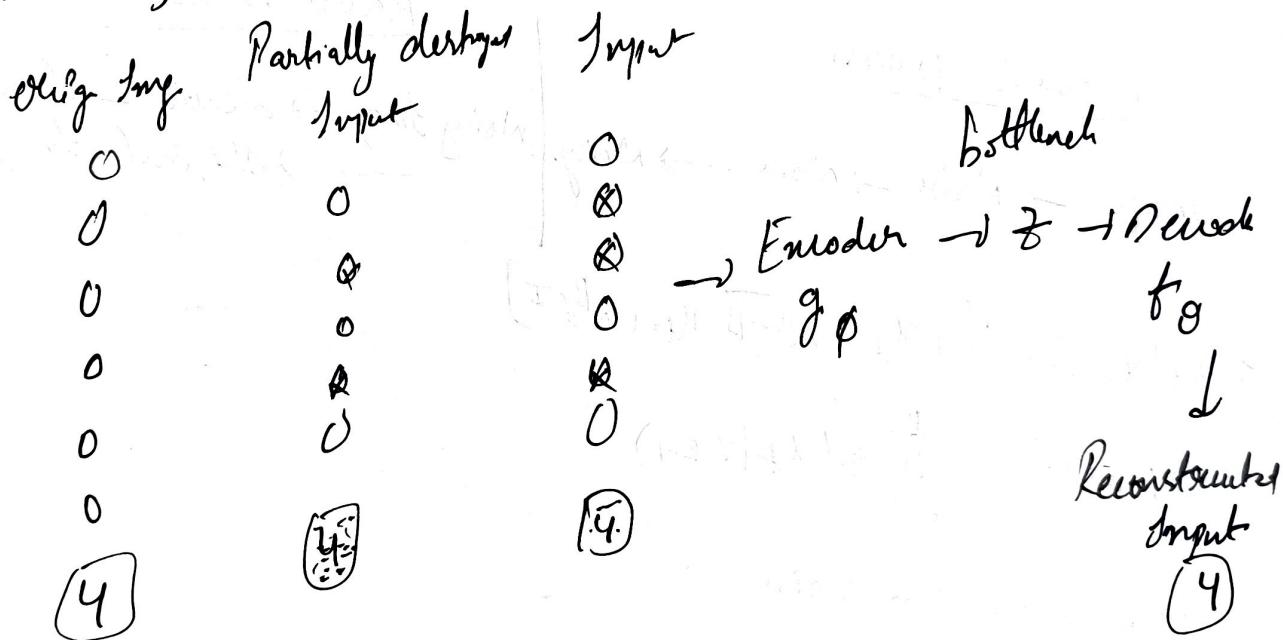


$$\text{cost function} = L(\theta, \phi)$$

$$L(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n [x^i - f_\theta [g_\phi(x^i)]]^2$$

Reconstructed Image

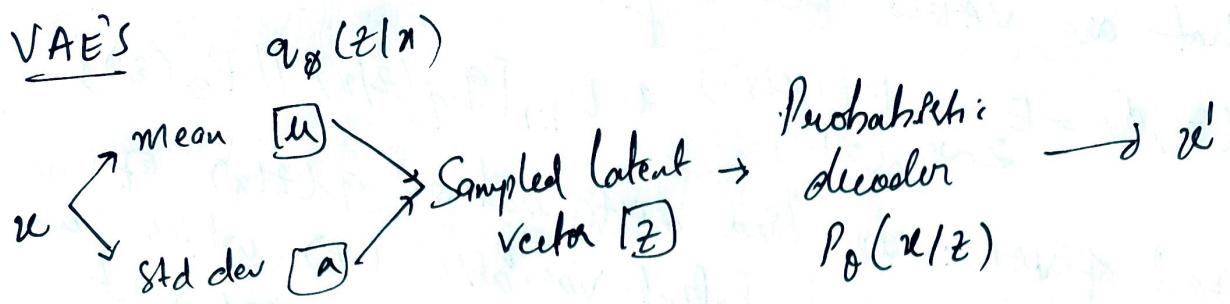
Denoising Autoencoder :-



$$L(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n [x^i - f_\theta [g_\phi(\tilde{x}^i)]]^2$$

$$\tilde{x}^i = x^i | \epsilon^i$$

VAE's



Probabilistic Encoder

$$z = \mu + \sigma \epsilon$$

$$\epsilon \sim \mathcal{N}(0, I)$$

$$L(\theta, \phi) = -E_{z \sim q_{\phi}(z|x)} [P_0(x|z)] + D_{KL}[q_{\phi}(z|x) || P_0(z)]$$

$$E(n) = \sum_{i=1}^n x_i P(x_i)$$

KL Divergence

Used to measure how 1 p.d. is different from other

P.d.

P.d. = Probabilistic distribution



KLD is high for this.

