

# Index

TOPIC	PAGE
Perceptron	1
Train single neuron	2
Train multi-layer	3
<b>Back-propagation</b>	4
<b>Activation functions</b>	5
Drop-out	6
<b>ReLU</b> < leaky ReLU variants	7
weight initialization	8
<b>Batch normalization algo</b>	9
SGD, momentum, adagrad	10
adadelta, adam opt	11
gradient clipping	12
<b>softmax classifier</b>	12,13
Auto Encoders	14
<b>W2V</b> < CBOW skip-grams	15,16
<b>CNN Basics</b>	17
LeNET + maxpooling	19
<b>AlexNet , VGGNet</b>	20
gradient checking	21
Resnets	21
Inception network	22
Transfer learning	23
<b>RNN</b> + architectures	24
<b>LSTM and GRU</b>	26
Deep & Bi-dir RNN	28
GAN's	29
<b>Encoder - Decoder</b>	31
Word embedding	33
Calculate CNN params	35
Functional API Keras, FC $\leftrightarrow$ CONN 2D	34
Image segment n FUNET	37
LSTM time series	39
Keras LSTM params	40
Nesterov Accelerated grad	42
YOLO v3 - architecture	49
Attention models	50
Transformers	52
BERT	55
CANET - img segment	57
Seq-2-Seq models	59
Seq models + attention	60
char-RNN + TD Dense Lay	62
BERT-1,2,3 / differences	63
LIME - explainable AI	65
SHAP	67
RCNN , FRCNN	69
OLS linear regn	71
PCA + tSne	72
XGBoost + pro/con	73
AdaBoost	77

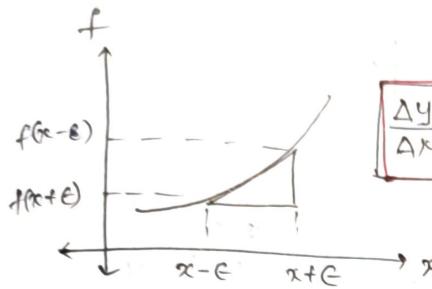
# ResNets Residual Networks

## Gradient Checking

\* Numerical approximation of derivative of  $f$ :

$$\lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$$

\* Actual der  $\rightarrow$  calculated by formula or backprop.



$$\frac{\Delta y}{\Delta x} = \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$$

because  $\epsilon$  is small,  
difference also is small

$$\text{grad-check} = \frac{\|d\mathbf{w} - d\mathbf{w}_{\text{approx}}\|_2}{\|d\mathbf{w}\|_2 + \|d\mathbf{w}_{\text{approx}}\|_2}$$

Euclidean distance

$\ll 10^{-7} \rightarrow \text{correct}$

$> 10^{-3} \rightarrow \text{false/not correct}$

$$\text{Ex: } \mathbf{W} = [W_1, W_2], \mathbf{x} = [x_1, x_2] \quad f(W_1, W_2, x_1, x_2) = w_1^2 x_1 + w_2^2 x_2$$

$$\frac{df}{dW_1} = 2 \cdot w_1 x_1 \rightarrow ① \quad \frac{df}{dW_2} = 2 \cdot w_2 x_2; \quad \left( \frac{df}{dW_1} \right)_{\text{approx}} = \frac{f(W_1 + \epsilon_1, W_2, \dots) - f(W_1 - \epsilon_1, \dots)}{2\epsilon}$$

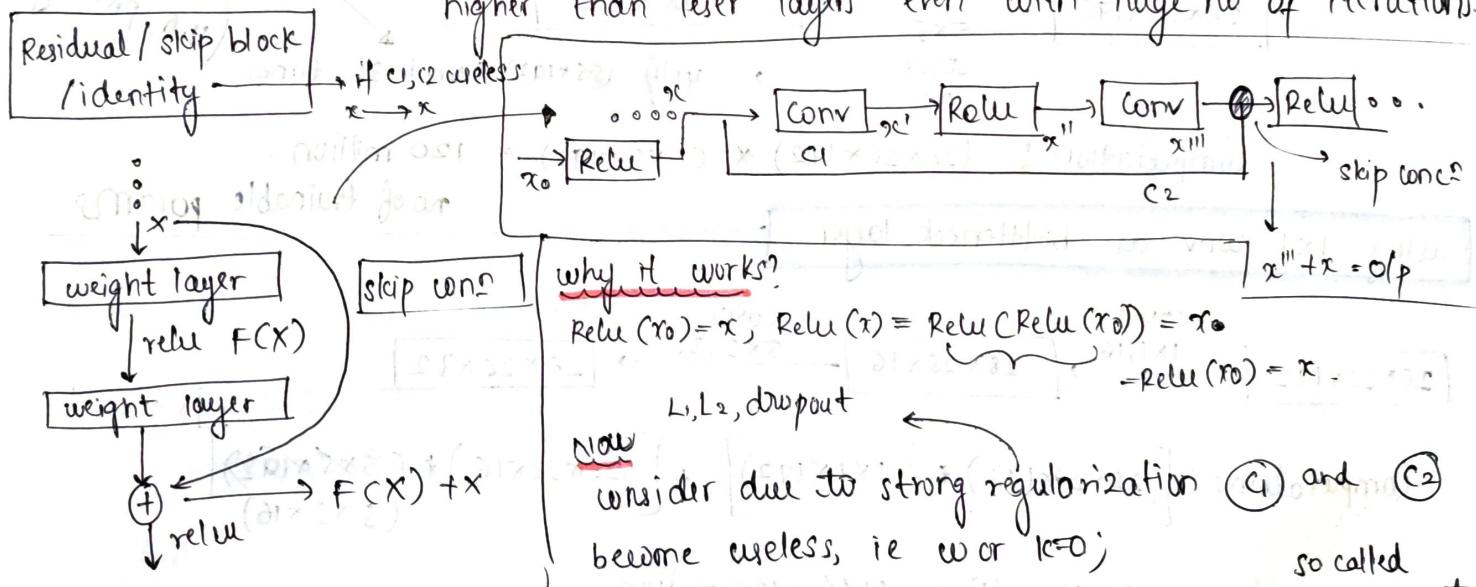
$$(d\mathbf{w}) = 2w_1 x_1 = 6 \rightarrow ①$$

$$(d\mathbf{w}_1)_{\text{approx}} = 5.9999999 \rightarrow ②$$

$$\text{grad-check} = \left( \frac{6 - 5.9999999}{6 + 5.9999999} \right) = 4.251268e^{-13}$$

$\therefore$  as  $\ll e^{-13} \rightarrow$  backprop  
is correct

**Resnets** :- motivation  $\Rightarrow$  when no of layers were increased, the loss was even with huge no of iterations



$$\text{ReLU}(\text{ReLU}(x)) = \text{ReLU}(x)$$

if +ve,  $\text{ReLU}(x) = \text{ReLU}(x) = x$   
-ve = 0

this wont propagate further

$\therefore \text{o/p} = x'' = 0$

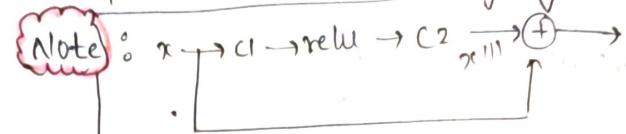
so previously,  $\rightarrow$  now,

$$\text{o/p} = x'' + x = x$$

so called skip connection

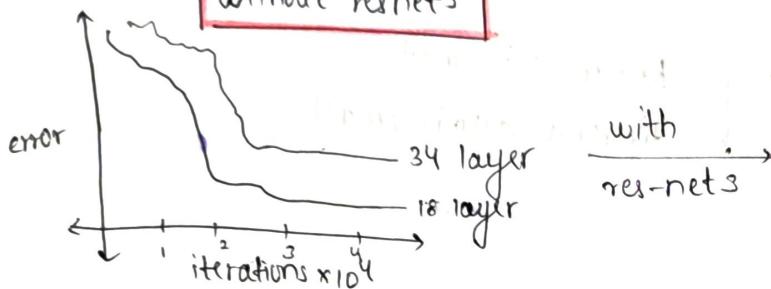
$\therefore c_1$  and  $c_2$  dont matter ie skip

- (22) key takeaway - guaranteed no decrease in performance, small increase in performance, adding additional layers would not hurt performance as they get skipped if they are found useless.
- \* if new layers are useful  
reg weights are non-zero.

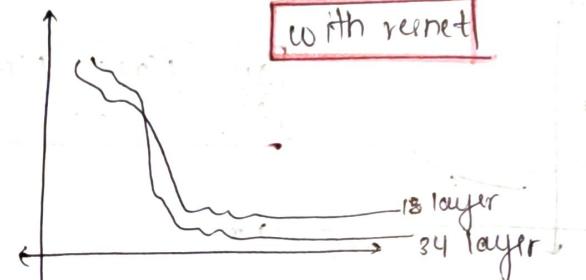


make sure  $x^{(1)}$  and  $x^{(2)}$  have same shape - use padding

**without resnets**

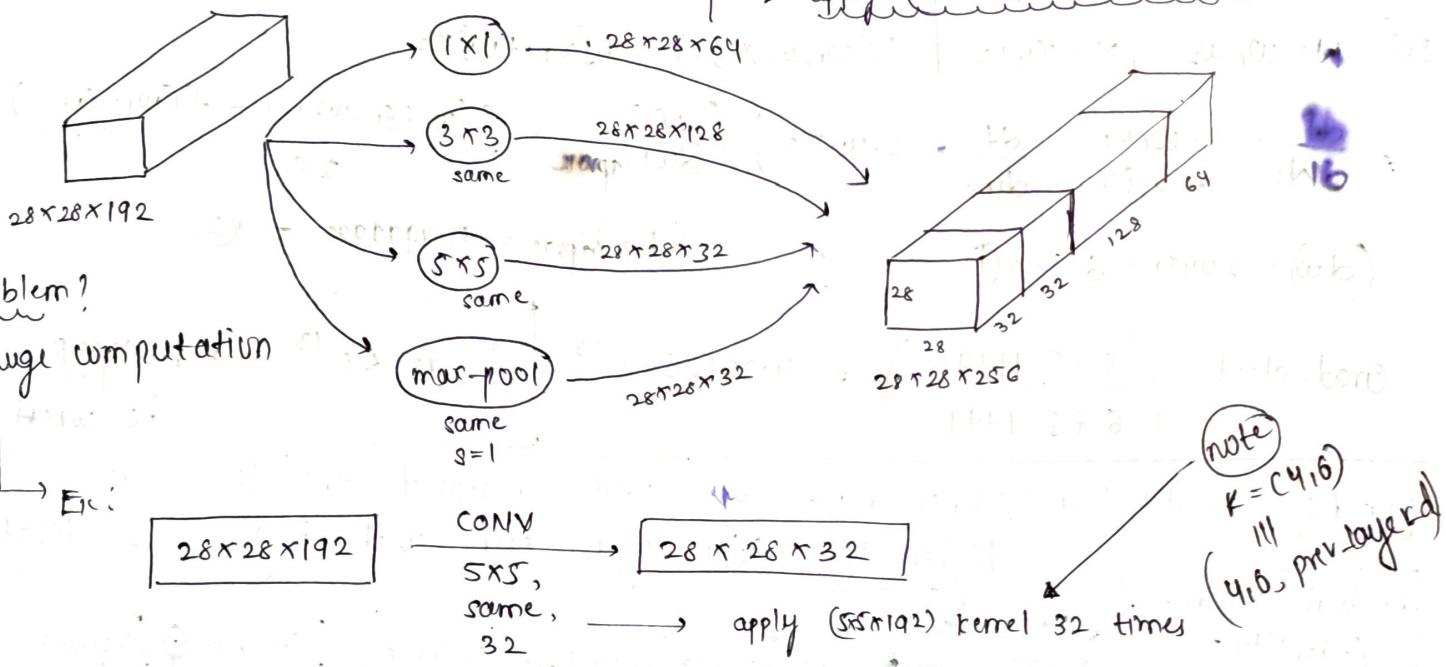


**with resnet**



### Inception Network

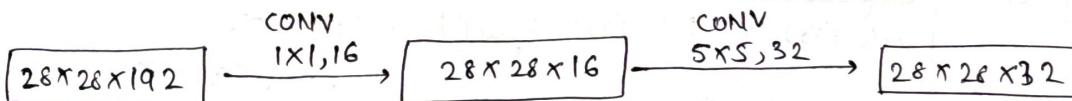
apply more than 1 type of kernels at once  
apply  $(1 \times 1 \times 192)$  kernels 64 times



$$\text{computation: } (28 \times 28 \times 192) * (5 \times 5 \times 192) = 120 \text{ million.}$$

no of trainable params

**using 1x1 conv as bottleneck layer**

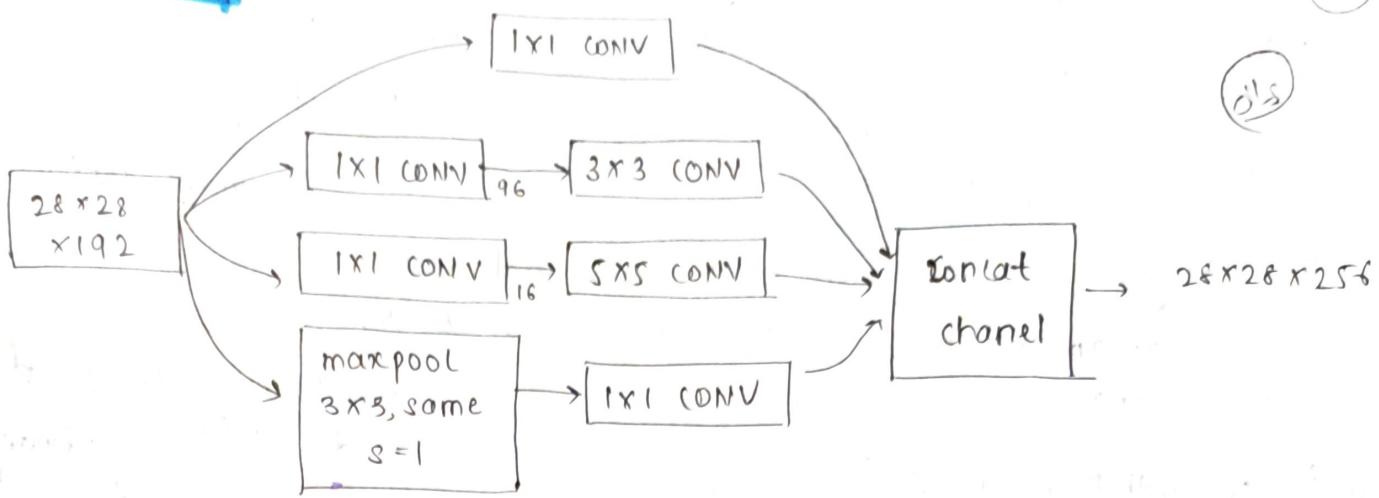


$$\text{computations: } [(28 \times 28 \times 192) * (1 \times 1 \times 192)] + [(28 \times 28 \times 16) * (5 \times 5 \times 32)] / (5 \times 5 \times 16)$$

$$= 12.4 \text{ million} <<< 120 \text{ million}$$

$\therefore (i/p \rightarrow \text{conv } 1 \times 1 \rightarrow \text{conv } 5 \times 5)$  is better computationally than  $(i/p \rightarrow \text{conv } 5 \times 5)$

## Inception module



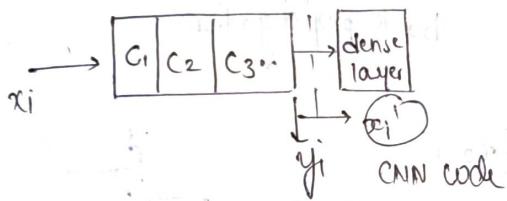
Q15

## Transfer Learning

→ using pretrained models

Ex: VGG16 trained with Imagenet dataset ; New dataset =  $D = \{x_i, y_i\}$

### Case 1 CNN codes



### Case 2 fine tune last layers

- Freeze most layers → initial layers are helpful
- fine tune only last few layers → use your dataset to train last layers.

they identify basics

### Case 3 tune all layers

- use VGG + ImageNet model as a initialization
- train the model using ur dataset

For  $\forall x_i$  in  $D$ ,

get  $(x_i)$  → vector obtained from last conv layer

$$\therefore D = \{(x_i, x_i^l, y_i)\}$$

↳ train any ml model

## Use cases

→ where to use?

↳ factors

→ similarity with Imagenet and size of dataset

**case 1**  $D \approx$  ImageNet  
IDL small

→ CNN codes

**case 3**  $D \approx$  ImgNet

IDL medium sized

**case 5**

|  $D$  | large

$D \neq$  ImgNet

size of  $D$

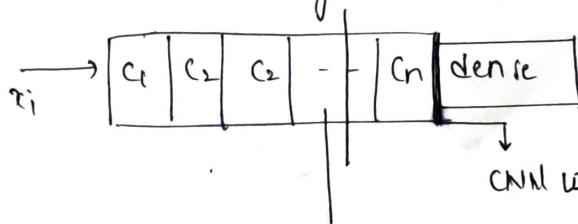
similarity to ImgNet

→ use pretrained model as init and train/tune whole model

**case 2**  $D \approx$  ImgNet  
IDL large

→ tune full n/w

**case 4** IDL is small  
 $D \neq$  Img Net



more generalized as  
cn will be useful iff similar

→ pick from later conv layers but not last

24

## Recurrent Neural Networks

• retains and leverages sequence info

Ex: Amazon fine food reviews

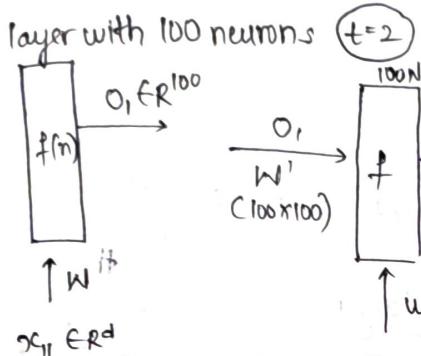
$$D = \{x_i, y_i\}$$

Each word is  $d$  dimensional vector

$$x_1 = \langle x_{11}, x_{12}, x_{13}, x_{14} \rangle$$

where  $d = \text{size of vocabulary}$

$$x_2 = \langle x_{21}, x_{22}, x_{23} \rangle$$

 $t=1$ 

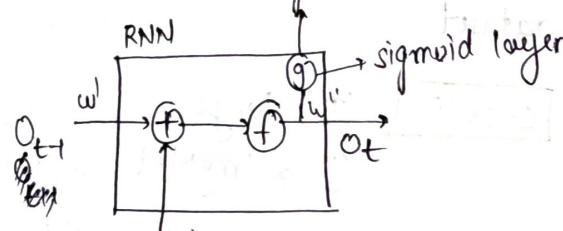
$$W: d \times 100$$

$f(m)$ : activation  $f^a$

$x_1 \rightarrow$   
initialise as  
zeros/random  
/any other

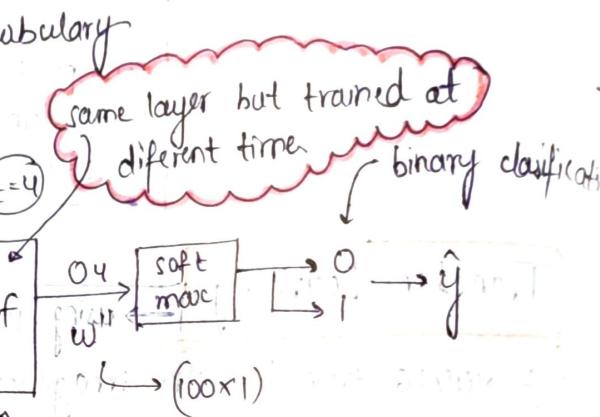
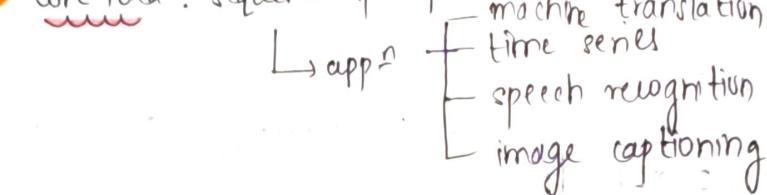
$$O_2: f(w' O_1 + w x_{12})$$

RNN layer

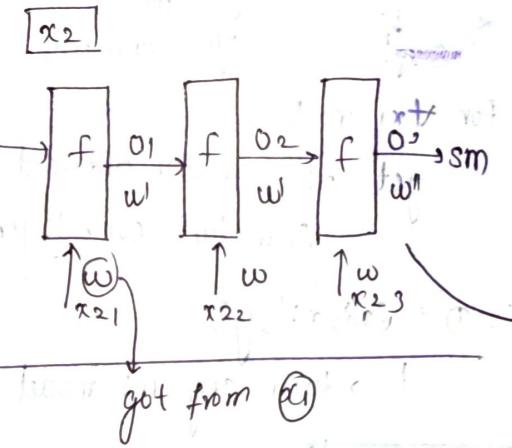
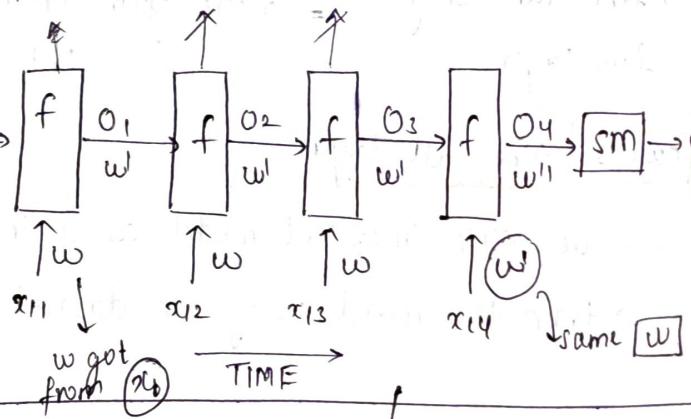


core idea: sequence of inputs

- ↳ app<sup>n</sup>
- machine translation
- time series
- speech recognition
- image captioning



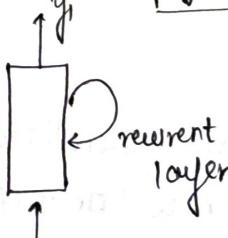
found out by forward and  
back propagation