μ_{KL} is almost or close to 0.

$$D_{KL}(P||Q) = \sum P(x=u) \log \left[\frac{P(x=u)}{Q(x=u)} \right]$$

$$\Rightarrow \sum P(n) \log \left[\frac{P(n)}{Q(n)} \right]$$

what are VAE's and why we need it?

$$L(\theta, \phi) = E_{z \sim q_\phi(z|x)} P(x|z) + D_{KL}[q_\phi(z|x) || P_\theta(z)]$$

goal of VAE: To find distribution $q(z|x)$ of some latent variable (z) which we can sample from $z \sim q_\phi(z|x)$ to generate new samples $x' \sim P_\theta(x|z)$

Approximate Inference

x - observed Variable

z - latent Variable

$$P(z|x) = \frac{P(x|z) \cdot P(z)}{P(x)}$$

(conditional prob)

difficult to obtain $P(x)$

$P(x) = \int P(x|z) P(z) dz \rightarrow$ not closed integral or intractable due to multiple integral

alternative?

approximate $P(z|x)$ by $q(z|x)$ where $q(z|x)$ are involved has tractable sol done by Variational Inference

Posing inference problem as optimization problem by modelling $P(z|x)$ using $Q(z|x)$ which has simple Gaussian distribution

$$D_{KL} [Q_\phi(z|n) \parallel P_\theta(z|n)] = \sum \downarrow Q(z|n) \log \left[\frac{Q(z|n)}{P(z|n)} \right]$$

↑
gaussian ↓ one we are trying
to approximate

$$= E_{z \sim Q(z|n)} \left[\log \left[\frac{Q(z|n)}{P(z|n)} \right] \right]$$

$$D_{KL} [Q_\phi(z|n) \parallel P_\theta(z|n)] = E_{z \sim Q(z|n)} [\log Q_\phi(z|n) - \log P_\theta(z|n)] \quad \boxed{B}$$

Substituting \boxed{A} in \boxed{B} , we get

$$D_{KL} [Q_\phi(z|n) \parallel P_\theta(z|n)] = E_{\mathbb{E}_{\mathbb{Q}_\phi}} \left[\log Q_\phi(z|n) - \log \frac{P(x|z) P(z)}{P(n)} \right]$$

$$= E_{\mathbb{Q}_\phi} (\log Q_\phi(z|n) - \log P_\theta(x|z) - \log P_\theta(z) + \log P_\theta(n))$$

\mathbb{Q}_ϕ is being sampled by $z \sim Q_\phi(z|n)$ = $\boxed{2}$

Since expectation is over z , $P(n)$ doesn't involve z , it can be moved out.

$$D_{KL} [Q_\phi(z|n) \parallel P_\theta(z|n)] = \log P_\theta(n) = E_{\mathbb{Q}_\phi} (\log Q_\phi(z|n) - \log P_\theta(x|z) - \log P_\theta(z))$$

Rearranging equations :-

$$\log P_\theta(n) - D_{KL} [Q_\phi(z|n) \parallel P_\theta(z|n)] = E_{\mathbb{Q}_\phi} [\log P(x|z)] - E_{\mathbb{Q}_\phi} [\log Q(z|n) - \log P(z)]$$

this is basically KL divergence between $Q(z|n)$ and $P(z)$

$$= E_{\mathbb{Q}_\phi} [\log P(x|z)] - D_{KL} [Q(z|n) \parallel P(z)]$$

thus is VAE objective function :-
 1st term = Reconstruction likelihood
 2nd term = ensure learned distribution Q is similar to prior distribution P

$$\text{loss} = -[\text{objective function}]$$

$$L = L(\theta, \phi) = -E_{\mathcal{D}} [\log P_{\theta}(x|z)] + D_{KL} [P_{\theta}(z|x) \parallel P(z)]$$

Also :- $h_{\theta, \phi}$, $P_{\theta, \phi}$.

$$\underbrace{\log P_{\theta}(x)}_{\text{Reconstruction loss.}} - D_{KL} [P_{\theta}(z|x) \parallel P(z|x)] = -L(\theta, \phi)$$

Inherent Regularizer

$\theta^*, \phi^* = \arg \min L(\theta, \phi) \rightarrow$ So our target is to find optimal θ, ϕ such that

weights of generative model
(density)

weights of encoder.

Reconstruction loss

$$-E_{z \sim P_{\theta}(z|x)} [\log P_{\theta}(x|z)]$$

Regularizer

$$D_{KL} [P_{\theta}(z|x) \parallel P(z)]$$

What are these both?

log likelihood [alias] $P_{\theta}(x|z) = N(\mu_{\theta}(z), \Sigma_{\theta}(z))$

Recognition Model
for this we get squared error
between data sample and
mean of gaussian distribution.

$$P_{\theta}(x|z) = \frac{1}{\sqrt{2\pi^k |\Sigma_{\theta}(z)|}} \exp \left[\frac{[x - \mu_{\theta}(z)]^T \Sigma_{\theta}^{-1}(z) (x - \mu_{\theta}(z))}{2} \right]$$

$$\log P_{\theta}(x|z) \propto \underbrace{x - \mu_{\theta}(z)^T \Sigma_{\theta}^{-1}(z) (x - \mu_{\theta}(z))}_{\text{Squared Reconstruction Error}}$$