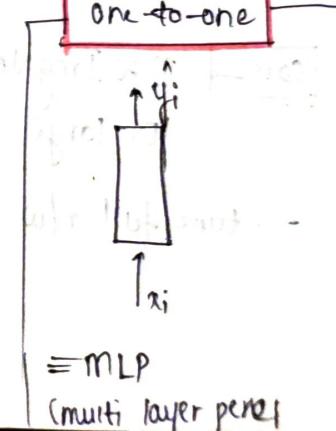
In the whole process we have only 3 matrices,  $w, w'$  and  $w''$

→ they are given different inputs according to time  $t$

$$\begin{aligned} O_t &= f(w x_t + O_{t-1} w') \\ y &= g(O_t w'') \end{aligned}$$

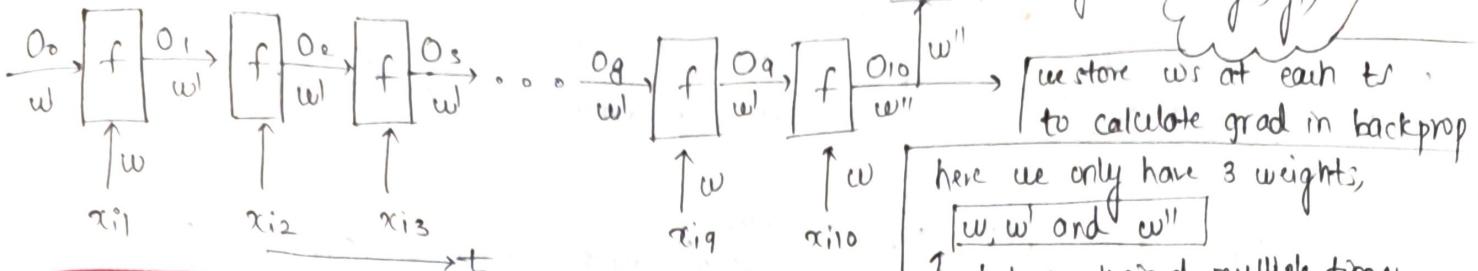


one-to-one



= MLP  
(multi layer perceptrons)

## Train $\rightarrow$ Backprop over time



## Forward prop

$$O_1 = f(O_0 w + x_{i1} w)$$

$$O_2 = f(O_1 w + x_{i2} w)$$

$$\vdots$$

$$O_{10} = f(O_9 w + x_{i10} w)$$

$$\hat{y}_i = g(O_{10} w'')$$

## Backward prop

$$\textcircled{1} \frac{\partial L}{\partial \hat{y}_i} \rightarrow \textcircled{2} \frac{\partial L}{\partial O_{10}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \rightarrow \textcircled{3} \frac{\partial L}{\partial w''} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial w''}$$

$$w_{t=10} = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial w} \rightarrow \left( \frac{\partial L}{\partial w} \right)_{t=1}$$

$$w'_{t=10} \approx \frac{\partial L}{\partial w'} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial w'} \rightarrow \left( \frac{\partial L}{\partial w'} \right)_{t=1}$$

## Problems with backprop over time

$\left( \frac{\partial L}{\partial w} \right)_{t=1}$  = multiplication and addition of multiple values

not bcoz of deep layers but. vanishing gradient Exploding gradient

## Types of RNN architectures

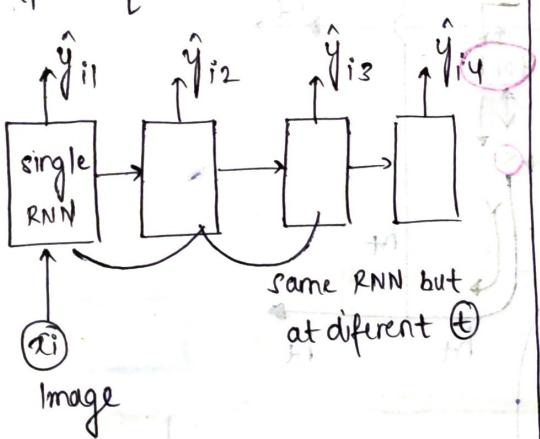
### ① many-to-one RNN

usu: sentiment analysis  
movie review, rating

### ② One-to-many RNN

inp: single

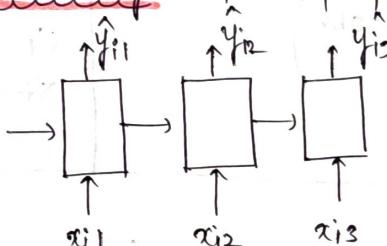
outp: sequence



Ex: Image captioning

### many-to-many (inp) = (outp)

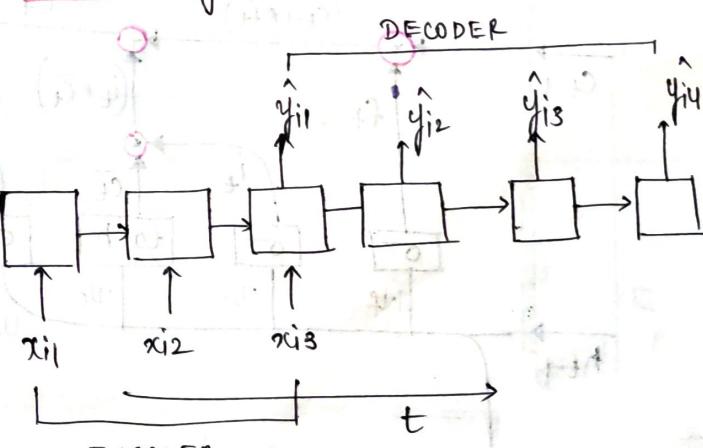
① same length : parts of speech



[parallel execution] can be done

### ② different length

#### Encoder-Decoder



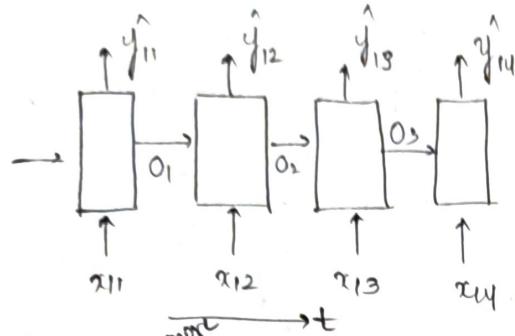
Ex: language translation = Eng → Spanish

a sentence in Eng may have only 8 words when with 10 word translated into Spanish

**LSTM**

why we need LSTM when RNN exists?

RNN can't take care of long term dependency  
ie later output depends on earlier input



→ Here,  
\*  $\hat{y}_{14}$  depends more on  $x_{14}$  and  $O_3$   
and depends very less on  $O_1$  and  $x_{11}$   
\* let's say we want  $\hat{y}_{14}$  to be more dependent  
on  $x_{11}$  which can't be accomplished

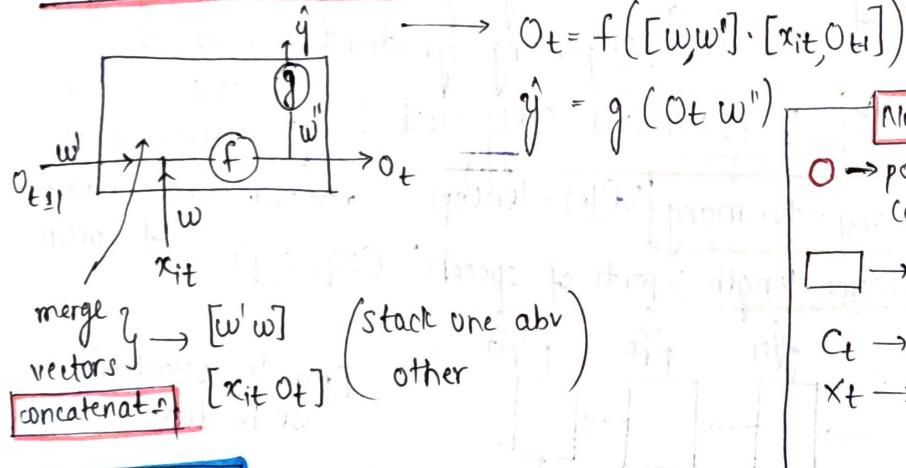
why  $\hat{y}_{14}$  depends more on  $x_{11}$ ?

→ there is lot of propagation &  
time b/w  $x_{11}$  and  $\hat{y}_{14}$

Ex- In some languages starting words might  
be of more importance when translating

**LSTM**

Improved simple RNN

**similarity**

$$c_t = o_t$$

$$x_t = X_{it}$$

$$h_t = \hat{y}_t$$

**Notations**

○ → point operation  
(apply to each)

□ → Neural network layer

$c_t$  → cell state.

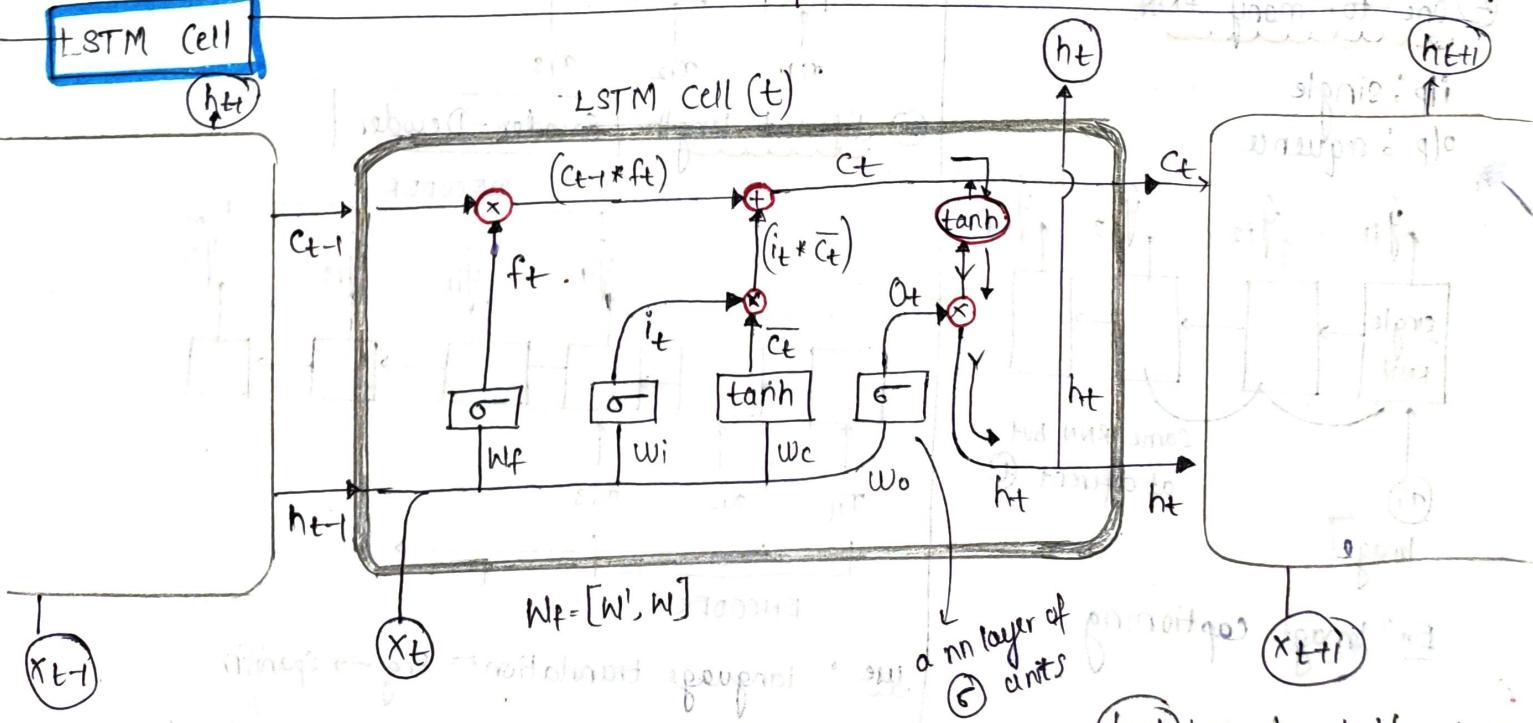
$x_t$  → ip point

$$\text{Ex: } [1 \ 2 \ 3] \otimes [1 \ 1]$$

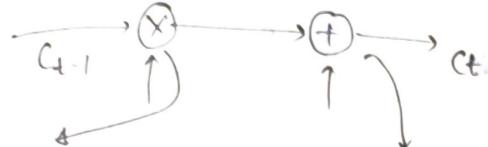
$$[1 \cdot 1 \ 2 \cdot 1 \ 3 \cdot 1]$$

→ concatenate

→ copy

**LSTM Cell**

skip / memory layer → if you op needs long term memory just skip .



how much to forget? how much more info to be added?

if i/p to  $\otimes$  =  $[1, 1, \dots]$  →  $[c_{t-1} \rightarrow c_t]$   
and i/p to  $\oplus$  =  $[0, 0, 0, \dots]$

cell / layer is skipped

### Forget part

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

$w_f$ : forget weights  
 $[w, w^T]$   
 $b_f$ : bias factor

$f_t$  → how much delta to forget

### Input gate

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

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

GRU

identity + forget + input gate

$$c_t = (f_t * c_{t-1}) + (i_t * \tilde{c}_t)$$

- \* few param → better with less data
- \* not good for long sequence

### Output gate

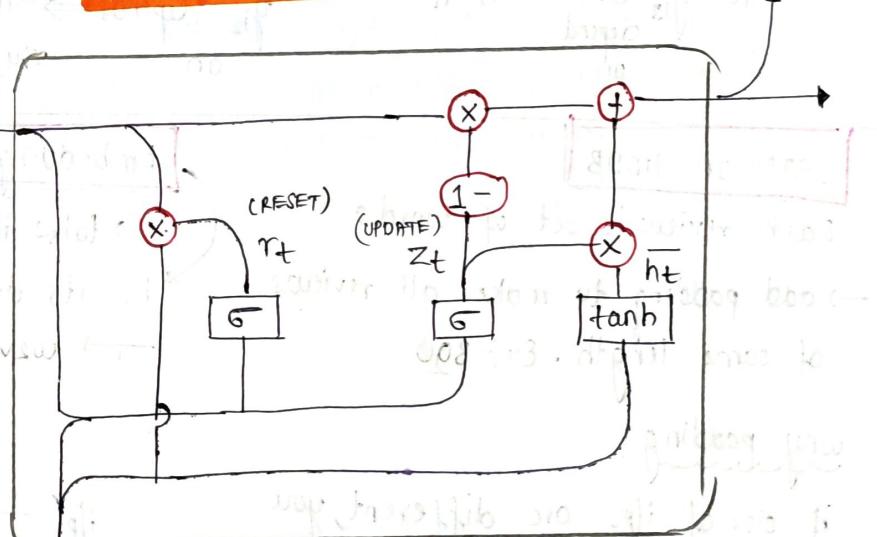
how much op to send?

$$o_t = \sigma(w_o \cdot [h_t, x_t] + b_o)$$

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

### GRU | Gated RNN unit

simplified LSTM



### GRU vs LSTM

→ LSTM have 3 gates

GRU have 2 gates - reset, update

∴ As fewer gates → fewer eqn

→ faster to train

$$z_t = \sigma(w_z \cdot [h_{t-1}, x_t])$$

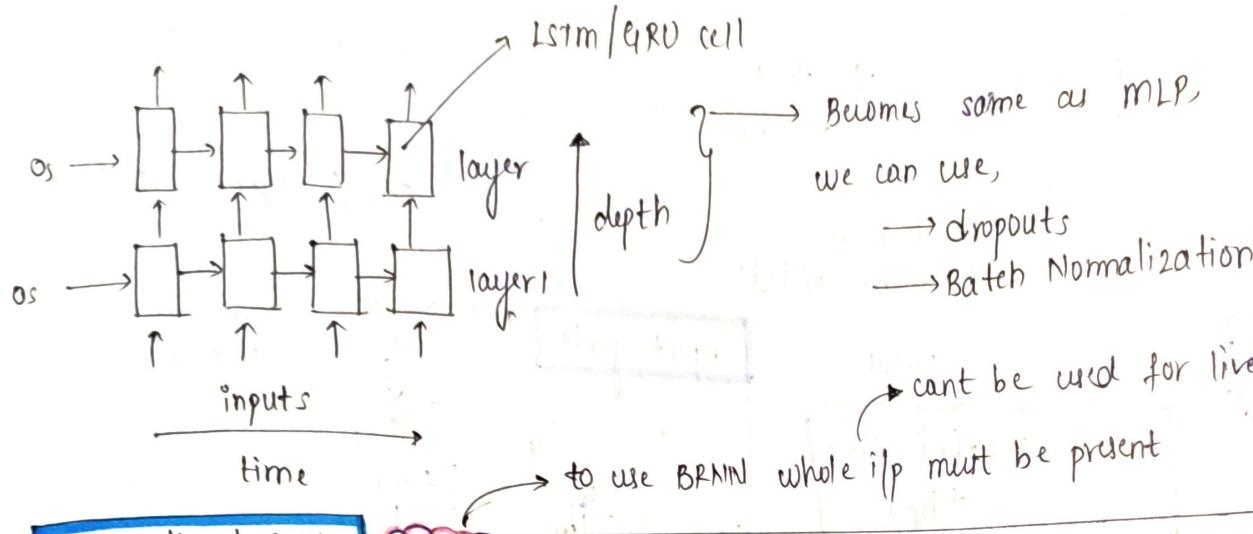
$$r_t = \sigma(w_r \cdot [h_{t-1}, x_t])$$

$$h_t = \tanh(w \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + (z_t * \tilde{c}_t)$$

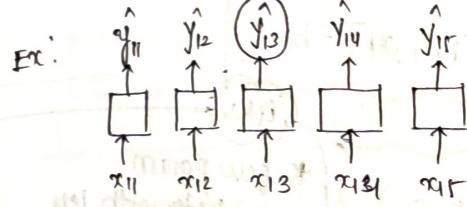
$[1, 1, \dots]$   
- [ip vector]

## (28) Deep RNN

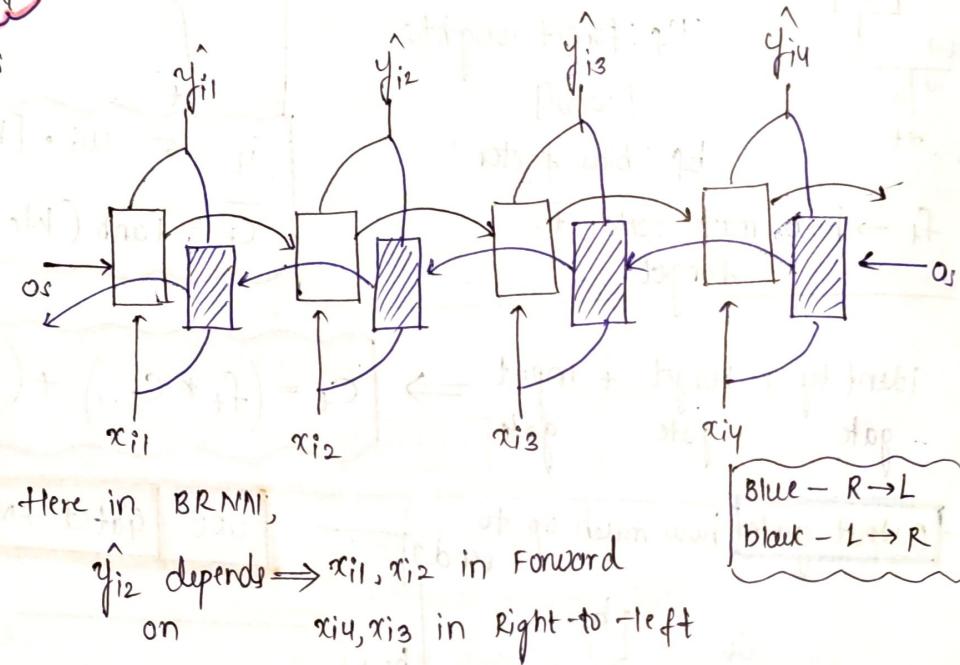


### Bi-Directional RNN

what if current O/p also depends on future i/p?



but can't depend on future i/p  
ie  $\hat{y}_{13}$  can't depend on  $x_{14}, x_{15}$



### LSTM on IMDB

Each review is set of words

→ add padding to make all reviews of same length. E.g. 390

why padding?

→ if size of i/p's are different, you can only pass one i/p at a time  
= stochastic i/p → very small

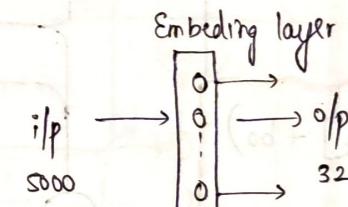
→ if i/p size is same, you can parallelly process inputs like in batch i/p → faster

### Embedding layer

→ takes i/p → returns vector

how its different from w2v?

→ w2v is static but here you can train



No of trainable params

$$= 5000 \times 32$$

$$= 160K \text{ params}$$

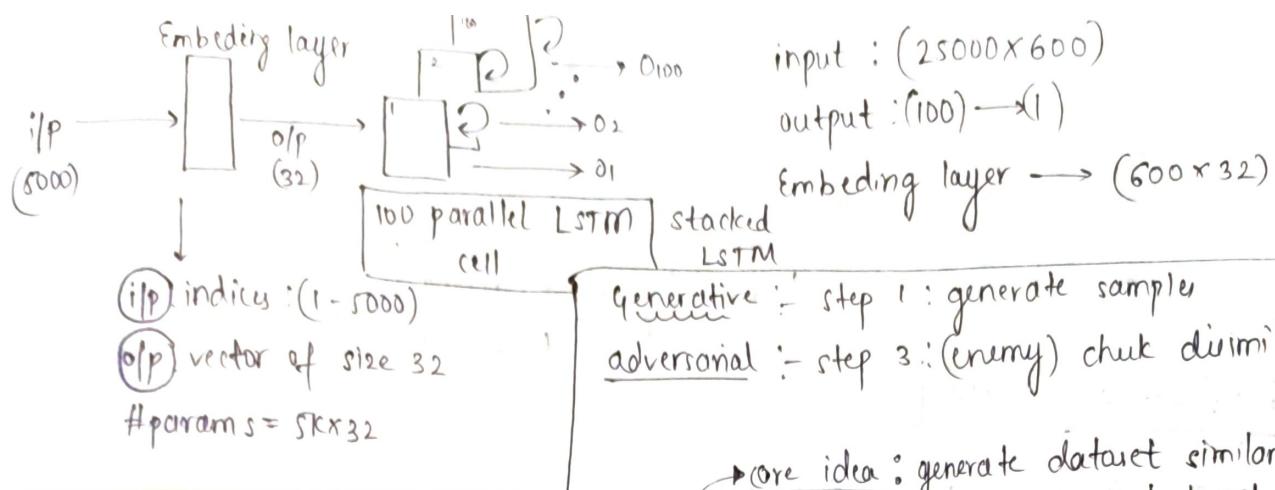
### No of train params in LSTM

$$m = \text{no of inputs} = 32$$

$$n = \text{no of o/p's} = 108$$

$$\# \text{params} = 4(m^2 + n^2 + mn)$$

→ Biases



## Generative Adversarial Networks

Simpler task: 1D scalar  $\rightarrow$  step 3: check how similar  $D$  and  $D'$

$D$  = {set of heights}

generate  $D'$  similar to  $D$

$\rightarrow$  s1: generate PDF of  $D$ ,  $P(D)$

\* find which dist<sup>r</sup> it is

ex: it is normal distr

$P_0(H)$   $\rightarrow$  vector which is distr params  $\mu, \sigma$

$$P_0(H) = N(\theta = [\mu, \sigma])$$

$\rightarrow$  s2: generate samples randomly from  $P_0(H) = D'$

① use statistical tests

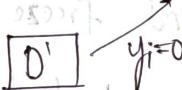
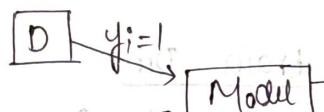
(prob model) JS-dist

KL-div

problem: it works only

with small dimension

② model approach



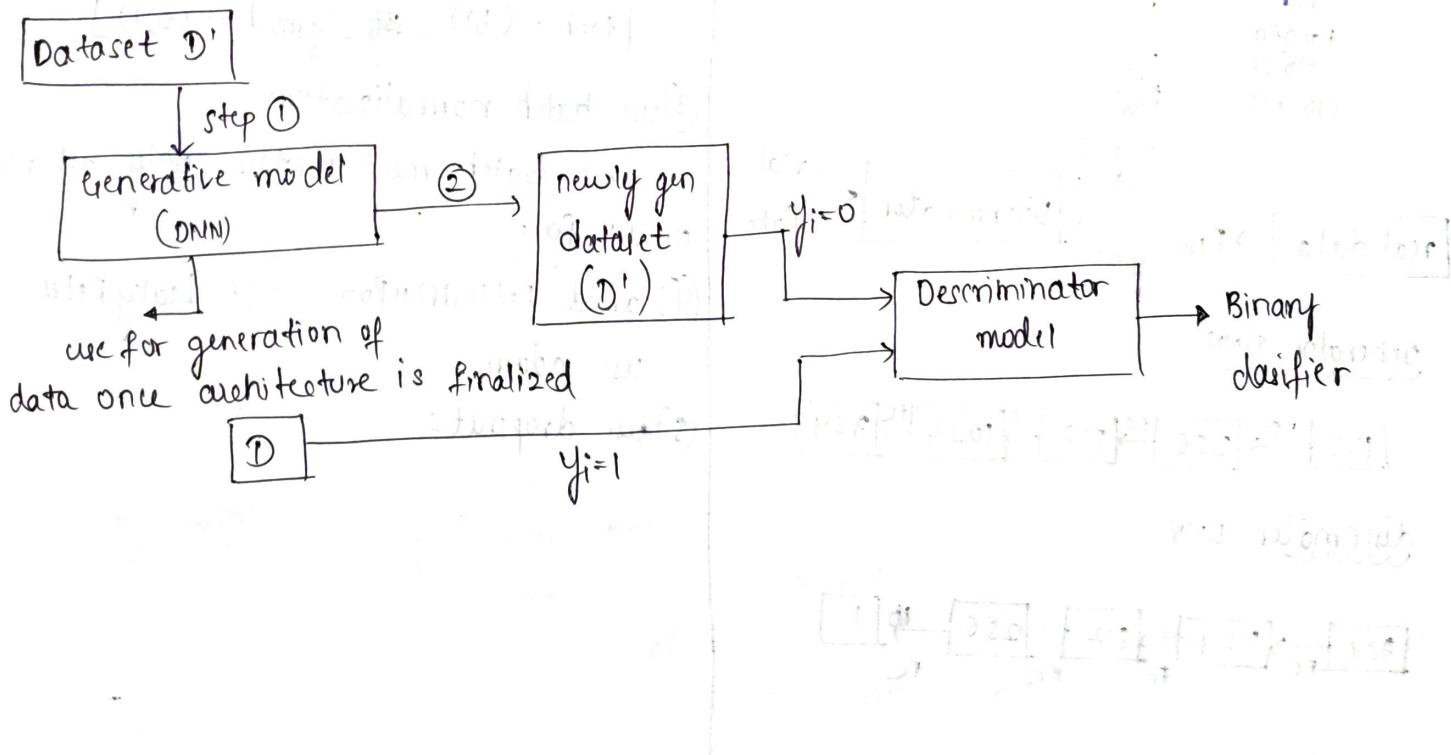
$\rightarrow$  if loss is less, then  $D$  and  $D'$  are separable

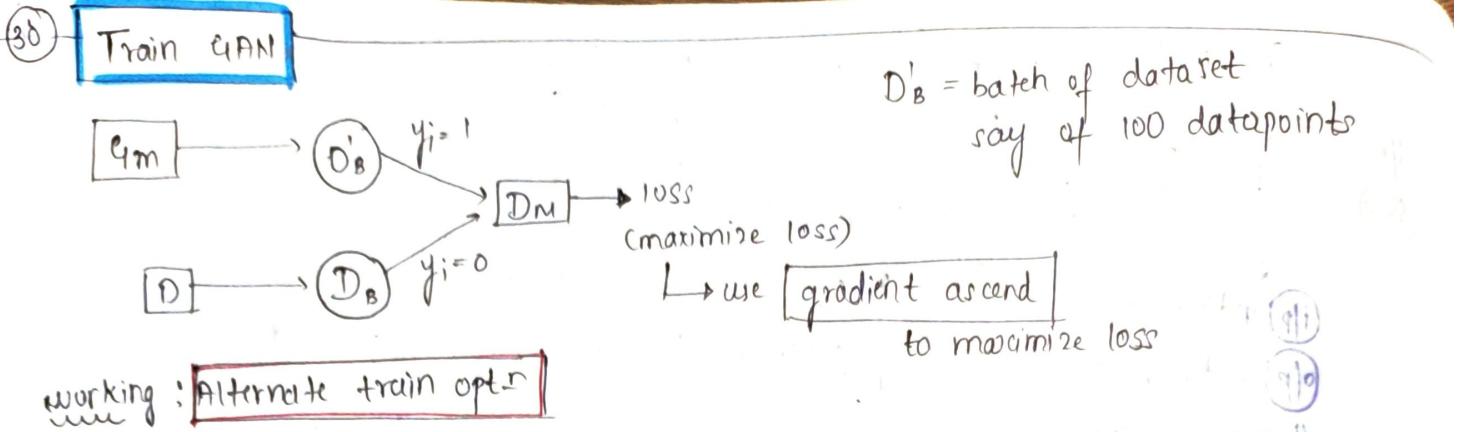
$\therefore D$  and  $D'$  are not similar

$\rightarrow$  if loss is large ie bad performing model

$\rightarrow D$  and  $D'$  are similar

Note: probabilistic models don't work for multivariate analysis. so use mH/DNN





### step 1 Keep $D_m$ fixed and train $G_m$

→ for  $D_m$  initialize we a random sample as train data  $\in \mathbb{R}^d$  normal dist<sup>r</sup> vector

train the model for it

\* Here  $G_m$  does up-sampling; i.e.  $(100d \rightarrow 784d)$  MNIST

\* Initially  $D'_B$  generated would be very-very random and thus different from  $D$ . So loss will be very less which is bad for us

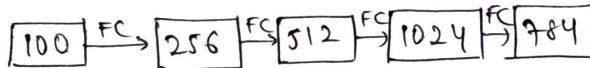
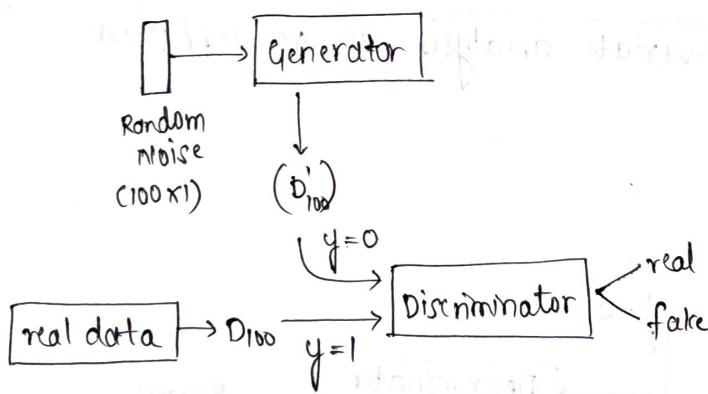
\* we need to find and maximize loss using gradient ascend

### step 2 Keep $G_m$ fixed and train $D_m$

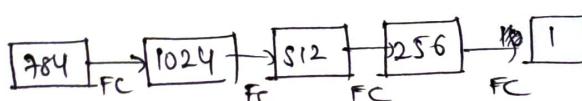
\* only update weights in  $D_m$  and freeze weights in  $G_m$

→ Alternatively do ① and ② until satisfactory loss is obtained

### GAN for MNIST



discriminator DNN



### Tips

① use normal distribution for sampling

② Normalize images b/w 1 and -1

→ use tanh as last layer of generator

[tanh - (-1, 1), sigmoid - (0, 1)]

③ use batch normalization

Each batch must contain only all real or all fake

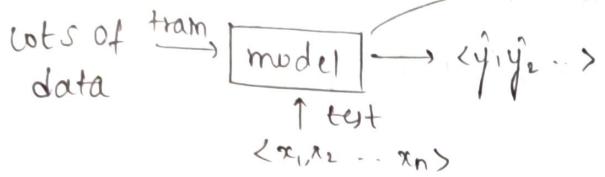
④ avoid ReLU, MaxPool ; use LeakyReLU

use adam

⑤ use dropouts

## Encoder - Decoder

Probabilistic repr of model

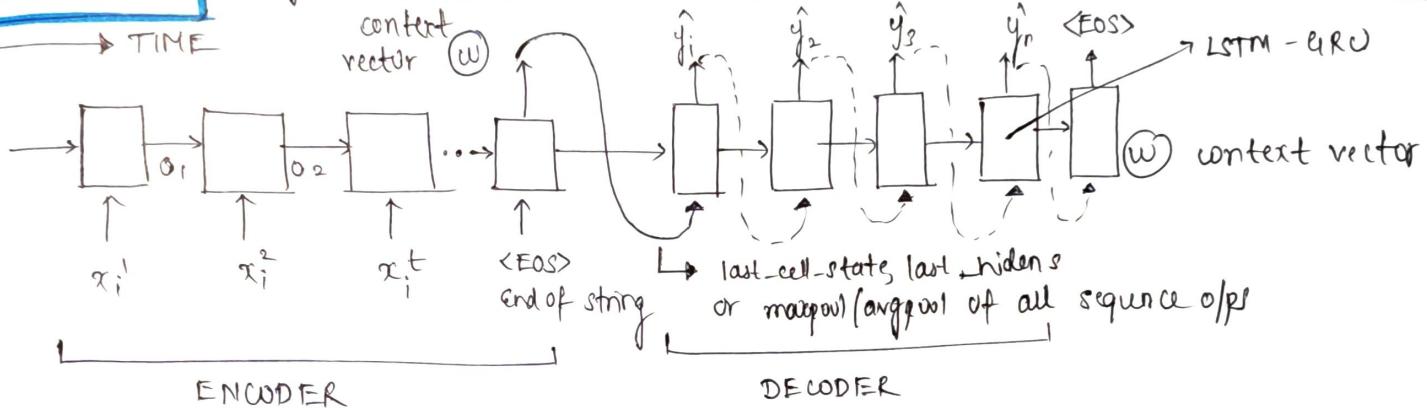


$$\text{argmax } p(y_1, y_2, \dots | x)$$

$$\hat{y} = \text{argmax } p(y | x)$$

find  $y$  which has max probability wrt data pt  $x$

### SEQ-2-SEQ $\xrightarrow{\text{Ilya Sutskever et al}}$

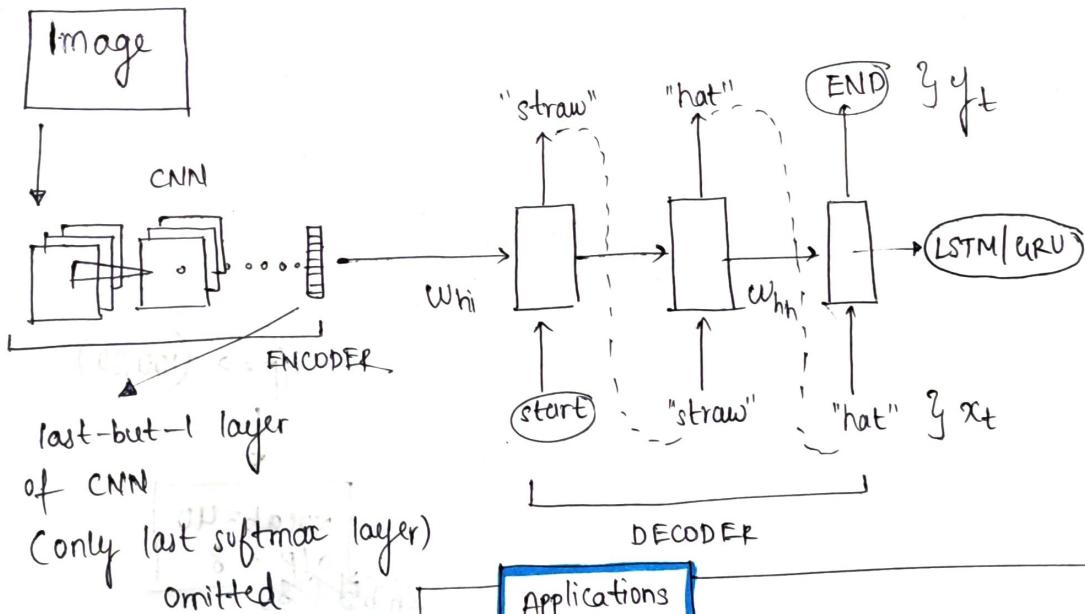


what is output sentence is long? [cho et al]

→ context vector is fed to all decoder cells, rather than only first decoder cell

but LSTM can't take 2 if at time stamp  
↓  
used modified LSTM

### Image-2-caption $\xrightarrow{\text{Karpathy et al}}$



### Applications

- ① translation of text
- ② Email auto complete and smart replies
- ③ Image captions and descriptions
- ④ statically finding code errors

## Char RNN

- \* many-to-many RNN → given large wrpw of sequence of data
  - \* one o/p wruponding to each i/p predict which char follows given char most probably
- (train) let sequence be  $c_1, c_2, c_3 \dots c_n$
- i/p if i/p is  $c_2$  if i/p is  $c_2$  learn this  
o/p should be  $c_3$  if i/p is  $c_2$  learn this

So if i/p is  $c_2$  then o/p will be  $c_3$   
and so on till the end of sequence

100 80

Input shape: (batch\_size, input\_length, input\_dim)  
For example: (100, 80, 1)

Output shape: (batch\_size, output\_length, output\_dim)

For example: (100, 80, 8)



i/p  $\Rightarrow (100, 80)$



~~Embed 8~~  
vocab = 40  
o/p = 8

Embed (40, 8, input\_len = 20)

Input shape: (batch\_size, input\_length, input\_dim)

Output shape: o/p  $\rightarrow (\text{None}, 100, 8)$

# params = vocab size \* 8

# Word Embedding Layer

DL + NLP

## word embedding?

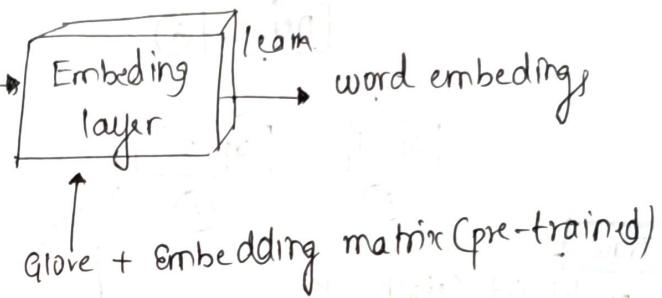
→ represent word with vector → TF-IDF, BOW, Word2Vec etc

word-embedding-layer → let nw learn embeddings by itself

## Process

sentence (seq of words) → Integer representation

Ex: docs = [ "sentence is one ...",  
"sentence - 2 ..."]



① Tokenizer (word → int) {{"word": index}}

→ t = Tokenizer().fit\_on\_texts(docs)

x\_tr = t.fit

→ x\_tr = t.texts\_to\_sequences(docs)

→ x\_tr is vector of integers.

Each word is represented by value indicating its index in vocabulary.

[t.word\_index]

let len(t.word\_index) = 300 (300 words)

padding make all sentences have same len

→ maxLen = for i in x\_tr find largest i

→ x\_tr = pad\_sequence(x\_tr, padding='post',  
dtype=int, maxlen=maxLen)

∴ we have each sentence/data point with len = maxLen

Ex → [1, 0 1 0 0 0 0 ... ]

use glove-vector / pretrained model

→ assume ur glove vector has 200 dim → vector for each weight

Create Embedding matrix

	0	1	2	...	200
0					
1					
2					
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					
14					
15					
16					
17					
18					
19					
20					
21					
22					
23					
24					
25					
26					
27					
28					
29					
30					
31					
32					
33					
34					
35					
36					
37					
38					
39					
40					
41					
42					
43					
44					
45					
46					
47					
48					
49					
50					
51					
52					
53					
54					
55					
56					
57					
58					
59					
60					
61					
62					
63					
64					
65					
66					
67					
68					
69					
70					
71					
72					
73					
74					
75					
76					
77					
78					
79					
80					
81					
82					
83					
84					
85					
86					
87					
88					
89					
90					
91					
92					
93					
94					
95					
96					
97					
98					
99					
100					

+ each word in vocab,  
have a v-dim vector

for word, idx in t.word\_index:  
embed[idx] = glove.get(word)

## Embedding Layer

total no of words in vocab / len(em)

e = Embedding(input\_dim=len(vocab),

output\_dim=200,

input\_length=maxLen,

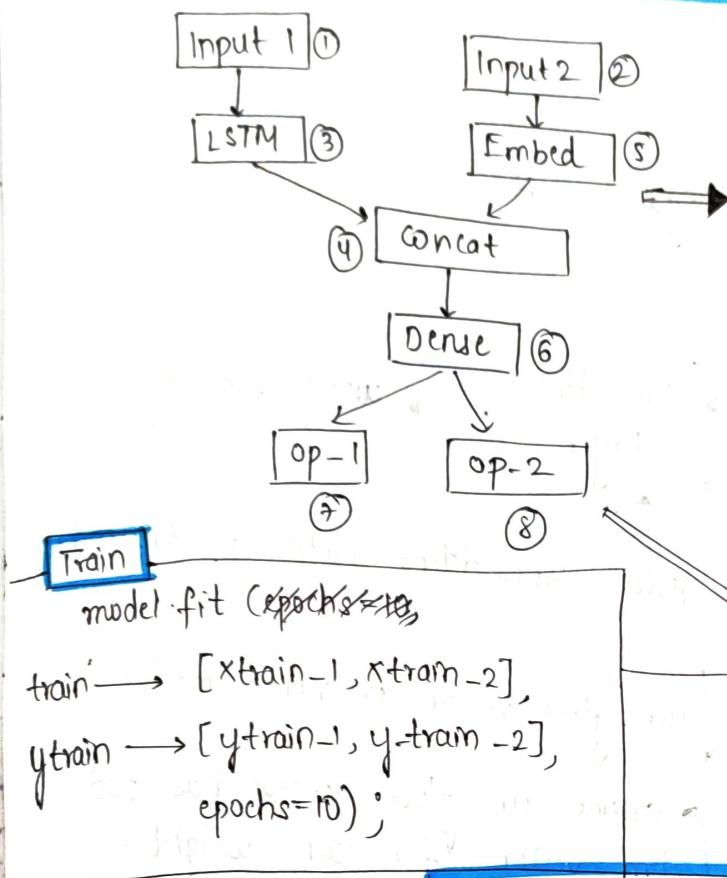
length of each training data.

trainable=False,

} for pretrained glove vector

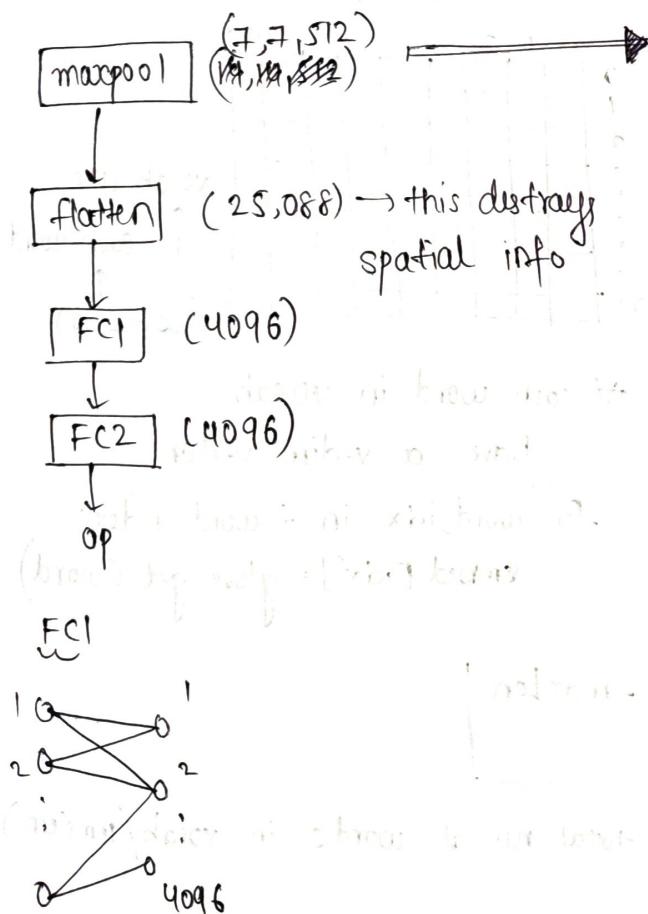
weights=[Embed\_matrix])