[data fidelity term] In this problem

$$P_\theta(z|x) = \frac{P(x|z) P(z)}{P(x)} \quad \stackrel{\text{VAE}}{=} \quad (A) \quad ①$$

$$D_{KL}[\alpha(z|x) || P(z|x)] = \sum \alpha(z|x) \cdot \log \frac{\alpha(z|x)}{P(z|x)}$$

$$= E_{\alpha} \left[\log \left[\frac{\alpha(z|x)}{P(z|x)} \right] \right]$$

$$D_{KL}[\alpha(z|x) || P(z|x)] = E_{\alpha} \left[\log \alpha(z|x) - \log P(z|x) \right] \quad ②$$

Putting ① in ②

$$= E_{\alpha} \left[\log \alpha(z|x) - \log \left[\frac{P(x|z) P(z)}{P(x)} \right] \right]$$

$$= E_{\alpha} \left[\log \alpha(z|x) - \log (P(x|z) P(z)) + \log P(x) \right]$$

$$= E_{\alpha} \left[\log \alpha(z|x) - \log P(x|z) - \log P(z) + \log P(x) \right]$$

$$\alpha = \alpha_\theta \quad ; \quad P = P_\theta \quad ; \quad \alpha = z \sim \alpha_\theta(z|x)$$

$\log P(x)$ can be moved out as it doesn't involve z

$$D_{KL}[\alpha(z|x) || P(z|x)] - \log P(x) = E_{\alpha} \left[\log \alpha(z|x) - \log P(x|z) - \log P(z) \right]$$

$$\log P(x) - D_{KL}[\alpha(z|x) || P(z|x)] = E_{\alpha} \left[\log P(x|z) \right] - E_{\alpha} \left[\log \alpha(z|x) - \log P(z) \right]$$

$$= E_{\alpha} \left[\log P(x|z) \right] - D_{KL}[\alpha(z|x) || P(z)]$$

Reconstruction likelihood

VAE objective fn

Ensures α (learned distribution)
similar to P (prior distribution)

$\text{loss} = -\text{objective fn}$

(2)

$$L(\theta, \phi) = -E[\log P(x|z)] + D_{KL}[Q_\phi(z|n) \parallel P(z)]$$

or

$$-L(\theta, \phi) = \log P_\theta(x) - D_{KL}[Q_\phi(z|n) \parallel P(z)]$$

$$D_{KL}[Q_\phi(z|n) \parallel P(z)] \quad \text{Regularizer}$$

$P_\theta(z)$ = latent variable distribution

Easiest choice = $N(0, I)$

want this to be close to $Q_\phi(z|n)$ to sample it easily
($P_\theta(z)$)

If we want $Q_\phi(z|n)$ to be gaussian with parameters $\mu_{\phi(n)}$, $\Sigma_{\phi(n)}$
then the KL divergence has the closed form as derived

when $P_\theta(z) = N(0, I)$

$$D_{KL}\left[N(\mu_n, \Sigma(n)) \parallel N(0, I)\right] = \frac{1}{2} \left[\text{tr}(\Sigma(n) + \mu_n \mu_n^T) - k - \log |\Sigma(n)| \right] \quad (2)$$

k = dimensionality of latent variables

$\text{tr}(\Sigma(n))$ = trace fn =
sum of diagonal
Matrix of $\Sigma(n)$

Proof =

$$\overline{E_P[(x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2)]}$$

$$\Rightarrow [E_P[(x - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2)]]$$

$$\Rightarrow (\mu_1 - \mu_1)^T \Sigma_2^{-1} (\mu_1 - \mu_2)$$

$$\Rightarrow 0 \quad (\text{proved})$$

$$KL\left[\frac{P(n)}{\Sigma_1} \parallel \frac{Q(n)}{\Sigma_2}\right] = \frac{1}{2} \left[\log \left[\frac{\Sigma_2}{\Sigma_1} \right] - k + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right] \quad (1)$$

$$\begin{array}{l|l} \Sigma_2 = I & \Sigma_1 = \Sigma_{\phi(n)} \\ \mu_2 = 0 & \mu_2 = \mu_{\phi(n)} \end{array}$$

by putting these values in (1) we get (2) (main fn)

$$\begin{aligned}
 D_{KL} [N(\mu_\phi(x), \Sigma(x)) // N(0, I)] &= \\
 &= \frac{1}{2} \left[\sum_k \Sigma_\phi(k) + \Sigma_k \mu_\phi^2(k) - \sum_k 1 - \log \prod_k \Sigma(k) \right] \\
 &= \frac{1}{2} \left[\sum_k \sum_p (k) + \sum_k \mu_\phi^2(k) - \sum_k 1 - \sum \log (\Sigma_\phi(k)) \right] \\
 &= \frac{1}{2} \sum_k \left[\Sigma_\phi(k) + \mu_\phi^2(k) - 1 - \log (\Sigma_\phi(k)) \right]
 \end{aligned}$$

Simplified closed form expression for KL div

Reparameterization

optimizing loss fn:

$$\theta^*, \phi^* = \arg \min_{\theta, \phi} [L(\theta, \phi)]$$

ELBO := evidence lower bound

$\boxed{KL \geq 0}$ property
therefore ; $L(\theta, \phi)$ = lower bound of $\log P_\theta(x)$

$$\log P_\theta(x) - \underbrace{D_{KL} [\phi(z|x) // P(z|x)]}_{\text{always } \geq 0} = -L(\theta, \phi)$$

$$\therefore \boxed{L(\theta, \phi) \leq \log P_\theta(x)}$$

\downarrow loss \uparrow lower bound of P (generating real data)