## Functional API Keras

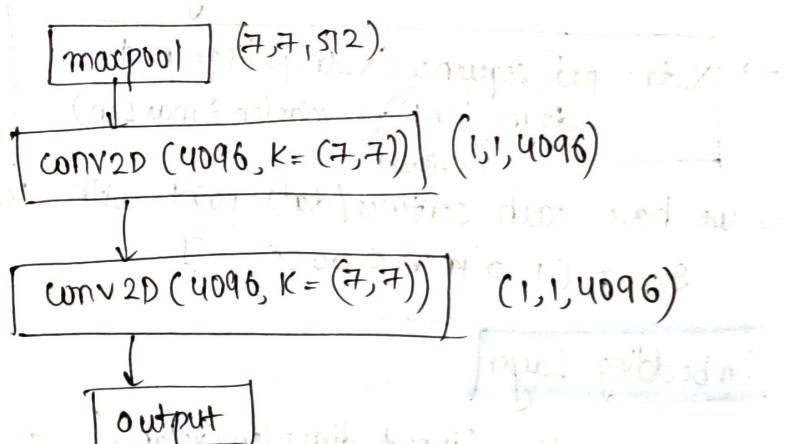


①  $\text{Inp-1} = \text{Input}((128,))$   
 ③  $\text{LSTM} = \text{LSTM}()$  ( $\text{Inp-1}$ )  
 ②  $\text{Inp2} = \text{Input}((10,))$   
 ⑤  $\text{Embed-L} = \text{Embedding}()$  ( $\text{Inp2}$ )  
 ④  $\text{concat} = \text{concatenate}([\text{LSTM}, \text{Embed-L}])$   
 ⑥  $\text{Dense} = \text{Dense}()$  ( $\text{concat}$ )  
 ⑦  $\text{op-1} = \text{Dense}(10, \text{softmax})$  ( $\text{Dense}$ )  
 ⑧  $\text{op-2} = \text{Dense}(2, \text{softmax})$  ( $\text{Dense}$ )  
  
 $\text{model} = \text{Keras}.\text{Model}([\text{Inp-1}, \text{Inp2}], [\text{op-1}, \text{op-2}])$

## Replace FC layer by Conv2D



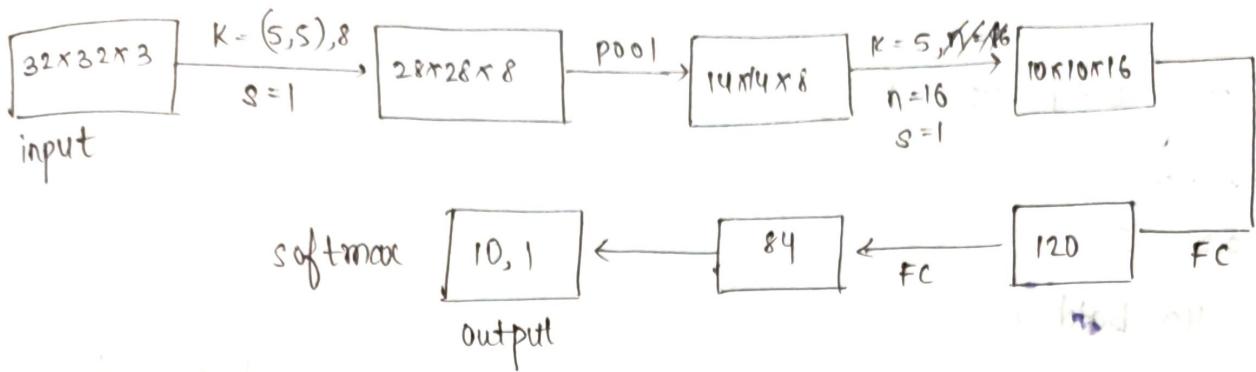
replace FC by CONV2D  
for FC1 we need  $4096 = (1, 1, 4096)$   
 $(7,7,512) \xrightarrow[s=1, f_n=4096]{K=(7,7)} (1, 1, 4096)$   
 $(\frac{(7-7+1)}{1}, \frac{(7-7+1)}{1}, 4096)$



25088 → number of channels in input  
 $\therefore 4096 * 25088$  com  
other parameters have to change  
ratio esp. bottleneck etc.

## Calculate CNN params

(35)



### \* conv layers

$$(n, n) \rightarrow \left( \frac{n-K+2P+1}{S}, \frac{n-K+2P+1}{S}, \text{nof filters} \right)$$

$$\# \text{params in conv layer} = \left[ \left( K \times K \times \text{nof filter in previous layer} \right) + 1 \right] * \text{nof filters in current layer}$$

$$\# \text{params in FC} = \text{FC(Fcn)} = \# \text{neurons in last layer}$$

\*  $\# \text{neurons in current FC} + \text{FCn}$

Bias term

### Count

layer	shape	# params	activation size	
Input	$(32 \times 32 \times 3)$	0	$32 \times 32 \times 3 = 3072$	
CONV1 $K=5, n=8$	$(28 \times 28 \times 8)$	$\left( \frac{5}{28 \times 28 \times 3} + 1 \right) 8 = 608$	$608 \times 6272$	$\begin{matrix} \uparrow \\ [(5 \times 5 \times 3) + 1] * 8 \\ \downarrow \end{matrix}$ prev-kernel bias cur kernel
POOL1	$(14 \times 14 \times 8)$	0	$1568$	
CONV2 $K=5, n=16$	$(10 \times 10 \times 16)$	$\left( \frac{5}{10 \times 10 \times 8} + 1 \right) 16 = 3216$	$1600$	
FC1 $(120)$	$(120, 1)$	$(10 \times 10 \times 16 \times 120) + 120 = 192120$	$120$	

### \* Embedding layer (vocab\_size, output\_dim)

o/p  $\rightarrow (\text{None}, \text{ip-size}, \text{output})$

#  $\rightarrow (\text{batch-size} * \text{output\_dim})$

36

Batch size

$\rightarrow \text{model}.\text{fit}(x, y, \text{batch\_size} = b);$

Let  $\text{len}(\tau) = 200$ .

batch-size = 5

$\therefore$  we get 40 batches

∴ we get 40 batches  
training → Epoch 1 → divide data into 40 batches → pass batch 1 & update weight  
→ pass batch 2 & update weight.

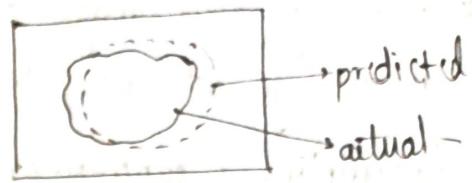
∴ In one epoch, weights are

update  $\left(\frac{\text{size of dataset}}{\text{batch-size}}\right)$  no of time. ie 40 in our case

# Image Segmentation

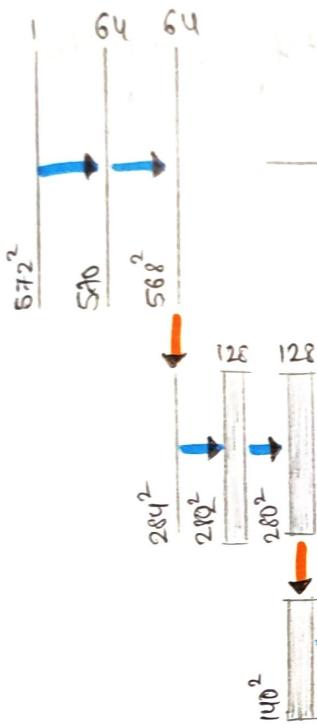
→ divide into areas / segments

→ measures :- ① jaccard similarity =  $\frac{R \cap R'}{R \cup R'}$   
(100)



② cross entropy = per pixel comparison

## U-net



COPY and CROP

crop (200x200)

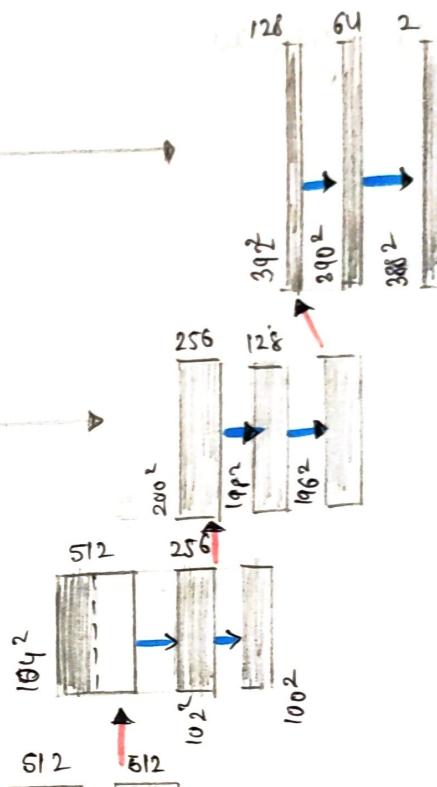
COPY and CROP

crop (64x64)

COPY AND CROP

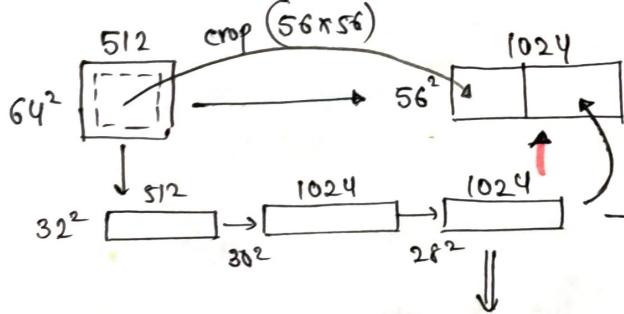
crop (56x56)

copy



- conv 3x3, ReLU
- max pool (2x2)  
stride = 2
- upconv (2x2)  
upsample + conv (2x2)

## upsample base layer



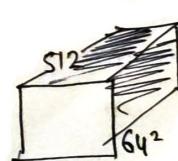
upsample

① upsample : duplicate rows and columns

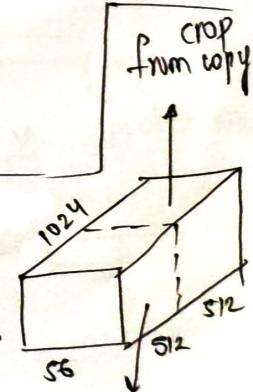
$$28 \times 28 \times 1024 \rightarrow 56 \times 56 \times 1024$$

② conv2D (f = 512, (2x2))  
(56x56x512)

$$(28 \times 28 \times 512) + (56 \times 56 \times 512)$$



crop  
from copy



from upsample

Q8

### ⑧ General idea

- \* as we go down, image size halves and no. of filters doubles  
ie we go into pixel level as we go down
- \* (upsampling + copy-n-crop) acts like a residual unit ie get information from lower level as well as from copy-crop unit
- \* Here o/p =  $(388 \times 388 \times 2)$   
i/p - grayscale
  - 2 masks are predicted  
ie binary segmentation to identify 2 regions

moving window average

consider you have data for  $t$  days

future data ( $t+1$ ) depends on average of previous  $K$  days

other ways

mean avg

median

exponentially weighted

dataset - weather temp each 10 minutes given

Univariate prediction

assume  $K=20$  we take window size = 20

\* Normalize all

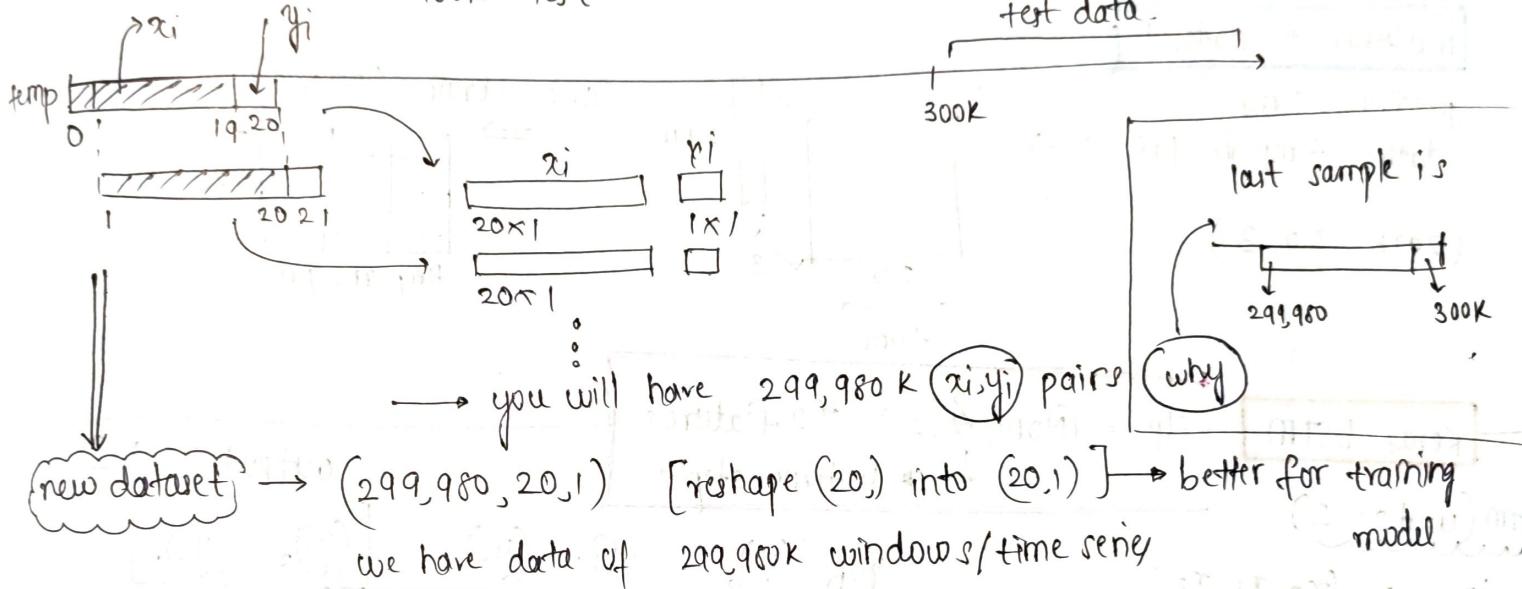
data

- only train data

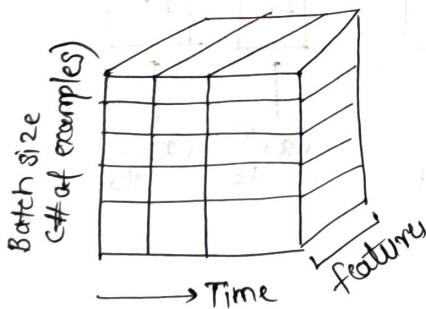
dataset creation :-

300K - train

100K - test

LSTM training diagram

\* LSTM expects data in this format → [batch, time, features]



height ( $y$ ) = samples

$x$  → time series

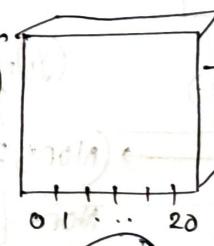
depth → features  
(univariate 1)

loss = mae

mean absolute error

our data

299K

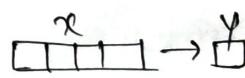
multivariate - single step

→ given a window's predict only one data at  $t_1$  -  $t_n$

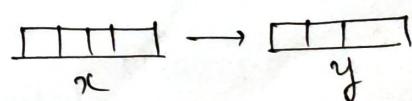
\* single step → predict only one pt at future

\* multi step → predict over multiple points

train



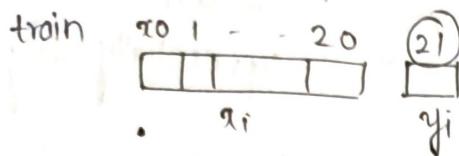
predict one



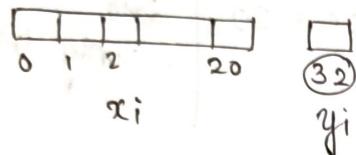
40

## Dataset generation hawks

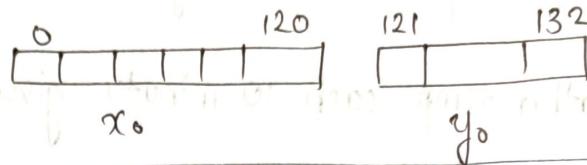
① predict next hour itself



② predict next 12 hours



② predict next 12 hours data using past 120 hours → multi step



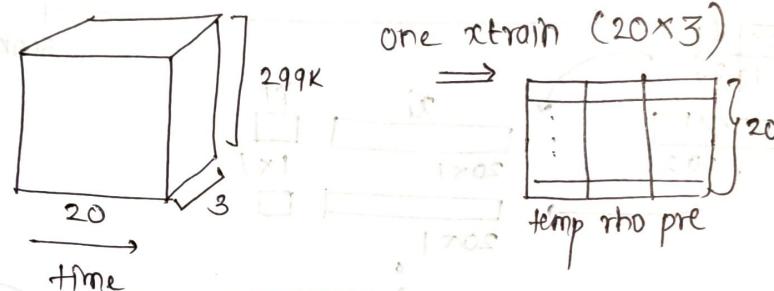
$$\begin{aligned} x_{\text{train}} &= (299K, 20, 1) \\ y_{\text{test}} &= (299K, 12, 1) \text{ or } (299K, 12) \end{aligned}$$

## Multivariate dataset

predict - temp

data - temp, rto, pre  $\Rightarrow$

(299K, 20, 3)



Keras LSTM

inp = (None, 3, 2)  $\rightarrow$  2 features  
 $\downarrow$  4 time-steps

ip  $\rightarrow$  ( $x_0 \ x_1 \ x_2$ )

②  $\rightarrow$  dim = no of units



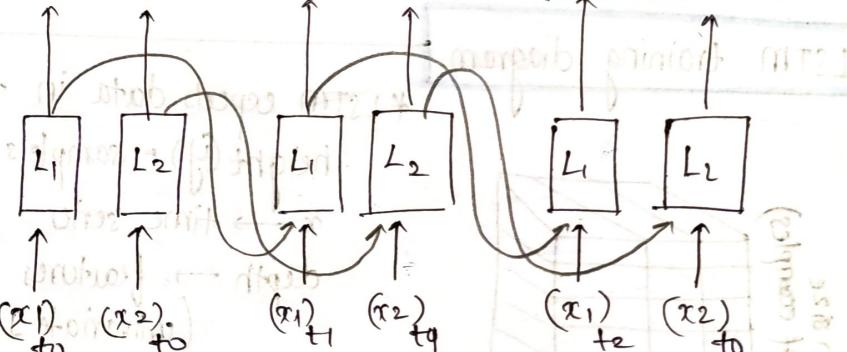
$(y_1)_2 \ (y_2)_2$

$(y_1)_0 \ (y_2)_0$

$(y_1)_1 \ (y_2)_1$

output  $\rightarrow$  ①

$(y_1)_2 \ (y_2)_2$



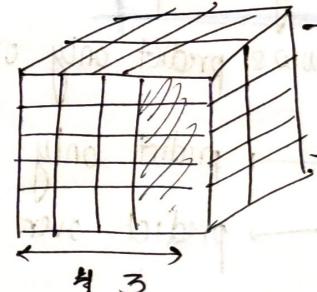
(None, 2, 2)  $\rightarrow$  (None, 2)

None  $\rightarrow$   $y_1, y_2$

units = 2  $\rightarrow$  generate 2 data for each

point  $\rightarrow$  one o/p from last unit

None, 2, 2

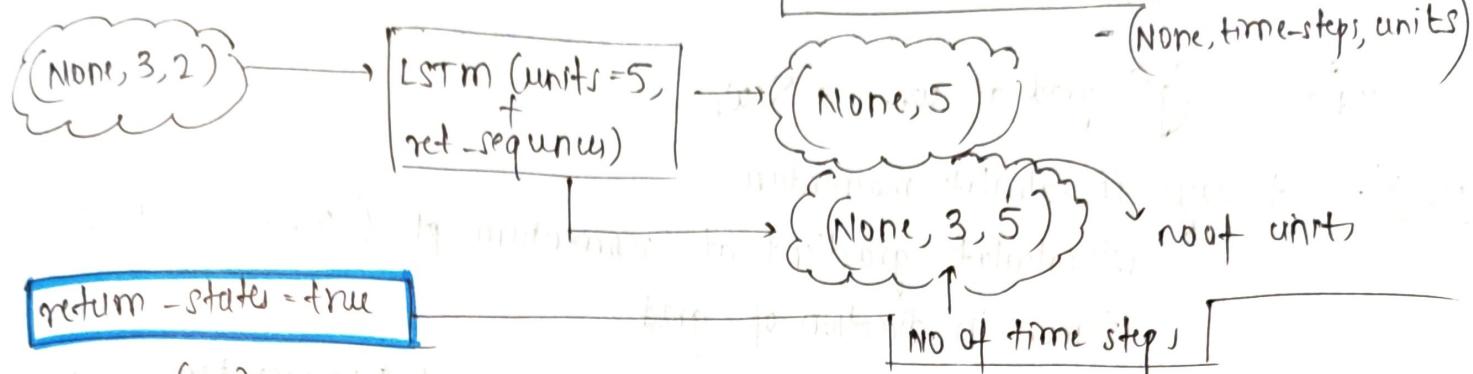


return-sequence = true

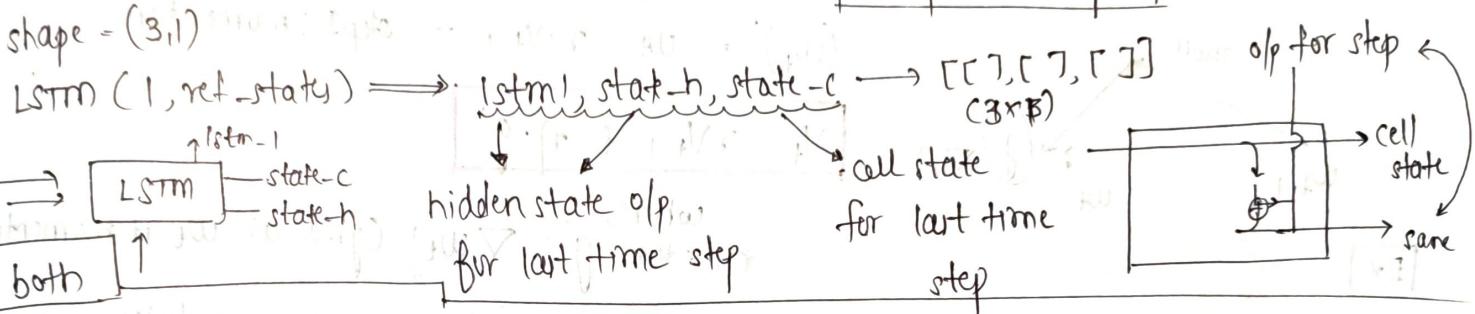
41

+ don't only return last-state o/p →  $[[y_{10}, y_{20}], [y_{11}, y_{21}], [y_{12}, y_{22}]]$   
but o/p for each state and time step