Recall :-

$$L(\theta, \phi) = -E_{\phi} \log P_\theta(x|z) + \frac{1}{2} \sum_k \left[\exp(\Sigma_\phi(k)) + \mu_\phi^2(k) - 1 - \Sigma_\phi(k) \right]$$

alternative optimization principle

①

$$\theta^*, \phi^* = \arg \min L(\theta, \phi)$$

$$\underbrace{\theta^*}_{\theta} = \nabla_{\theta} [L(\theta, \phi)] ; \phi = \text{constant}$$

$$\underbrace{\phi^*}_{\phi} = \nabla_{\phi} [L(\theta, \phi)] ;$$

$$\theta^{t+1} = \nabla_{\theta} [L(\theta, \phi)]$$

the cycle keeps repeating for certain iterations

Problem will occur while deriving w.r.t ϕ to learn $\phi_{\theta}(z|x)$ and $P_{\theta}(x|z)$ at same time

$$\underset{\theta, \phi}{\text{Min}} L(\theta, \phi) = \underset{\theta, \phi}{\text{Min}} \left\{ -E_{\theta} (\log P_{\theta}(x|z)) + \frac{1}{2} \sum \left[\exp(\sum_{\theta}(w)) + f_{\phi}^2(z) - 1 - \sum_{\theta}(w) \right] \right\}$$

$$\hat{\theta}_i = \nabla_{\theta} \nabla_{\theta} [L(\theta, \phi)]$$

$$\hat{\theta}_i = \frac{1}{2} \sum_{j=1}^L \nabla_{\theta} \log P_{\theta}(x|z^j) + 0 : z^j - \phi_{\theta}(z|x)$$

$$\hat{\phi}_i = \nabla_{\phi} \{ L(\theta, \phi) \}$$

= can't differentiate 1st term + can do for 2nd term

∇_{ϕ} is harder to

estimate

so here we use reparametrization technique

$$\text{i.e., } \nabla_{\phi} E_{\phi_{\theta}(z|x)} (f(z)) \neq E_{\phi_{\theta}(z|x)} (\nabla_{\phi} f(z))$$

If we can somehow rewrite the expectation such that ϕ appears inside the expectation. Then we can differentiate.

i.e.,

$$E_{\theta_p(z|x)}[f(z)] = E_{P(\epsilon)}[f(g_\phi(\epsilon, x))]$$

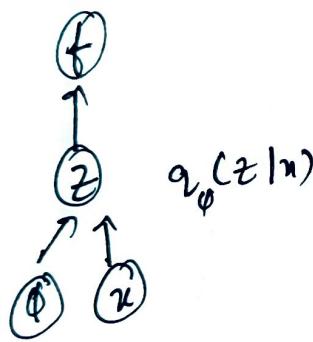
such that $z = g_\phi(\epsilon, x) \sim \epsilon \sim N(0, 1)$

linear transformation $\Rightarrow g_\phi(\epsilon, x)$

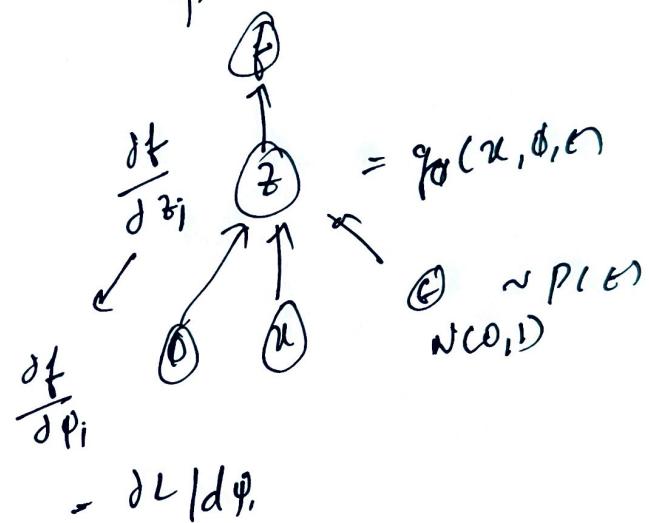
$$g_\phi(x, n) = \mu_\phi(n) + \epsilon \odot \Sigma_\phi^{1/2}(n) = z \sim N(\mu(n), \Sigma(n))$$

Mean Covariance

Original



Reparameterized



finally

$$\hat{\phi}_i = -E \left[\nabla_\phi (\log P_\phi(u|z^1)) \right] + \underbrace{\nabla_\phi \left[\frac{1}{2} \left(\sum_n [\exp(\Sigma(n)) + \psi_\phi^2(n) - 1 - \Sigma(n)] \right) \right]}_{II}$$

$$= -\frac{1}{S} \sum_{s=1}^S \log P_\phi(u|z^1) +$$

Monte Carlo

estimate

$$z^1 = \mu_\phi(n) + \epsilon \odot \Sigma_\phi(n) \quad \text{as } \epsilon \sim N(0, 1)$$

Viet

End

GANs START

Generative model or

Regularizer

$$D_{KL} [Q_\phi(z|u) \parallel P_\theta(z)]$$

Recognition model

that takes mapping
from u and produce
output z

gaussian $\{ \mu \text{ mean}$
 $\sigma^2 \text{ variance}$

parameters are known.

We want the generated
data z , not to be too far from Normal distribution

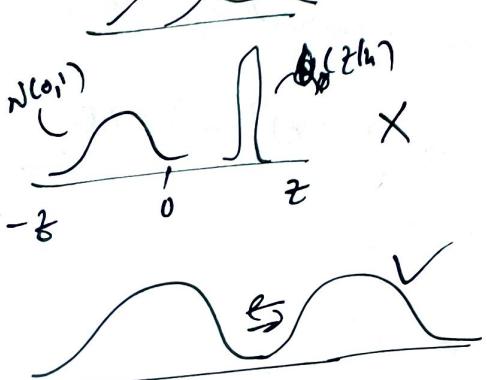
$$\sim N(0, 1)$$

this is basically allowing the latent variable to move away from $(0, 1)$. (Acting like softmax activation for)

Pdf of latent Variable $\rightarrow Q_\phi(z|u)$

$\sim N(0, 1)$ lets it stay in this range.

$$L = \text{data fidelity} + \lambda L_2 \text{ divergence}$$



\times → deviates from distribution

✓ → Network cheats by learning narrow distribution Variance is very small

✓ → Sufficient Variance is ensured attraction b/w 2 is there

Simplifying Regularizer Term $D_{KL} [Q_\theta(z|x) \parallel P_\phi(z)]$

$P_\phi(z)$ = latent variable

mean = 0

variance > 1

→ here, we need $P_\phi(z)$ to be as close to $P_\theta(z|x)$ to sample it easily.

→ Having $P_\phi(z) = \mathcal{N}(0, 1)$ adds another benefit

optimization of loss function

$$\theta^*, \phi^* = \{\arg \min \mathcal{L}(\theta, \phi)\}$$

ELBO Proof :

$$\log P_\theta(x) - D_{KL} [Q_\theta(z|x) \parallel P_\theta(z|x)] = -\mathcal{L}(\theta, \phi) \quad (\text{ELBO})$$

and we know ; $D_{KL} [Q_\theta(z|x) \parallel P_\theta(z|x)] \geq 0$

hence ; $\mathcal{L}(\theta, \phi) \leq \log P_\theta(x)$ by putting this in

In Variational Bayesian Method, this loss for ($\mathcal{L}(\theta, \phi)$) is known as Variational lower Bound or Evidence lower bound (ELBO). Thus "lower Bound" part comes from the fact that KL Div is always Non-Negative and thus $\mathcal{L}(\theta, \phi)$ is the lower bound of $\log P_\theta(x)$. Therefore we are minimizing the loss and maximizing lower bound of probability of generating real data samples $\hookrightarrow [\mathcal{L}(\theta, \phi) \leq \log P_\theta(x)] \rightarrow \underline{\text{ELBO}}$

Reparameterization Trick

Problem is. doing derivative with ϕ . Doing derivative with respect to θ is absolutely fine.

Solution :-

$$L(\theta, \phi) = -E_{\theta}[\log P_{\theta}(z|x)] + \frac{1}{2} \sum_k \left[\exp[\epsilon_{\phi}(x)] + P_{\phi}^2(x) - 1 - \epsilon_{\phi}(x) \right]$$

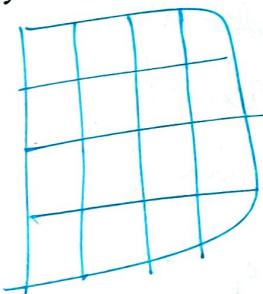
optimization is carried out wrt θ and ϕ at same time

Run algo upto fixed number of iterations.

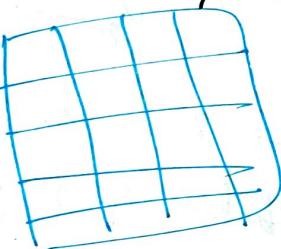
Niterm khapra

GAN's Math

true distribution



$P_{Gn}(x)$



$P_{data}(x)$

distribution of data generated by model

we want to prove $P_{Gn}(x) = P_{data}(x)$

we don't have $P_{Gn}(x)$ formula
we can say this in different versions
→ if the discriminators loss is at its minimum
it can happen only if $P_{Gn}(x) = P_{data}(x) \rightarrow$ obj fn

If $P_{Gn}(x) = P_{data}(x) \rightarrow$ global minimum is achieved
 $c(x) = \min_D V(G, D)$ achieved only if $P_G = P_{data}$

Only If
Step 1 :- we need to find $V(G, D)$

① find when $P_G = P_{data}$

② find when $P_G \neq P_{data}$

$$③ \text{ S-T :- } a < b \quad \# P_a = P_{\text{data}}$$

obj' fn :-

$$\min \max \left[E_{n \sim P_{\text{data}}} \log D_\theta(n) + E_{z \sim P(z)} \log [1 - D_\theta(a_{\text{gen}}(z))] \right] \quad ①$$

expanding to its integral :-

$$\min \max \left[\int_n P_{\text{data}}(n) \log D_\theta(n) \right] + \int_z P(z) \log [1 - D_\theta(a_{\text{gen}}(z))] \quad \text{this is generated } x$$

~~$$\min \max \left[\int_n P_{\text{data}}(n) \log D_\theta(n) + \int_z P(z) \log [1 - D_\theta(n)] \right] \quad ②$$~~

$\therefore G_\theta(n) = x$ because $P_{\theta}(n)$ denotes distribution of x generated by generator

Rule → Since x is function of z , we can replace it.