(21)  
2



return - states = true



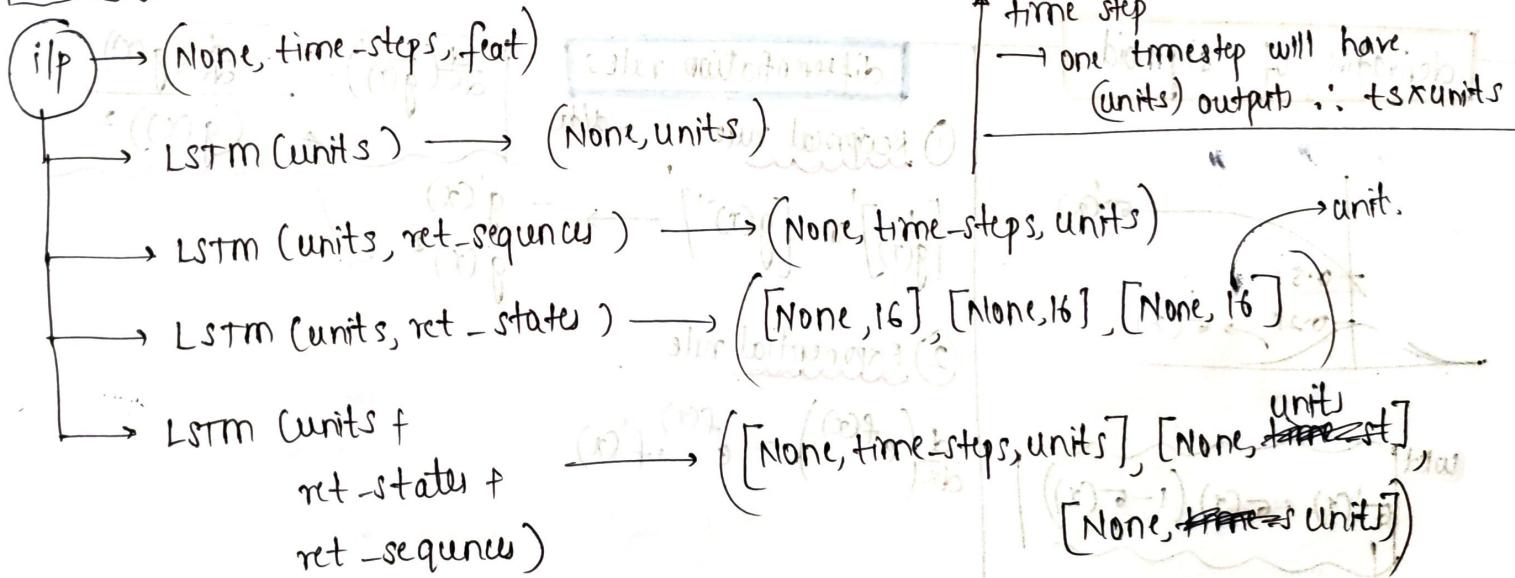
\* LSTM (1, ret-state) = true, but following code is wrong  
when ret-seq = true) → [lstm1, state-h, state-c] → [ [0, 0, 0], [ ] ]

if ip = (3, 1) → o/p of all layers/time steps → [ [ ] ]

o/p at last step → [ ]

o/p cell state of last cell → [ ]

Standard



\* Only ret-states

lstm-1 and state-h are same  
(None, units) (None, units)

ret-state with ret-sequence

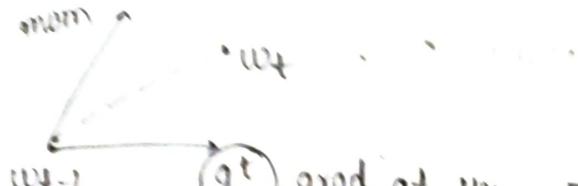
LSTM-1 = (None, ts, units)

state-h = (None, units)

## Nesterov Accelerated gradient

momentum = cma of previous steps

$$\text{step} + \text{momentum} = w_t - w_{t-1} - [\gamma v_{t-1} + \eta g_t] \rightarrow \text{gradient}$$

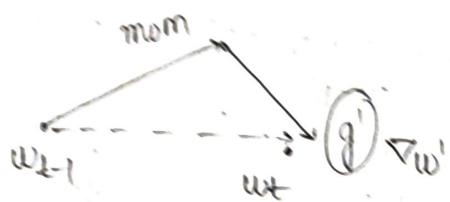


$$(g^t) \text{ grad at } w_{t-1} = \nabla w_t$$

→ it says, ① calculate momentum

② calculate gradient at momentum pt ( $w_{t-1} - \gamma v_{t-1}$ )

③ move in direction of grad

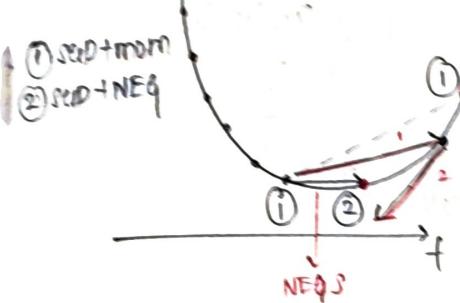


$$(w_t)' = w_t - \gamma \cdot v_{t-1} \rightarrow \text{step1: momentum}$$

$$v_t = \gamma v_{t-1} + \eta \cdot g'$$

where,  
\*\*\*\*  $g' = \nabla w_t'$  (grad at  $w_{t-1}$  momentum)

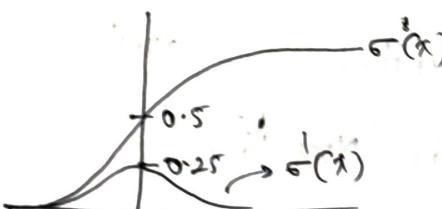
Ex



\* at point ①, we previously had many +ve update  
 $\therefore$  velocity/momentum will be high and it may overshoot minima

\* NEG: - it moves first to momentum point which is on other side, and then gradient is calculated which is  $\text{fve}$  as slope is -ve  $\therefore$  it overshoots less than SEGD + momentum

## derivative of sigmoid



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

max = 0.25 at  $x=0$   
min = 0

## differentiation rules

① Reciprocal rule

$$\left[ \frac{1}{g(x)} \right]' \text{ or } \left[ g(x)^{-1} \right]' \rightarrow -\frac{g'(x)}{g(x)^2}$$

$$\frac{d}{dx} \left( \frac{1}{g(x)} \right) = -\frac{d}{dx} \left( \frac{1}{g(x)} \right) / (g(x))^2$$

② Exponential rule

$$\frac{d}{dx} \left( e^{f(x)} \right) = e^{f(x)} \cdot f'(x)$$

$$\rightarrow \sigma'(x) = \frac{d}{dx} \left[ \frac{1}{1+e^{-x}} \right] = \boxed{\frac{-\frac{d}{dx}(1+e^{-x})}{(1+e^{-x})^2}} \rightarrow ① \text{ using reciprocal rule}$$

→ solving Numerator of ①

$$-\frac{d}{dx}(1+e^{-x}) = -\left[ \frac{d}{dx}(1) + e^{-x} \cdot \frac{d}{dx}(-x) \right] = -\left[ 0 + (e^{-x})(-1) \right] \\ = e^{-x}$$

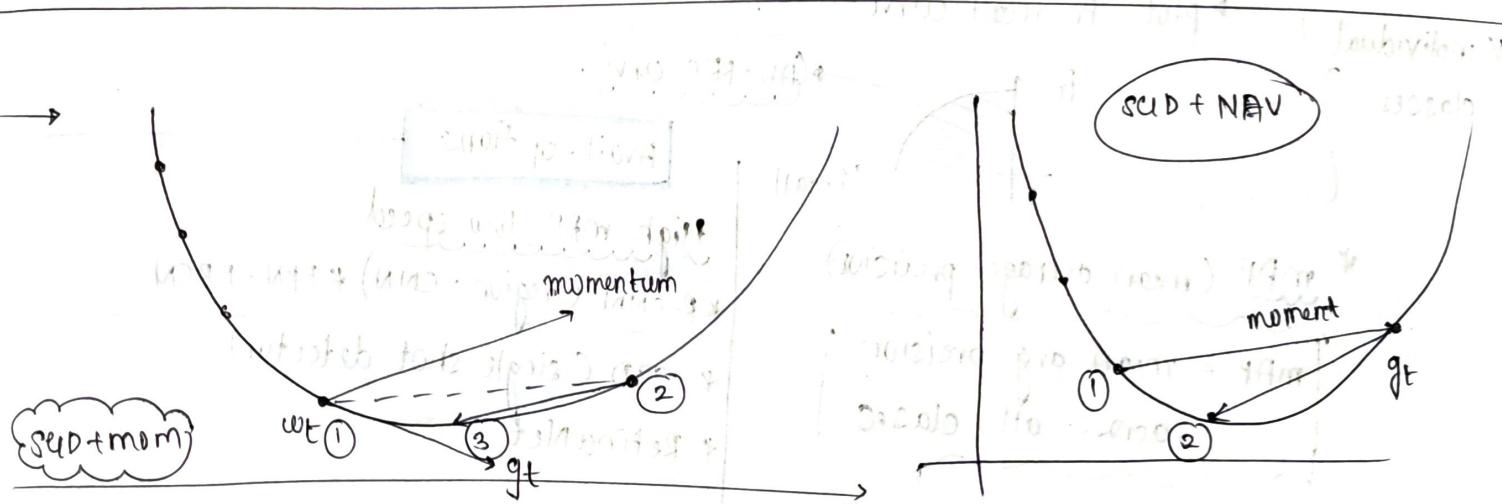
$$\therefore ① \text{ becomes } \Rightarrow \boxed{\frac{e^{-x}}{(1+e^{-x})^2}} \Rightarrow ② \quad \therefore \sigma'(x) = \frac{e^{-x}}{(1+e^{-x})^2}$$

Simplifying ② further

$$\frac{e^{-x} \cdot 1}{(1+e^{-x})(1+e^{-x})} \rightarrow \frac{(e^{-x} + 1 - 1) \cdot 1}{(1+e^{-x})(1+e^{-x})} \rightarrow \boxed{\left[ \frac{1}{(1+e^{-x})} \right] \cdot \left[ \frac{e^{-x} + 1 - 1}{1+e^{-x}} \right]}$$

$$\rightarrow \sigma(x) \left[ \frac{e^{-x} + 1 - 1}{(1+e^{-x})} \right] \rightarrow \sigma \cdot \left[ \frac{e^{-x} + 1}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right] \quad \left\{ \begin{array}{l} \text{split into} \\ \text{fractions} \end{array} \right\}$$

$$= \sigma(x) \left[ 1 - \sigma(x) \right] \quad \boxed{\therefore \sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))}$$



\* at pt wt, gt/momentum ↑ as we made lot of tve increments previously

∴ due to this point overshoots to right side

(in case of NAEQ, we calculate grad at momentum)

right side ∴ it will be closer to minima and

which is won in negative direction

# YOLO Algorithm

object detection

## Object Detection

\* focus on bounding boxes

and not on pixels → use image segmentation (slower and not suitable for real-time apps)

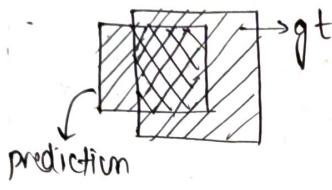
\* IP → algo → o/p image

with bounding boxes with confidence level (0-1)

## Performance Metrics

$$IoU = \frac{\text{overlap area}}{\text{union area}}$$

gt (ground truth) - human labelled data



\* prediction is correct if  $IoU \geq 0.5$

steps

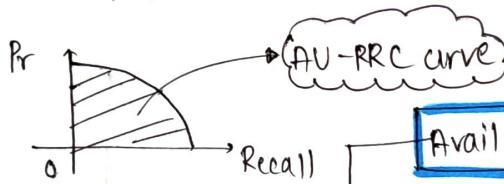
① convert o/p to binary ie True or False using IoU & bounding boxes

② Image may have multiple bounding boxes,

so,

\* for each class (chairs person...)

- calculate avg-precision = area under pre-recall curve
- plot Pr-recall curve



\* mAP (mean average precision)

$$mAP = \text{mean avg precision across all classes}$$

## Avail-options

\* high mAP, low speed

\* R-CNN (region-CNN) \* FPN - FRCN

\* SSD (single shot detector)

\* RetinaNet

YOLO → faster and good mAP

- YOLOv3 - 320 → input image size
- YOLOv3 - 416
- YOLOv3 - 608

Note: as image size ↑ time ↑

## Trade off

→ speed vs mAP

video based  
(self-dr cars)  
real time face

medical diagnosis  
or  
OCR

# YOLO v3

## Architecture

only used for initialize

## i) Feature Extractor / Backbone (Pre-Training)

### DarkNet - 53

53 conv-layer (fully convolutional w)

YOLOv2  
V3  
DarkNet - 19  
DarkNet - 53

Layer	Filters	Size	Strd
conv	32	(3x3)	256x256
conv	64	(3x3)	128x128
		stride=2	
IX	conv 32	(1x1)	
	conv 64	(3x3)	
	residual		128x128
1X	conv 128	3x3/2	64x64
2X	conv 64	(1x1)	
	conv 128	(3x3)	
	residual		64x64
2X	conv 256	3x3/2	32x32
	conv 128	1x1	
	conv 256	3x3	
	residual		32x32
8X	conv 512	3x3/2	16x16
	conv 256	1x1	
	conv 512	3x3	
	residual		16x16
8X	conv 1024	3x3/2	8x8
	conv 512	1x1	
	conv 1024	3x3	
	residual		8x8
4X	Avgpool	global	
	converted	100	
	softmax		

416x416x3 → 13x13x1024  
(416/32) last filter size

## Characteristics

- ① 53 conv layers
- ② fully conv neural net
- ③ Each conv = conv → Batch Norm → Leaky ReLu
- ④ conv with stride=2 → acts as max-pool
- ⑤ no maxpooling

whole model is trained on ImageNet and save weights

→ this only detects image objects but not bounding boxes

why we DarkNet when we had ResNet?

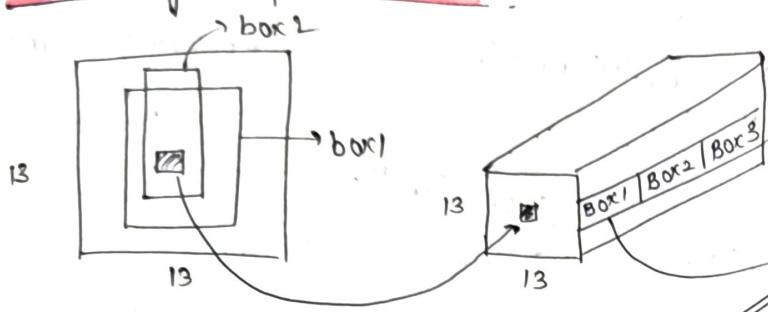
Backbone	Top-5	B Ops	BFCOps/s	FPS
Res-152	77	2930	1090	57
Dark Net	77.2	18.7	1457	78

- darknet almost gives same or little less accuracy comp to resNet
- \* BN Ops is less (Billion of ops needed)
- \* BFCOps/s is high (Bill of float op/sec).
- \* and FPS is high

## ② Bounding boxes and output

$$416 \times 416 \times 3 \xrightarrow{\text{DNet}} 13 \times 13 \times 1024 \xrightarrow{\text{Box}} 13 \times 13 \times 425$$

How do you represent boxes?



O/p

(Each cell may belong to many boxes)  $\#$  of boxes

Each bounding box

$x$	$y$	$tw$	$th$	$P_o$	$C_0$	$C_1$	$C_2$	$C_3$	$C_4$	$C_5$	$C_6$	$C_7$	$C_8$	$C_9$
-----	-----	------	------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

$(x,y)$  → centre of box with dimensions  $tw \times th$

$P_o$  → objective score → 0/1  
whether it contains an obj/not

$C_0 \rightarrow C_K$ : class scores for know

$$\# \text{box} = P_o * C = [P_o, C_1, P_o, C_2, P_o, C_3, \dots, P_o, C_{80}]$$

$$\# \text{score} = P_o * C = \begin{bmatrix} P_o, C_1 \\ P_o, C_2 \\ P_o, C_3 \\ \vdots \\ P_o, C_{80} \end{bmatrix} = \begin{bmatrix} 0.12 \\ 0.21 \\ 0.80 \\ \vdots \\ 0.10 \end{bmatrix} \rightarrow \text{probability that object bounded } \in \text{class 2}$$

$$P(C=c) = P_o * P_K \text{ or}$$

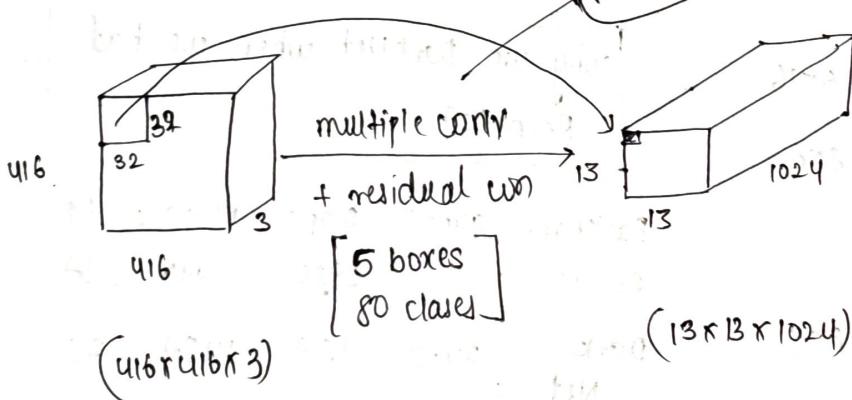
$$P_o * C_K$$

assume we get  $P_o, C_3 = 0.80$  as max

∴ the box bounds object

belonging to class = 3

not darknet



$$\frac{416}{13} = 32$$

∴ Each pixel in O/p is of size ( $13 \times 1024$ ) is captured from ( $32 \times 32 \times 3$ ) grid input.

Effective Receptive Field

height, width

depth, channels

pixels

pixels

pixels

pixels

pixels

pixels

③ **Anchor Boxes** → prior boxes where you highly probable to find

\* Tried to predict  $(tx, ty, tw, th)$  directly but couldn't get good results

so,

① Fix some predefined boxes called anchor-boxes :  $[Cx, Cy, Pw, Ph]$

② Define bounding boxes in terms of anchor boxes →

: given box coordinates and anchor boxes you can find  $[bx, by, bw, bn]$

Learned box + anchor box → actual bounding box

**Objectness score**  $P_o$

Is there any object or not? 0/1

→ have a log-regression for each box

\*  $[C_0 \dots C_K]$

→ tried softmax — didn't work because of hierarchy presence

⇒  $\text{softmax}$

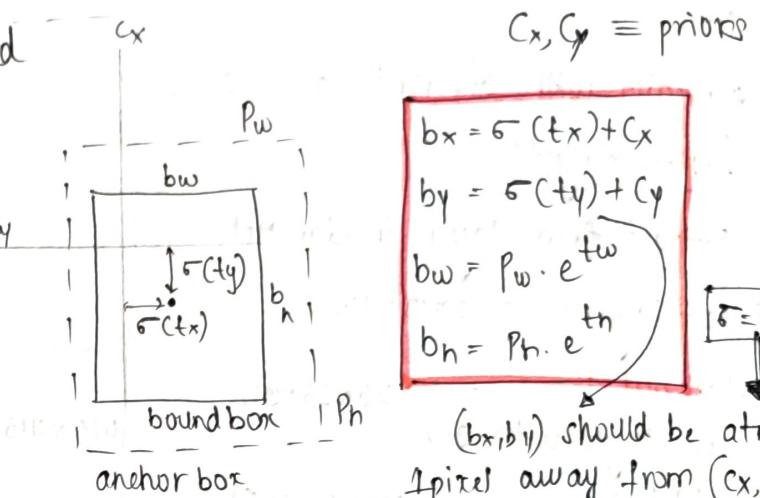
Build one logistic regression for each class i.e. u need  $P_o + [C_0 \dots C_K]$

log-regressions

**Loss**

$$\begin{aligned} & \text{log-loss for } P_o \\ & + \text{log-loss for } (C_0 \dots C_K) \\ & + \lambda_{\text{co-ord}} * \text{sq-loss for } tx, ty, \sqrt{tw}, \sqrt{th} \end{aligned}$$

all bounding boxes containing object in  $y_{\text{train}}$



How to find how many anchor boxes and location?