{Change of
Variables.
Assume its
invertible.}

here P_a is the function of z
which is many to many function

Rewised obj fn - ② :-

$$\min \max \int_n P_{\text{data}}(n) \log D_\theta(n) + P_a(n) \log [1 - D_\theta(n)] \, dn \quad ②$$

when will it get minimum?

when every term in sum is minimized

take derivative and set it to 0

$$\frac{d}{d[D_\theta(x)]} P_{\text{data}}(x) \log D_\theta(x) + P_a(n) \log (1 - D_\theta(n)) = 0$$

$$P_{\text{data}}(n) \frac{1}{D_\theta(n)} + P_a(n) \frac{1}{1 - D_\theta(n)} (-+) = 0$$

$$\frac{P_{\text{data}}(n)}{D_\theta(n)} = \frac{P_a(n)}{1 - D_\theta(n)}$$

$$P_{\text{data}}(n) (1 - D_\theta(n)) = P_a(n) D_\theta(n)$$

$$D_\theta(n) = \frac{P_{\text{data}}(n)}{P_a(n) + P_{\text{data}}(n)} \quad \left. \begin{array}{l} \text{we get optimal} \\ \text{discriminator} \end{array} \right\}$$

→ it is a score given x , the best score that can be assigned.
for any given generator, the optimal discriminator is given

$$\text{by } D_\theta(n) \quad \therefore D_\theta(n) = \frac{P_{\text{data}}}{P_{\text{data}} + P_a} = 1/2$$

If $P_{\text{data}} = P_a$ \therefore half probability is real / false

$$V(G, D) = \int P_{\text{data}}(n) \log D_\theta(n) + P_a(n) \log (1 - D_\theta(n)) dn$$

$$= \int P_{\text{data}}(n) \log \frac{1}{2} + P_a(n) \cdot \log (1 - \frac{1}{2}) dn$$

$$P_a = P_{\text{data}}$$

$$= \log 2 \underbrace{\int_1^2 P_a(n) dn}_{1} - \log 2 \underbrace{\int_1^2 P_{\text{data}}(n) dn}_{1}$$

$$V(G_1, D) = -2 \log 2$$

$$V(G_1, D) = -\log 4$$

when $P_a \neq P_{data}$

to show this, we get rid of assumption that

$$P_a = P_{data}$$

$$C(G_1) = V(G_1, D) = \int P_{data}^{(n)} \log \left[\frac{P_{data}^{(n)}}{P_a + P_{data}} \right] + P_a \log \left[-\frac{P_{data}}{P_a + P_{data}} \right] dn$$

$$= " + \underbrace{\left[\log^2 - \log^2 (P_{data} + P_a) \right]}_0 dn$$

refer slides.

rearranging terms:

$$- \log 4 + \text{Some term}$$

$$= - \log 4 + K^2 \text{ div}$$

$$= - \log 4 + [\geq 0]$$

min attained only when $(-\log 4)$

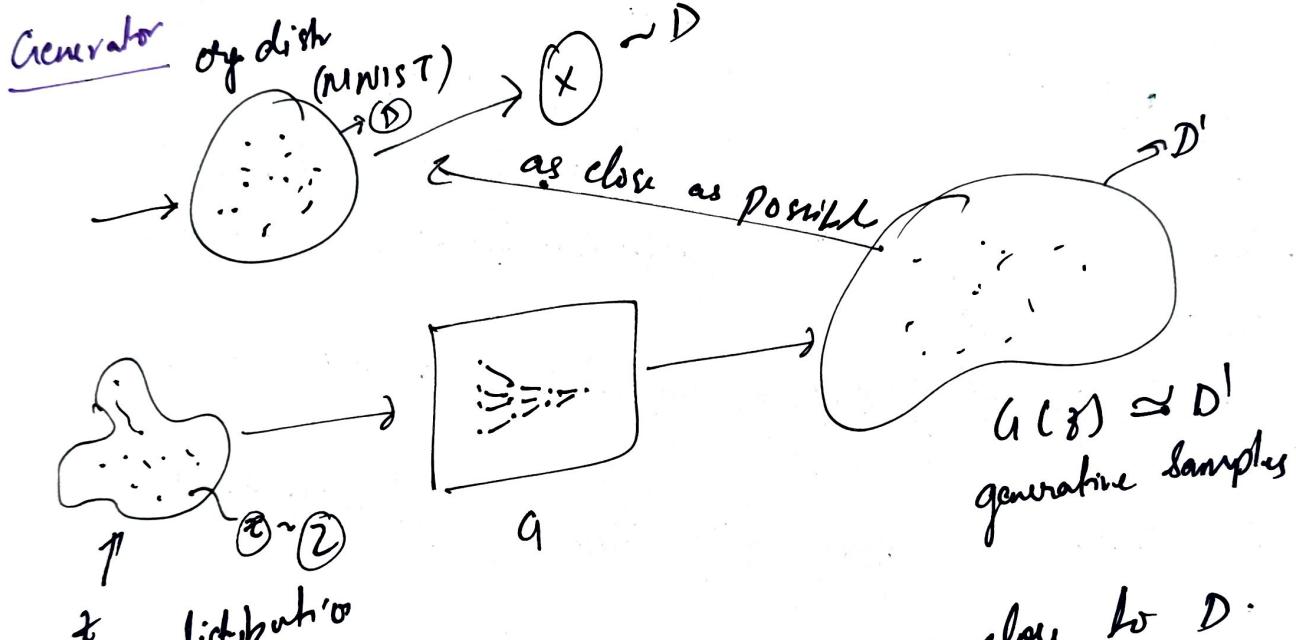
$$C(G_1) \geq - \log 4$$

lower bound

$$\text{when } P(G_1) = P_{data}$$

GAN

generator \rightarrow fight with each other } deep NN / CNN / Vanilla
 discriminator (discriminates b/w real and fake) $(0|1)$
 trained on x , Sampled on distribution D , produce D'
 given random distribution z
 (true)
 false real



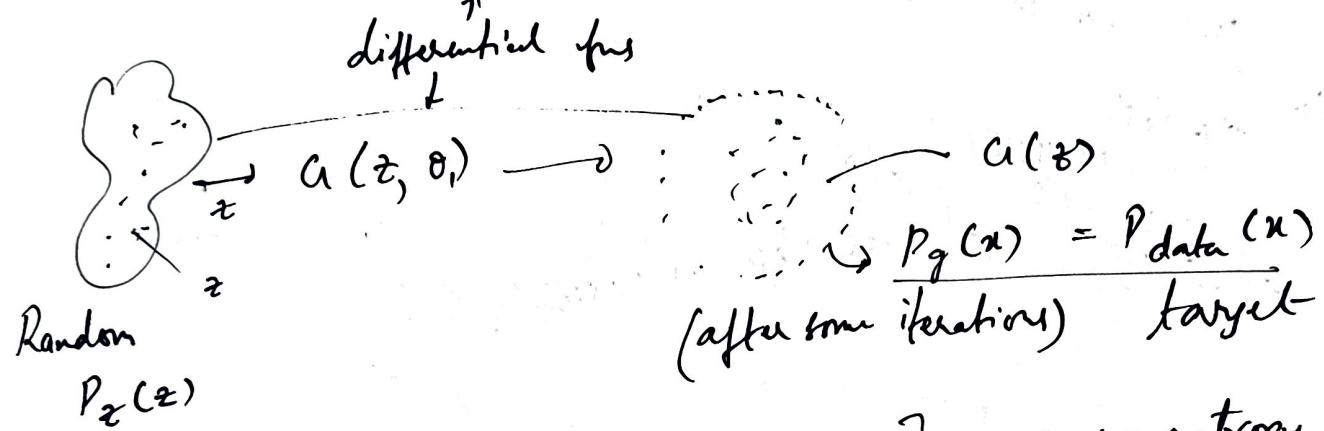
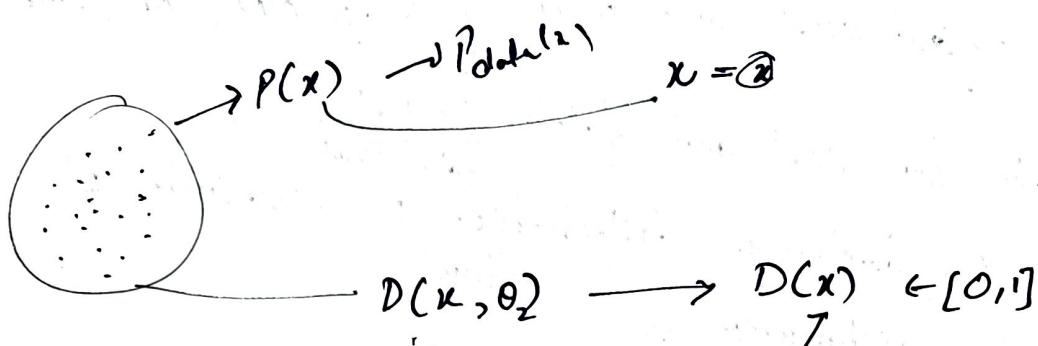
Random distribution
 Objective = producing D' which is as close to D .