→ apply k-means clustering on COCO dataset which is set of all bounding boxes. → K=5 was found to be best which indicates that there are 5 anchor boxes have high probability region

if for any Bbox if  $P_o$  is very small  
→ it has very less chance of enclosing any object  
∴ can be ignored

$$\begin{aligned} \lambda_{\text{co-ord}} = 5 & \rightarrow tx, ty \rightarrow 0 \dots \infty \\ \lambda_{\text{noobj}} = 0.5 & \rightarrow P_o = (0 \dots 1) \end{aligned}$$

$$\begin{aligned} & \sum_{\text{all boxes in } y_{\text{train}}} \lambda_{\text{no-obj}} * \text{log-loss for } (P_o) \\ & + \sum_{\text{predicted which don't have obj}} \end{aligned}$$

(predicted boxes not present in  $y_{\text{train}}$ )

## 4) multiscale prediction

\*  $13 \times 13$  image feature can only detect large objects

### Idea

\* predict with different grid sizes

↳ How to get?

Extract from layers in darkNet

① last layer  $\rightarrow 13 \times 13 \times 1024$

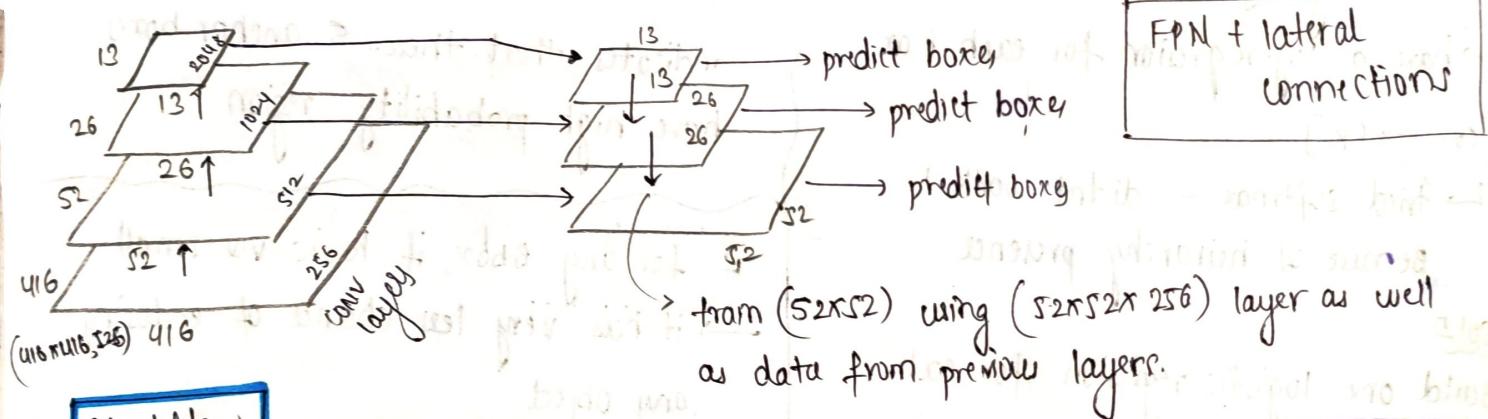
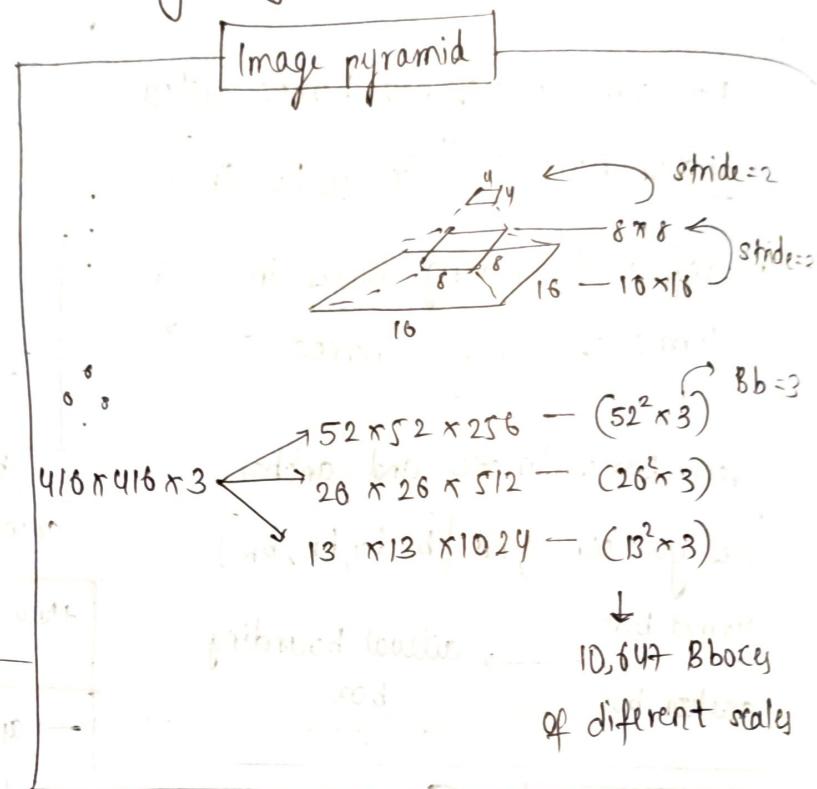
② last but-1  
conv layer with  $s=2$   $\rightarrow 26 \times 26 \times 512$

③ 2<sup>nd</sup> last  
conv layer  $s=2 \rightarrow 52 \times 52 \times 256$

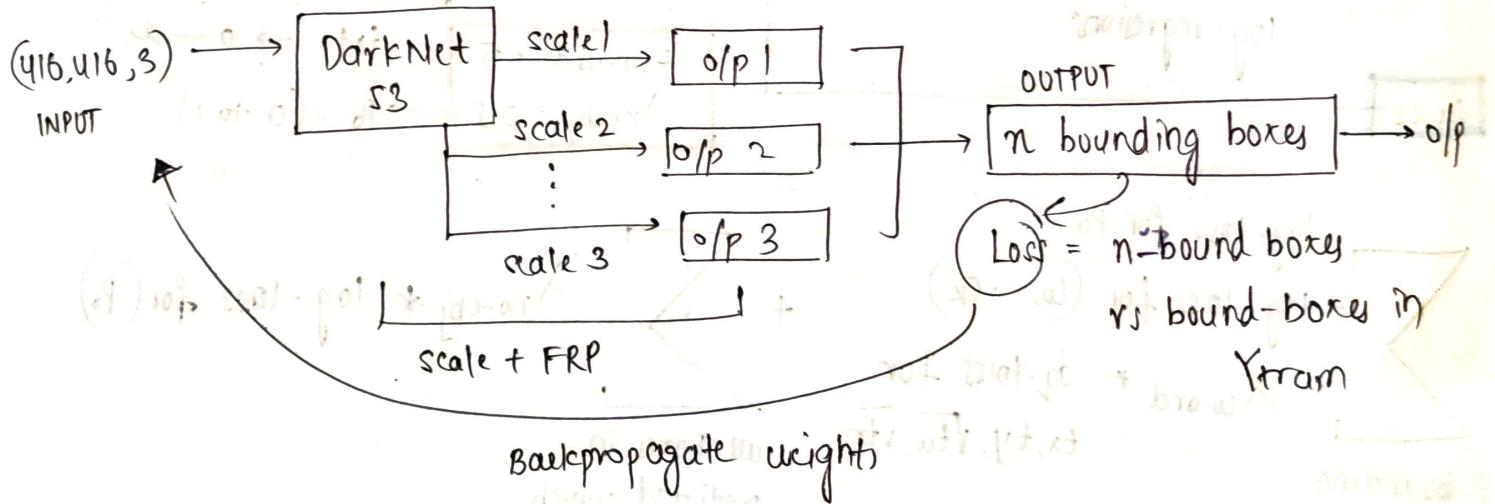
### Feature pyramid networks

#### FPN

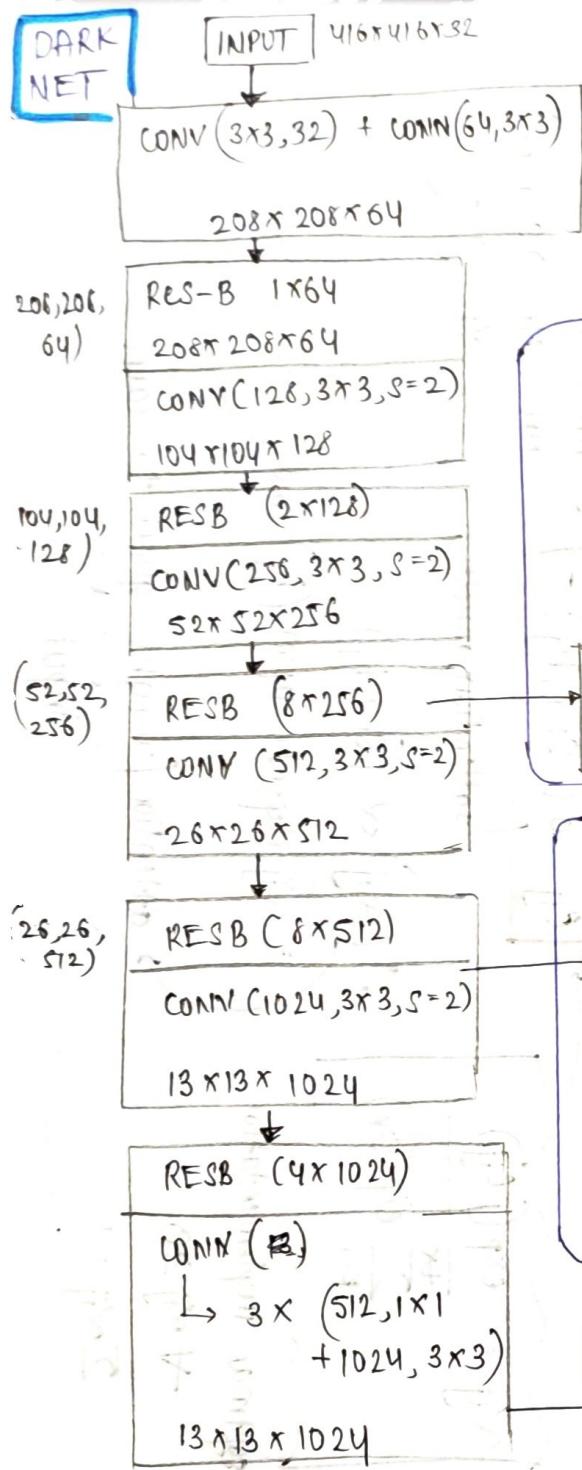
→ not only use different grid sizes  
to predict but also cross-grid information to use



### Workflow

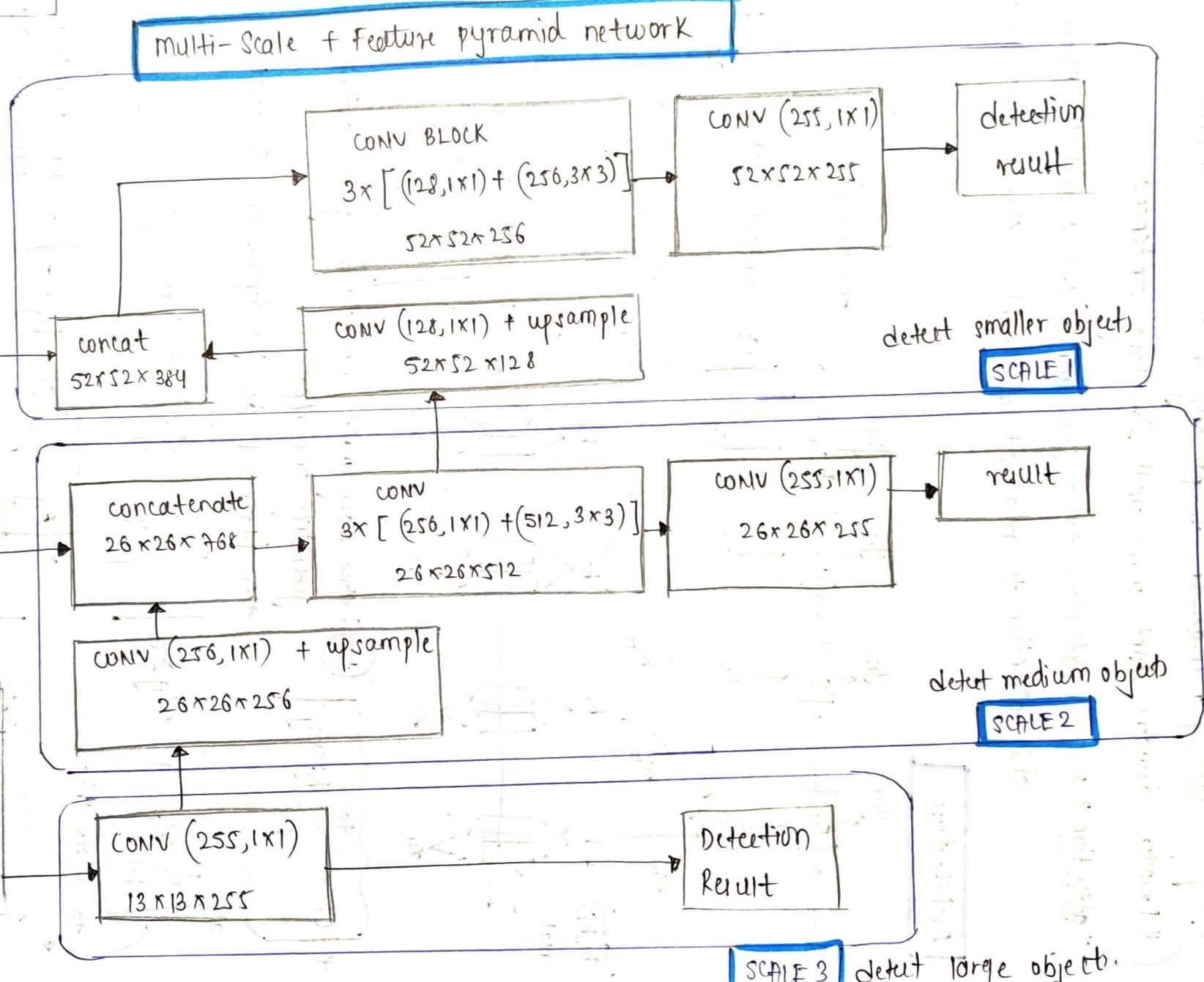


## DARK NET



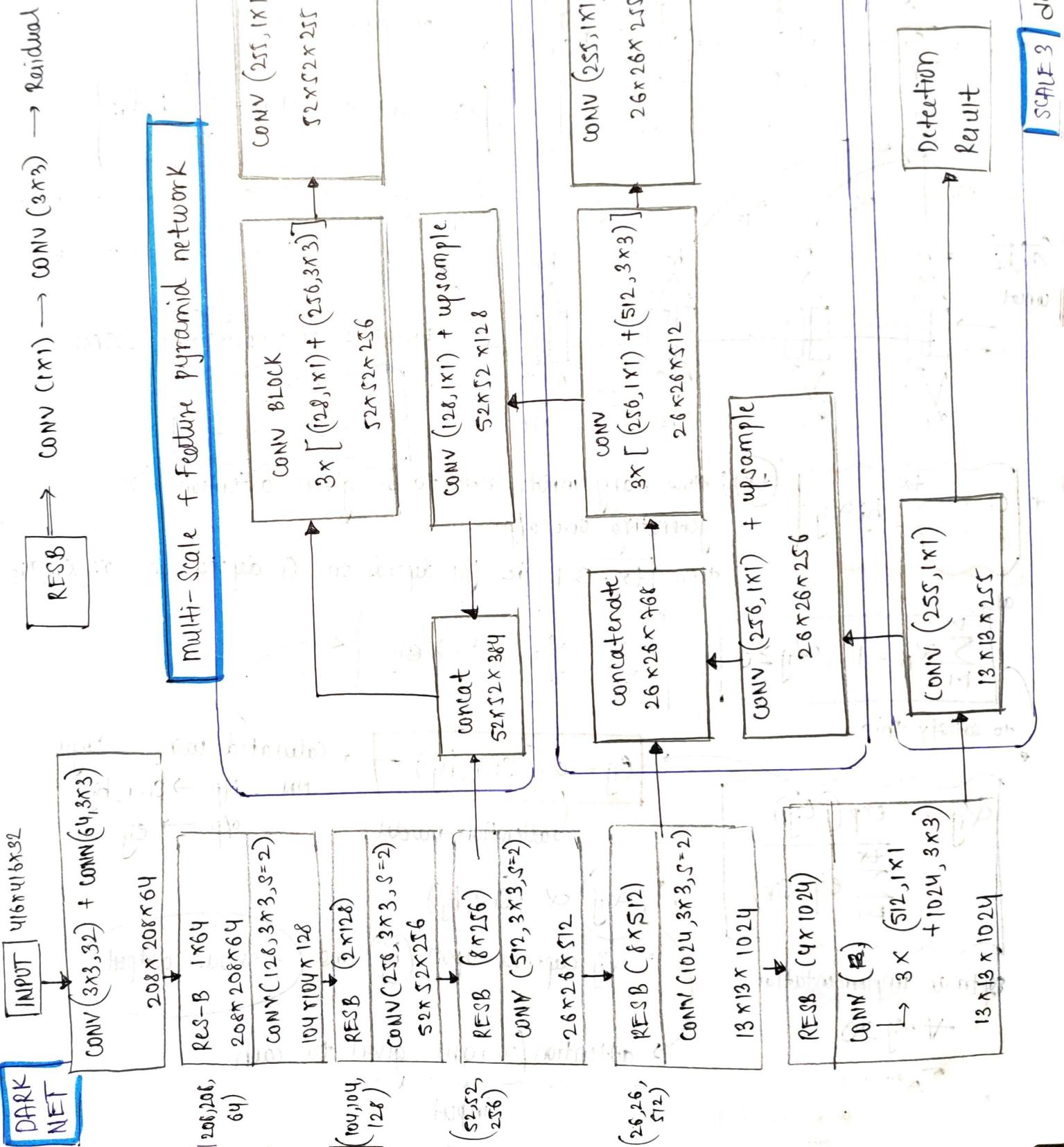
**RESB**  $\rightarrow$  CONV (1x1)  $\rightarrow$  CONV (3x3)  $\rightarrow$  Residual

## YOLO V3 ARCHITECTURE



# YOLO V3

## ARCHITECTURE



## Attention Models

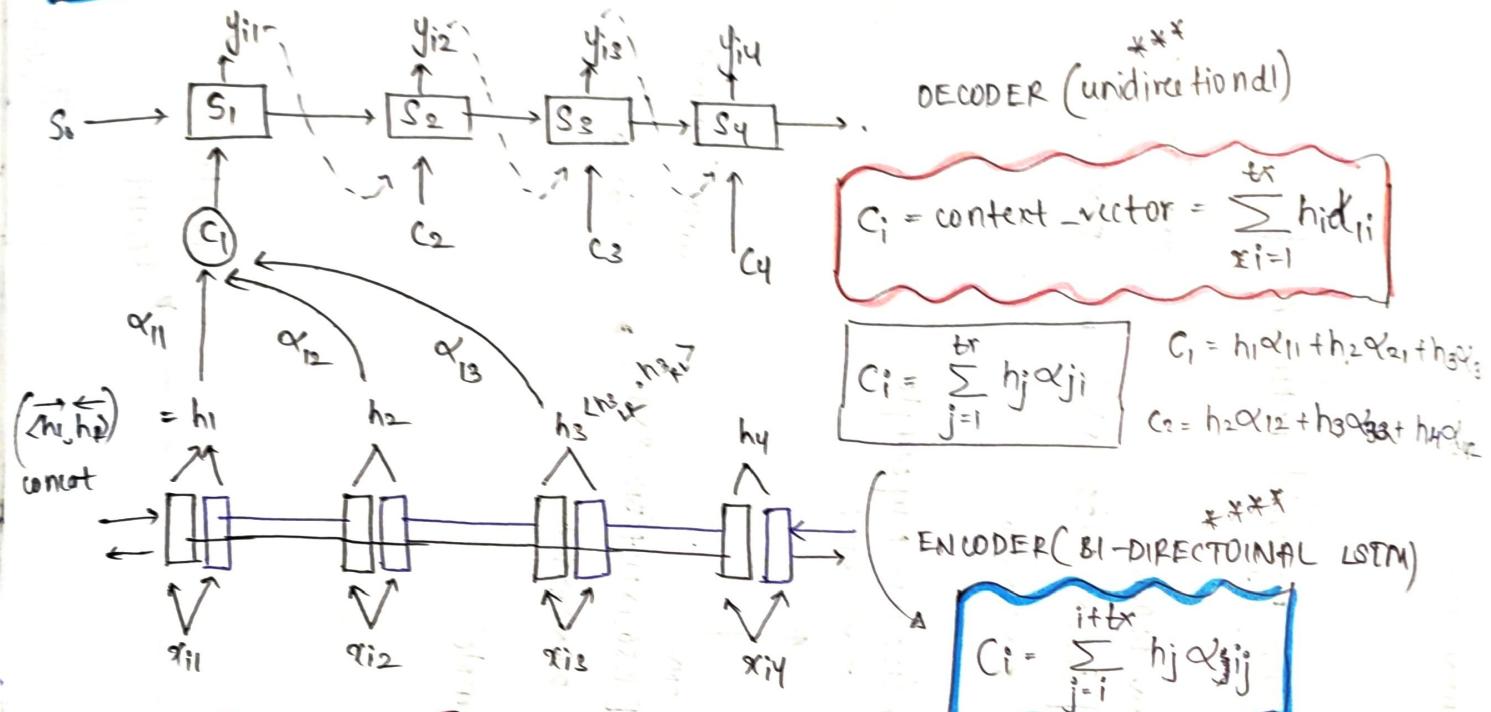
this may not capture the essence of whole sentence

why you need? In enc-decoders we pass only one o/p ( $\vec{w}$ ) from encoder

\* Encoders decoders don't work well for long sentences

\* attention? We give attention to a few words ( $t_x$ ) translate them and shift attention to next set of words (Human characteristic)

### Architecture



$$c_i = \sum_{j=1}^{t_x} h_j \alpha_{ij}$$

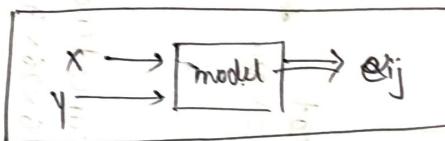
$t_x$ : How many inputs need to be given attention for generating one o/p.

Here  $t_x = 3$  ie  $y_{i1}$  depends on  $c_1$  depends on  $x_1, x_2, x_3$

also,

$$\sum_{i=1}^{t_x} \alpha_{ri} = 1, \alpha_{ij} \geq 0$$

to satisfy this



$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{t_x} \exp(e_{ik})}$$

softmax implementation

$$\alpha_{ij} \geq 0$$

$$\sum \alpha_{ij} = 1$$

$$\sum \alpha_{ik} = 1 \quad \text{all } \alpha \text{ correctly}$$

$$e_{ij} = a(s_{i-1}, h_j) \rightarrow \begin{array}{l} \text{calculated using 2-layer NN} \\ \text{i/p} \rightarrow s_{i-1}, h_j \\ \text{o/p} \rightarrow e_{ij} \end{array}$$

$$\alpha_{ij} \propto (s_{i-1}, h_j)$$

$\alpha_{ij}$  depends upon  $(h_j)$  and previous output so

attention/weight given to each input

# Time complexity

$K_1 = \text{len(input)}$

$K_2 = \text{len(output)}$

$$TC = O(K_1, K_2)$$

if  $T_x = K_2$ ,

$$\text{then } TC = O(K_1, K_2)$$

i.e consider whole sentence as  
attention to generate o/p.

## $\alpha_{ij}$ importance

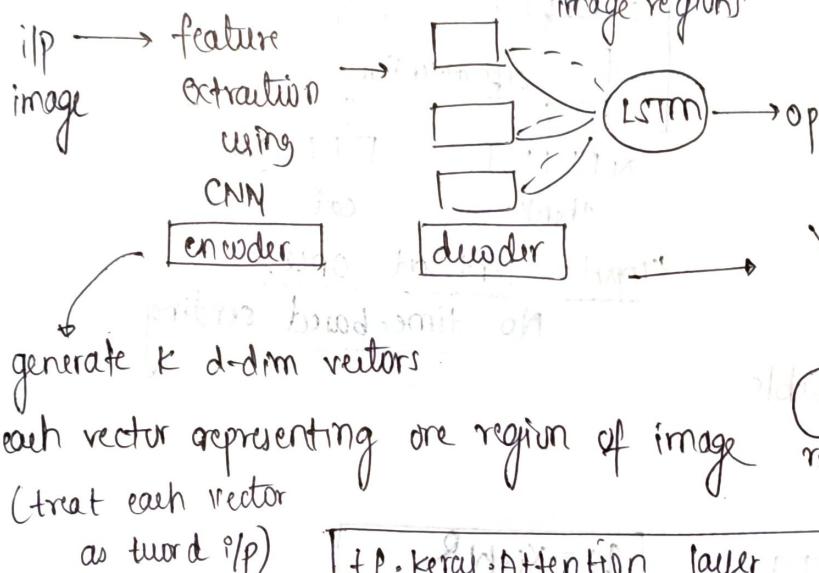
Ex: for calculating  $\alpha_{ij}$  you use  $\alpha_1, \alpha_2, \alpha_3$ ,  
in this case if  $\alpha_3$  is large than  
that particular word had more impact  
on prediction

Ex: french  $\rightarrow$  english

	you	and	call	o/p
ti	0.9	0	1	
ans	0.1	2	0	
mé	0	2	10	

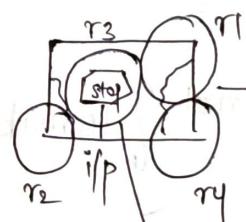
to predict  $w_3 = \text{call}$   
 $mé$  was most useful.

## Image Captioning



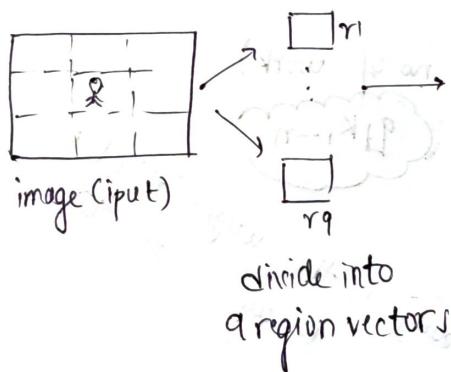
tf.keras.Attention layer  
→ inbuilt into keras  
other blogs

visualize  $\alpha_{ij}$



A mountain road with  
stop sign  
see highest  $\alpha_{ij}$   
say highest is  $\alpha_{31}$   
∴ region ③ contributed  
most

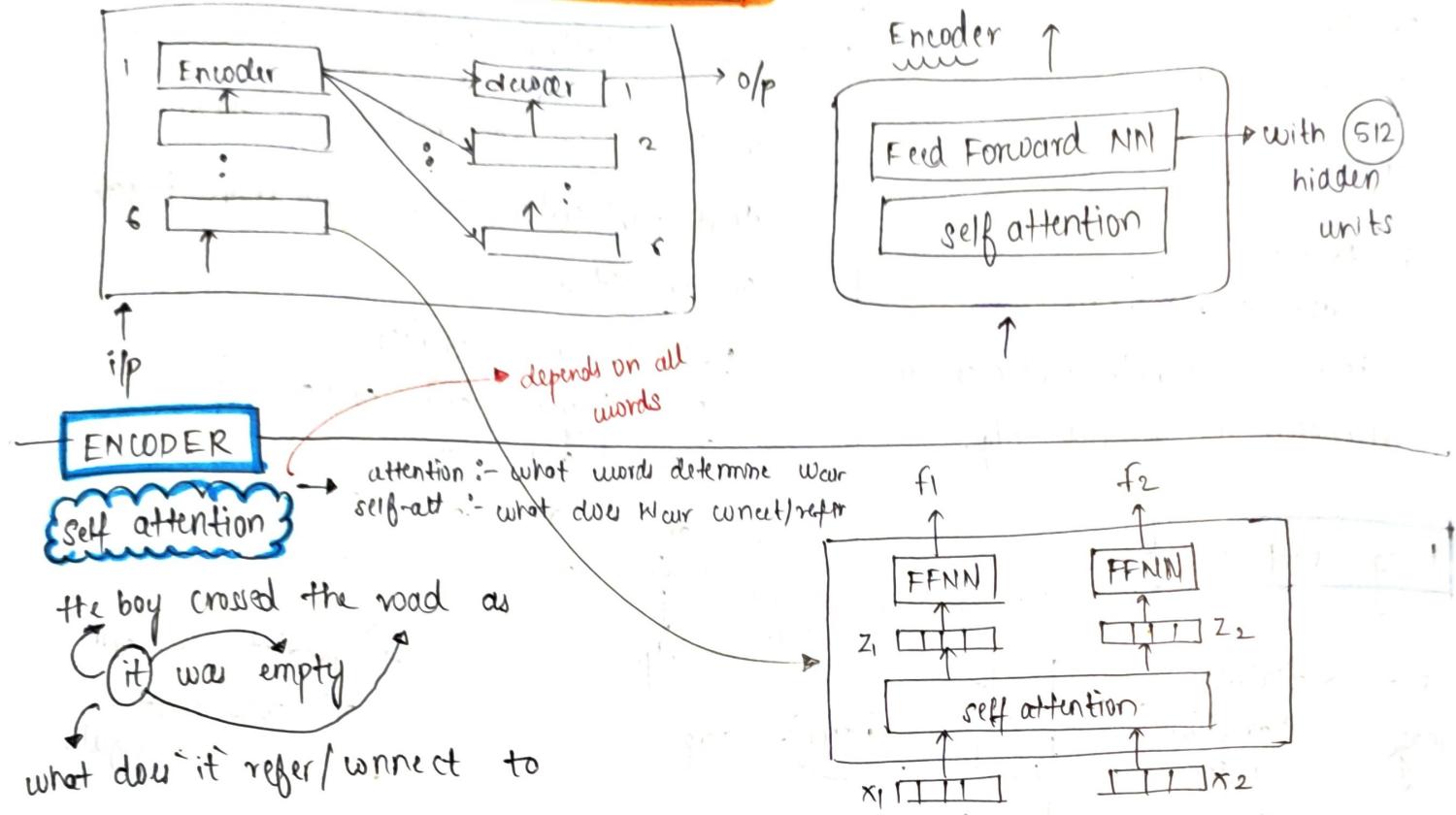
## Visualize



visualize  $\alpha_{ij}$  as feature importance  
(region with face will have high  $\alpha_{ij}$ )

# Transformers

→ Not RNN (!time based)

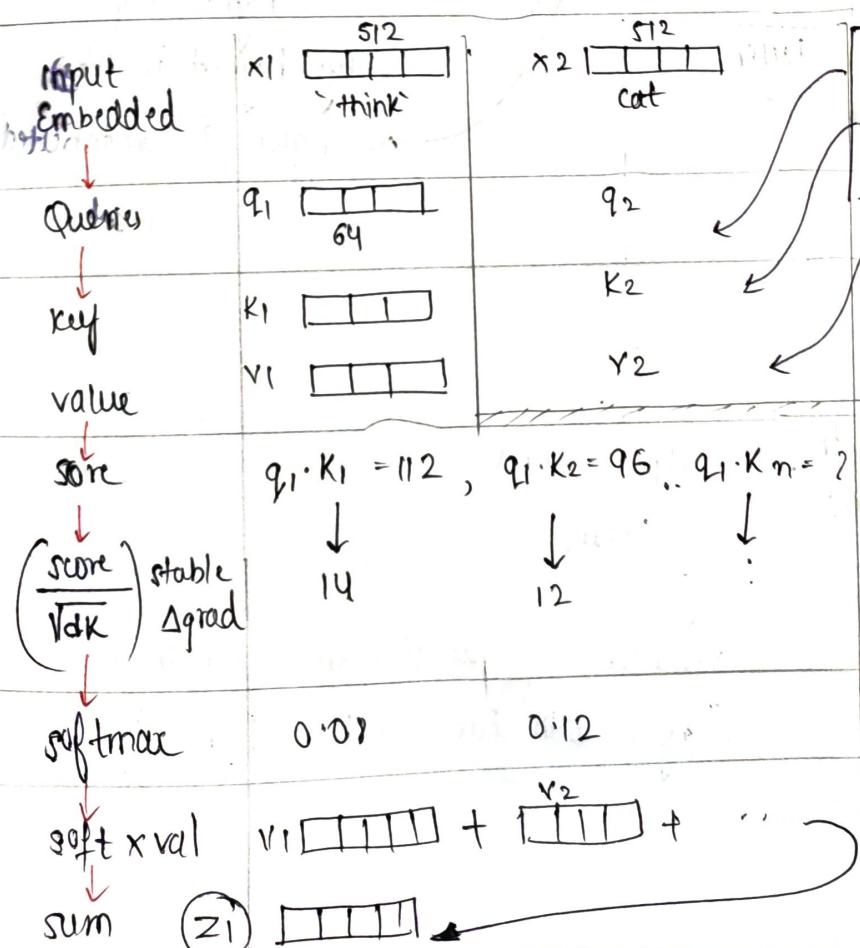


## steps

hyper-params :

$d_k = 64$ ,  $W^Q$  - Query mat } learnable  
 $W^K$  - keys  
 $W^V$  - value }

\*inputs sent at once  
 No time-based sending



$$\begin{aligned} q_i &= x_i \cdot W^Q \\ k_i &= x_i \cdot W^K \\ v_i &= x_i \cdot W^V \end{aligned}$$

$$\text{score} = \frac{q_i \cdot k_i}{\sqrt{d_k}}$$

( $n$  - no of words)  
 $= q_i \cdot k_i - n$

here  $q_i$  is the weight/attention given to current word to word  $i$

here  $v_1 > v_2$  as  $v_1$  is of same word

## matrix computation

X

$$\begin{array}{|c|c|c|} \hline & 1 & 2 \\ \hline 1 & & & \\ \hline 2 & & & \\ \hline 3 & & & \\ \hline \end{array} \times W^Q = Q \Rightarrow \text{softmax} \left( \frac{Q \times K^T}{\sqrt{dk}} \right) V$$

$$X \times W^K = K$$

$$X \times W^V = V$$

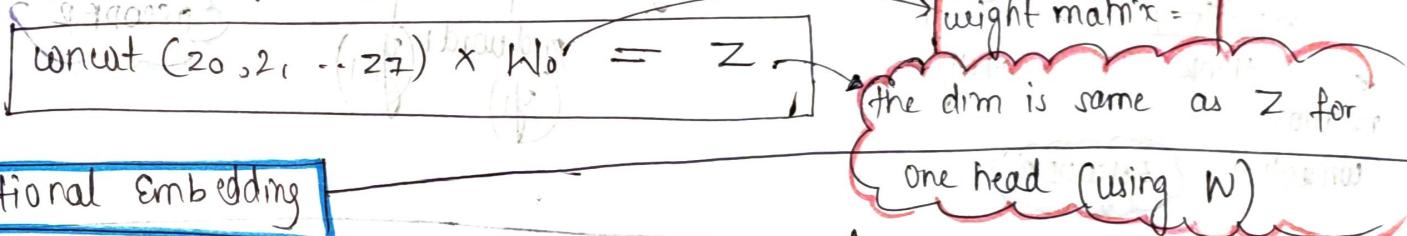
$$= \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} \quad \begin{array}{|c|c|c|} \hline & 1 & 2 \\ \hline 1 & & & \\ \hline 2 & & & \\ \hline 3 & & & \\ \hline \end{array}$$

$\therefore$  self-attention focuses on words which matter more in its neighborhood (i.e. same sentence)

One-headed :- only 1 set of matrix

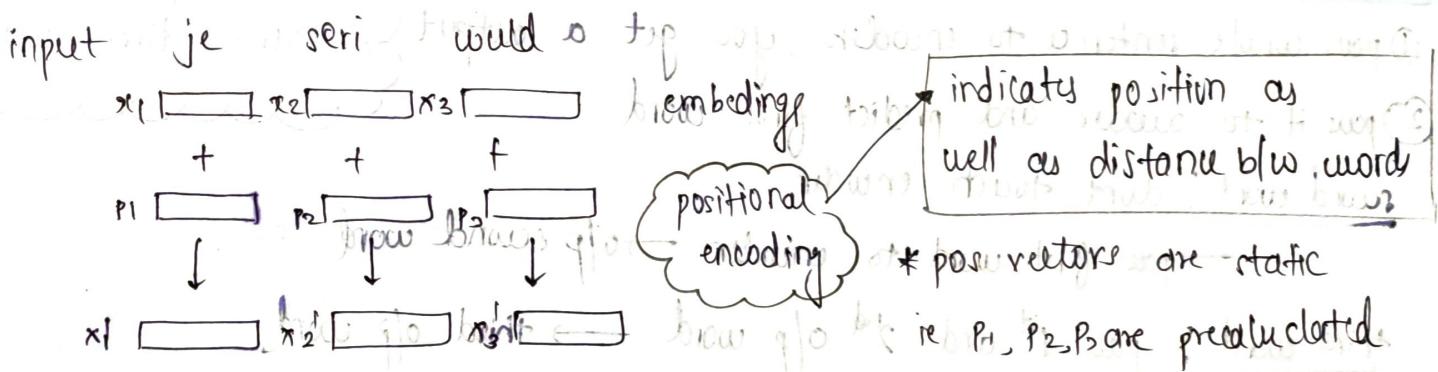
multi-headed ( $8$ -pairs of  $W^Q, W^K, W^V$ )  $\rightarrow$  same steps but with  $(W_Q, W_K, W_V)$  to  $(W_Q^8, W_K^8, W_V^8) = 8$

$\therefore$  for  $X$  we get  $(z_0, z_1, z_2, \dots, z_7)$



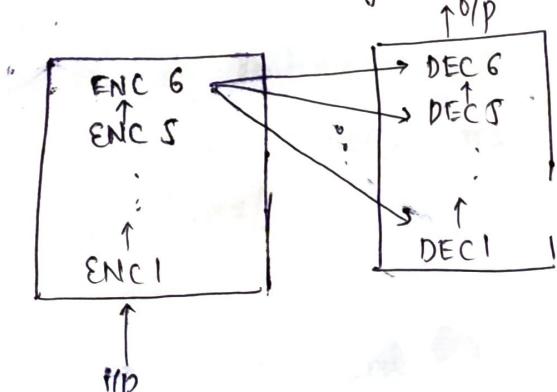
## positional Embedding

we pass all words at once, how to maintain time/position



## original research paper

(Attention is all you need)

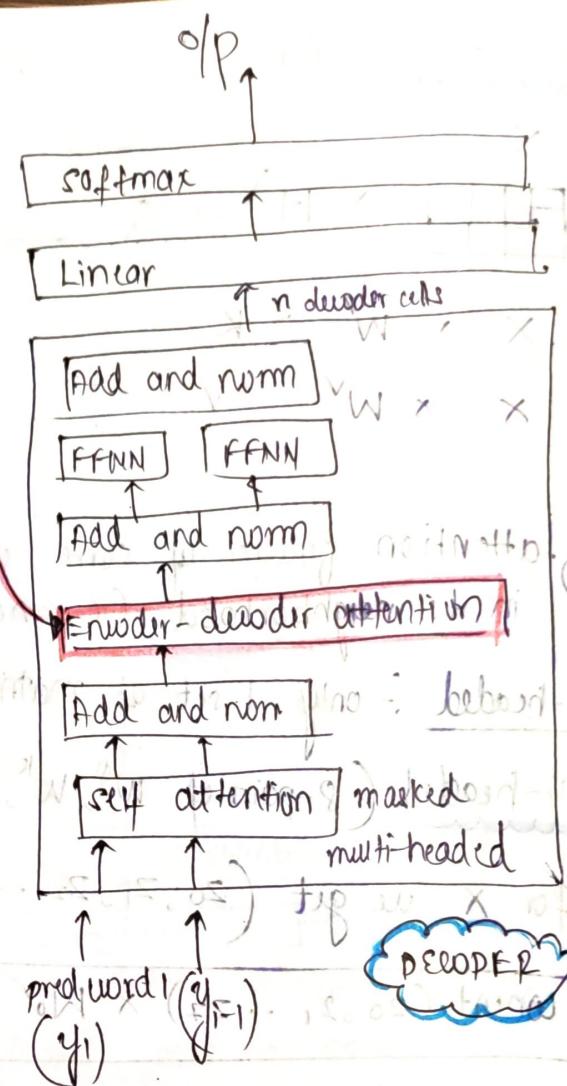
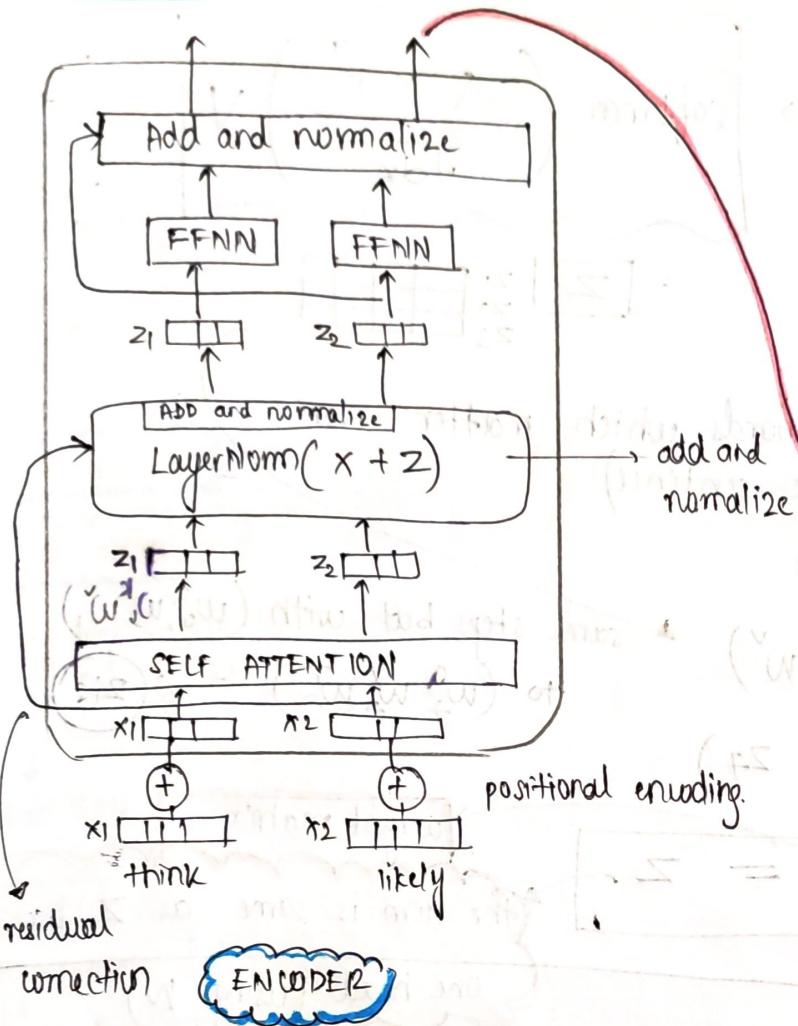


learnable matrices

①  $W_Q^{1-8}$  ②  $W_K^{1-8}$  ③  $W_V^{1-8}$

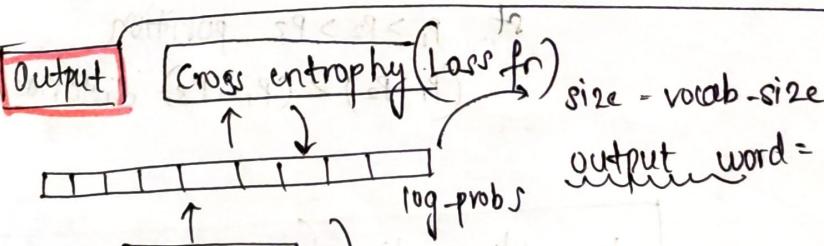
④  $W$  simplified  $\rightarrow$   $W'$

## ENCODER architecture



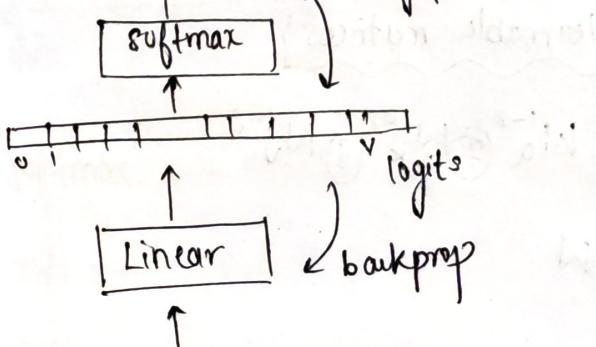
## Train phase

- ① pass whole sentence to encoders, you get a output
- ② pass it to decoder and predict first word
- second word :- dont touch encoder
- pass first word to encoder → o/p second word
- third word :- pass 1<sup>st</sup> and 2<sup>nd</sup> o/p word → third o/p word