Optimal Discriminator = 0.5

when it is 0.5, the discriminator fails to differentiate between real and fake data

While training generator, the discriminator is constant, vice versa
 [hard] [easy]

loss / cost function



$$L(\hat{y}, y) = [y \log \hat{y} + (1-y) \log (1-\hat{y})] \rightarrow \text{cross entropy}$$

$$\begin{aligned} \hat{y} &= D(x) \\ y &= 1 \end{aligned} \quad \left. \begin{array}{l} \text{labeled from } P_{\text{data}}^{(x)} \\ \text{Discriminator} \end{array} \right.$$

$$L(D(x), 1) = \log D(x) \rightarrow A$$

$$\begin{aligned} y &= 0 \\ \hat{y} &= D(a(x)) \end{aligned} \quad \left. \begin{array}{l} \text{data from Generator} \\ \text{Discriminator} \end{array} \right.$$

$$L(D(a(x)), 0) = \log [1 - D(a(x))] \rightarrow B$$

We need A and B to be maximized
 By doing so, A will force $D(x)$ to be 1
 and B will force $D(a(x))$ to be 0

$$\text{Max} [\log D(x) + \log (1 - D(a(x)))] = D$$

$$\min_a \max_D \left[\log D(x) + \log 1 - D[G(z)] \right]$$

obj fn

Original Eq:-

$$\min_a \max_D V(D, G) = \min_a \left[E_{x \sim P_{\text{data}}} [\log D(x)] + E_{z \sim P_g(z)} [\log 1 - D[G(z)]] \right] \quad //$$

Obj of generator = fool the discriminator

$$y=0 \quad \hat{y} = D[G(z)]$$

$$L(D[G(z)], 0) = (1) \log [1 - D[G(z)]]$$

$D[G(z)] = 1 \rightarrow$ Should be
this is possible when we minimize this

Finding best Discriminator

$$D_a^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)}$$

derivation

$$D_a^* = \arg \max V(D, G)$$

$$D_a^* = \arg \max \left\{ E[\log D(x)] + E[\log (1 - D(x))] \right\}$$

$$E_{P(x)}(x) = \int x P_x(x) dx \rightarrow \text{we know}$$

GANS END

DIFFUSION MODELS

START

Diffusion Models

$$\text{VAE} = L_{\text{VAE}}(\theta, \phi) = -\log P_\theta(x) + D_{KL}[q_\phi(z|x) \parallel P_\theta(z|x)]$$

$$= -\mathbb{E}_{z \sim q_\phi(z|x)} \log P_\theta(x|z) + D_{KL}[q_\phi(z|x) \parallel P_\theta(z|x)]$$

$$\theta^*, \phi^* = \arg \min L_{\text{VAE}}$$

$$- L_{\text{VAE}} = \log P_\theta(x) - D_{KL}[q_\phi(z|x) \parallel P_\theta(z|x)] < \log P_\theta(x)$$

~~GAN~~ → take multiple small steps known as
diffusion module → take multiple small steps known as
Markov chain

Forward Process

Image → Noise → Noise → ... → Noisy Image → denoise → clean
→ denoised Image

$$q(x_t|x_{t-1}) = N[x_t; \sqrt{1-\beta}x_{t-1}, \beta I]$$

$$q(x_{t:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

Noise → gaussian noise

Reverse Process

$$V(q, \pi) = \int p_{\text{data}}(w) \log D(w) dw + \int p_z(z) \log (1 - D(G(z))) dz$$

diffusion model

$$q(x_t | x_{t-1}) = N \left[x_t, \sqrt{1 - \beta_t} x_{t-1}, \beta_t^{-1} \right]$$

mean
variance
output

$\beta = 0.0001$ to 0.02
to apply this forward formula multiple times:

$$\alpha_t = 1 - \beta_t$$

$$\bar{x}_t = \sum_{s=1}^t \alpha_s$$

hence

$$q(x_t | x_{t-1}) = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon$$

$$= N \left[x_t, \sqrt{1 - \beta_t} x_{t-1}, \beta_t^{-1} \right]$$

$$= \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \epsilon$$

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon$$

$$= \sqrt{\alpha_t \alpha_{t-1} \alpha_{t-2}} x_{t-3} + \sqrt{1 - \alpha_t \alpha_{t-1} \alpha_{t-2}} \epsilon$$

Calculating $P_G(x_0)$ is hard hence we can use
Variational lower bound

$$-\log P_\theta(x_0) \leq -\log (P_\theta(x_0)) + D_{KL}\left[q(x_1 \mid \cdot \mid x_0) \parallel P_\theta(x_{1:T} \mid x_0)\right]$$

diffusion models ~ Variational auto encoders

after computing this we get forward process

$$-\log P_\theta(x_0) \leq \log \frac{q(x_{1:T} \mid x_0)}{P_\theta(x_{0:T})}$$

$$P_\theta(x_{0:T}) = p(r_T) \prod_{t=1}^T P_\theta(x_{t-1} \mid x_t)$$

$$\text{forward} = q(x_{0:T}) = q(x_0) \prod_{t=1}^T q(x_t \mid x_{t-1})$$

$$\text{Reverse} = q(x_{0:T}) = \prod_{t=1}^T q(x_{t-1} \mid x_t) q(x_t)$$

$$q(x_{t-1} \mid x_t) = \mathcal{N}(x_{t-1} \mid \tilde{\mu}_t(x_t, t), \sigma_t^2)$$

reverse doesn't have closed form, so

$$\tilde{\mu}_t(x_t, t) \approx \frac{1}{\delta t} \int_{x_t}^{x_t} x_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_t$$

DIFFUSION MODELS
END