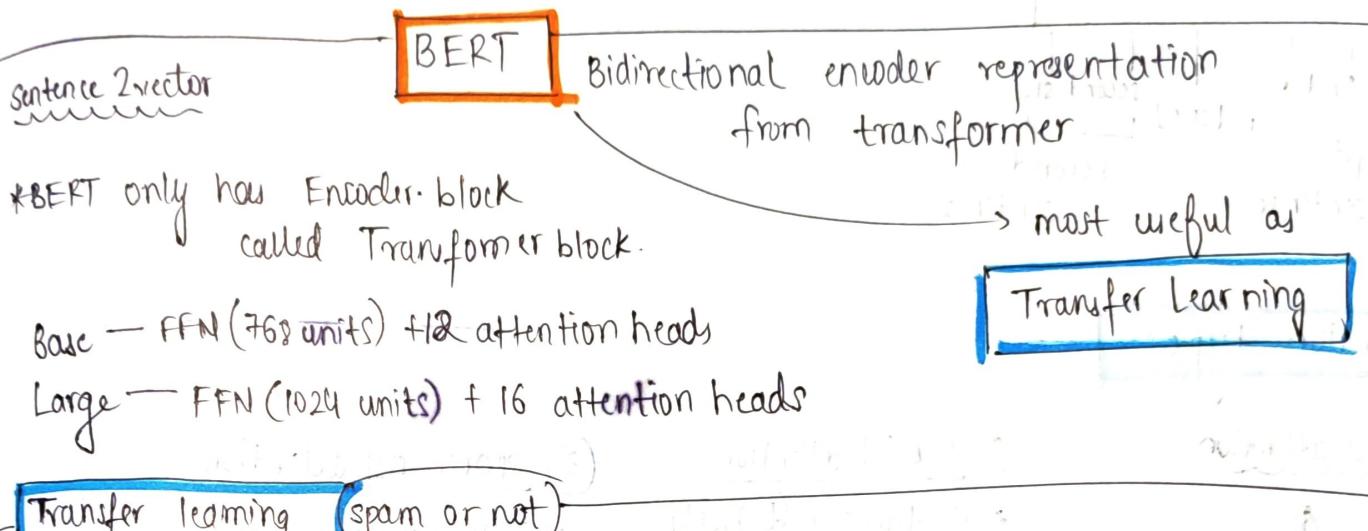
size = vocab-size

output word = vocab [arg max op]

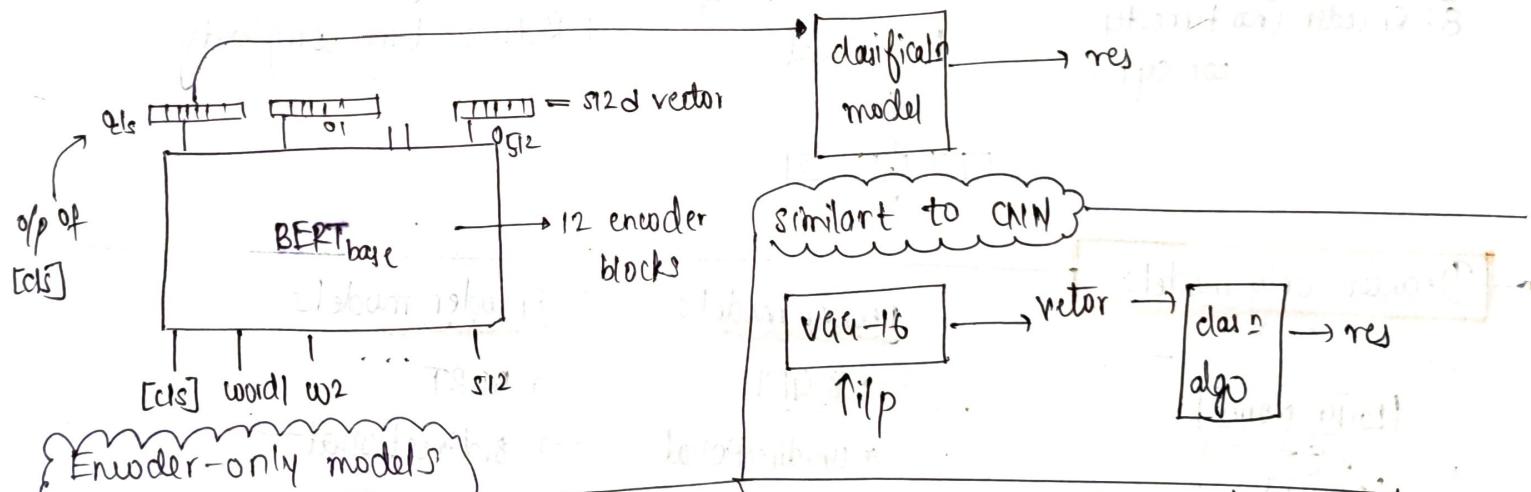


## Advantages of Transformers

- \* Faster than RNN models as all data is ingested once.
- \* state-of-the-art (2018)
- \* Inspired from CNN (self-attention/neighborhood  $\approx$  convolutions in CNN).



Transfer learning (spam or not)  
 Base BERT; can take upto 512 worded sentence and produces 512 o/p each of dimension (768)



**BERT Training**  $\approx$  word-2-vec training

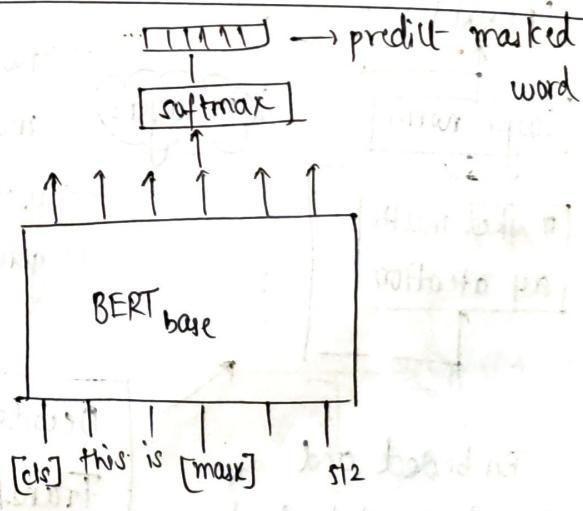
\* Bidirectional models, while training

look @ words before and after

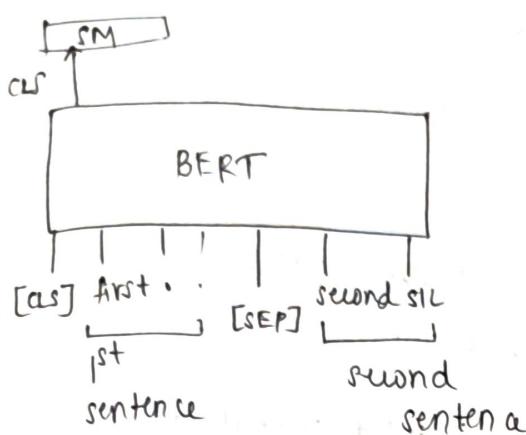
**steps**

\* provide many sentence as ilp and while doing so randomly mask 15% of words/tokens

\* try to predict masked tokens given other words and optimize



BERT with 2 sentences

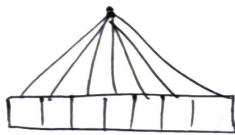


## BERT training

NWP = next word prediction

## Attention types

## ① self attention



t depends on all words

## ② Masked attention

- + at timestep ( $t_i$ ), self attention calculated on  $x_{t_i}$  to  $x_{t-1}$
- + used in decoders as

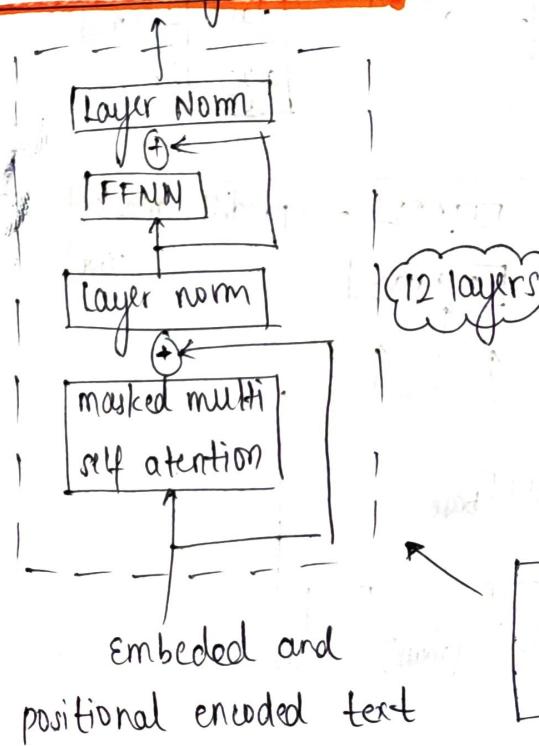
+ used in deoders as  
time board

### ③ sparse self attention

\* if  $\text{len}(\text{sentence}) = \text{NP}$

choose only  $\sqrt{N}$  words for calculating self attention

## Decoder only models



## Deoder models

→ BS Egypt

→ unidirectional

→ normally use  
masked intent

→ used more in generative tasks

## Encoder models

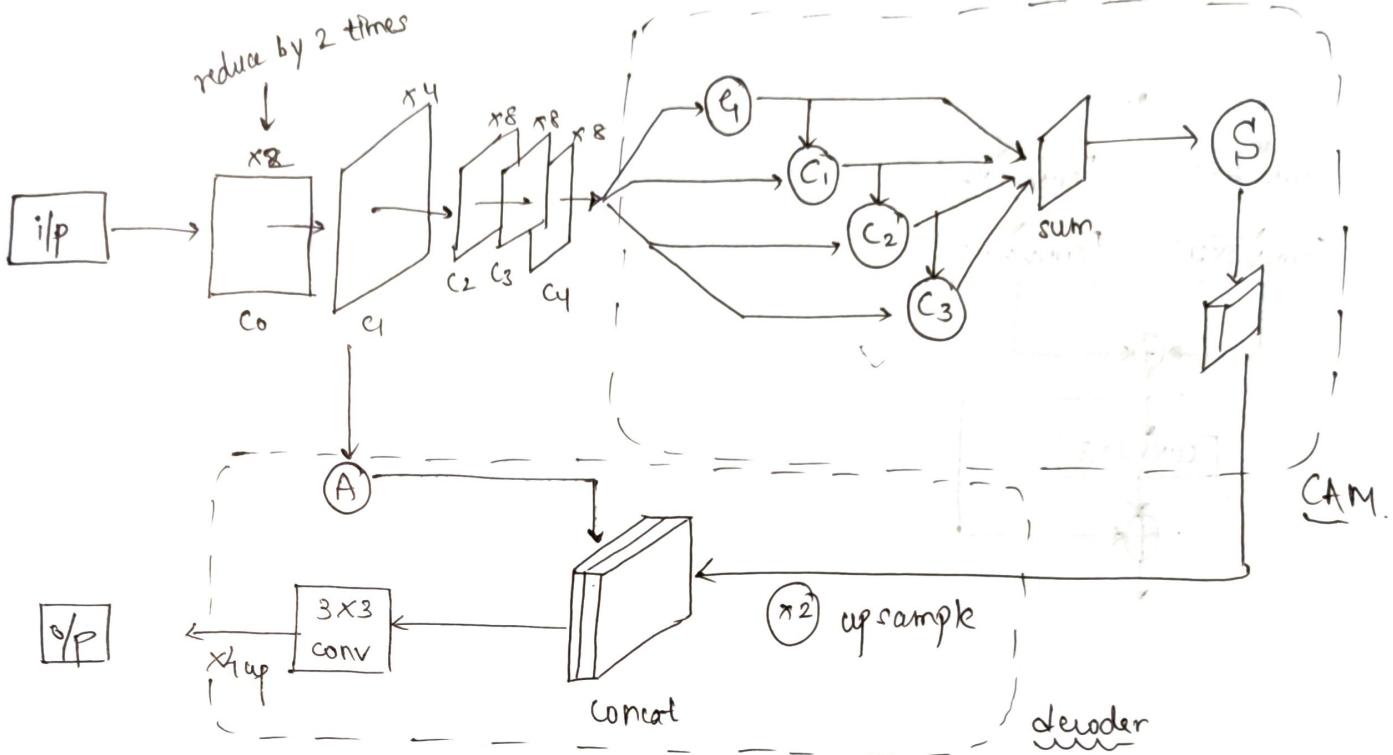
→ BERT

→ Bidirectional

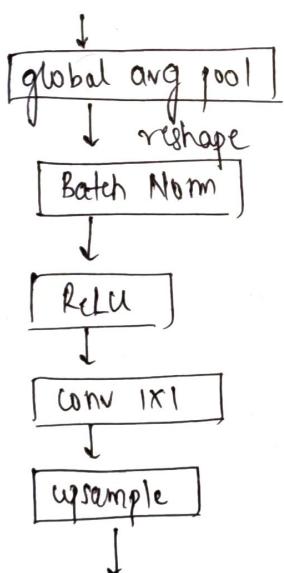
→ simple self attention

→ used mostly for  
tent representation /  
embedding tasks

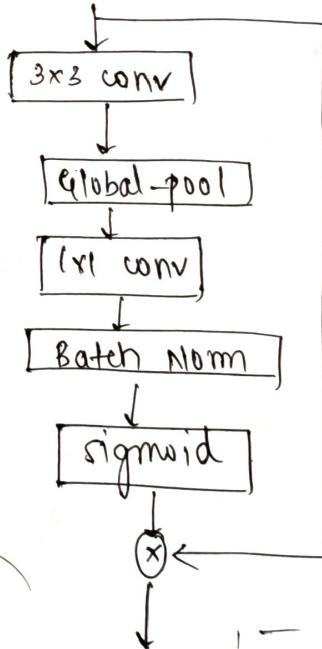
# CANET Image segmentation



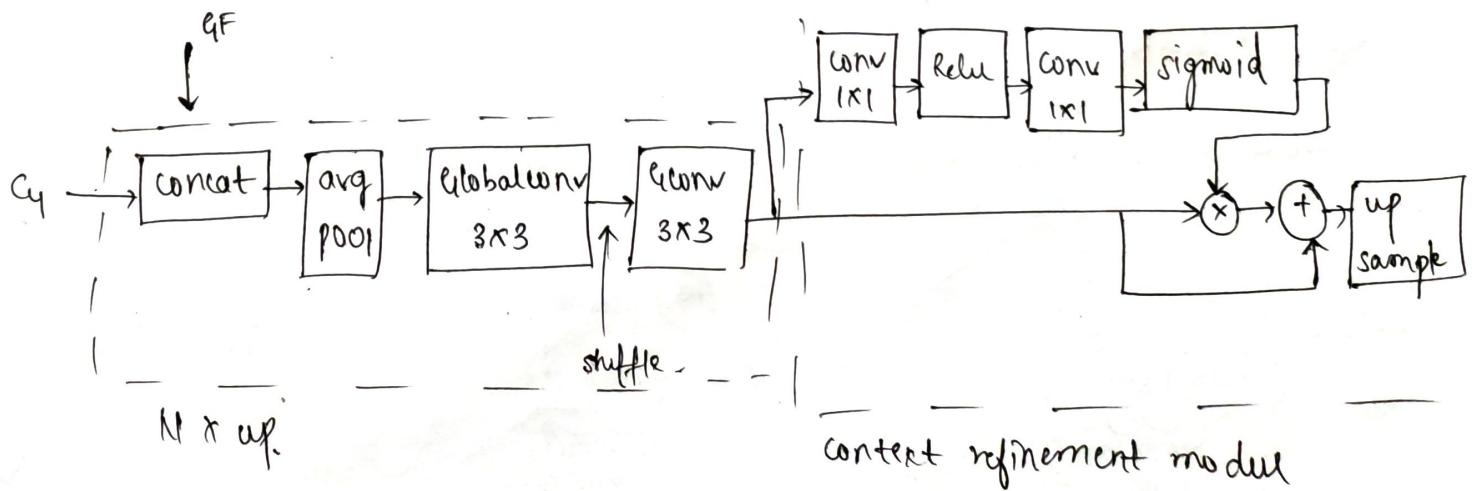
① ④ Global flow



③ ⑤ Feature selection module

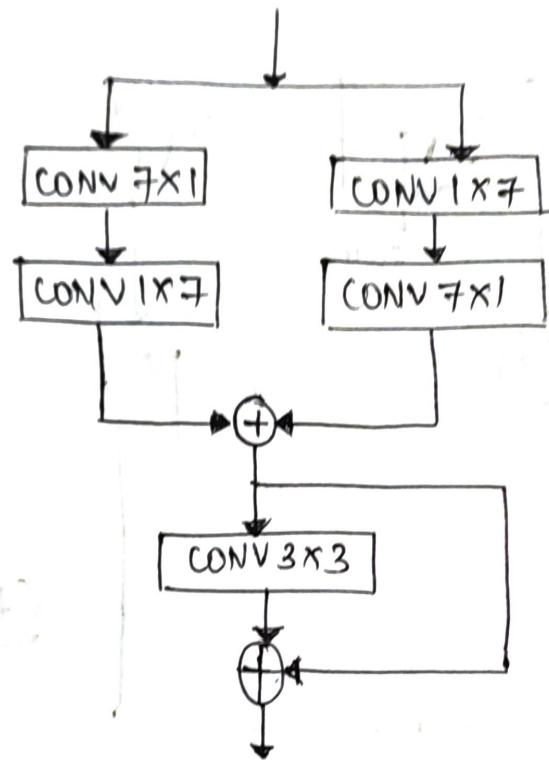


② ⑥ Content fusion module



(58)

AGCN Adapted global convolution n/w



downsample

upsample

target

global context nodes

target

local nodes

local

global

global

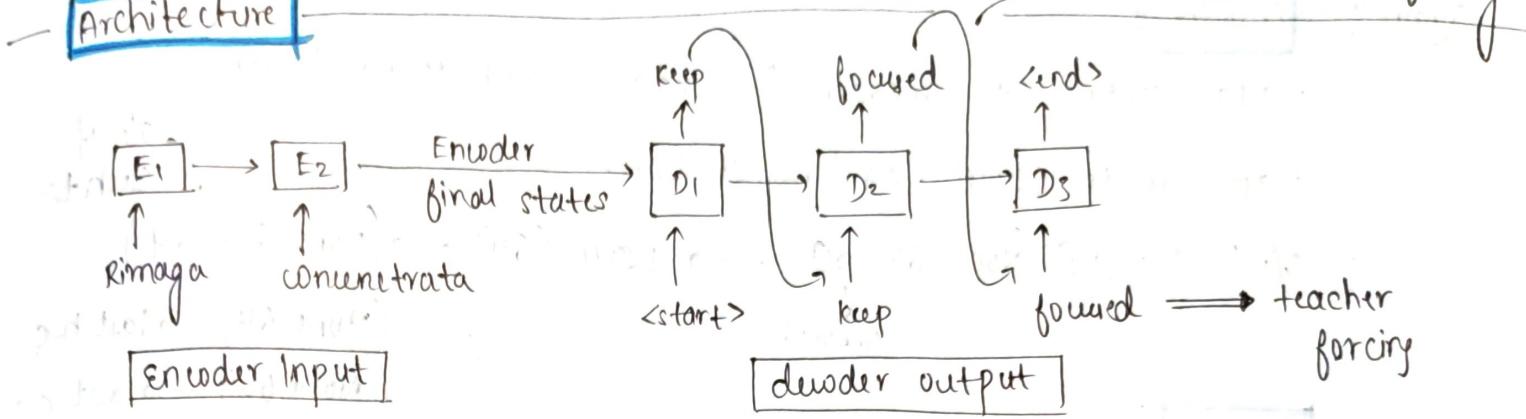
## Sequence - 2 - Sequence

(59)

use case

→ Italian to English translation

### Architecture



**data**

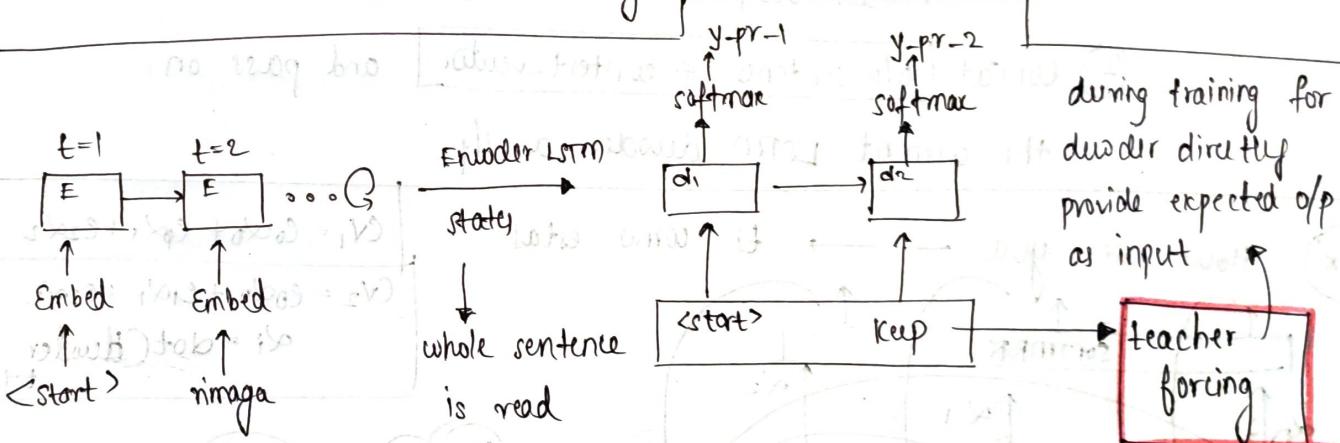
① rimaga concentrato → Encoder input

② keep focused → decoder i/p : <start> keep focused  
decoder o/p : keep focused <end>

↳ apply tokenizer to convert text into integers.

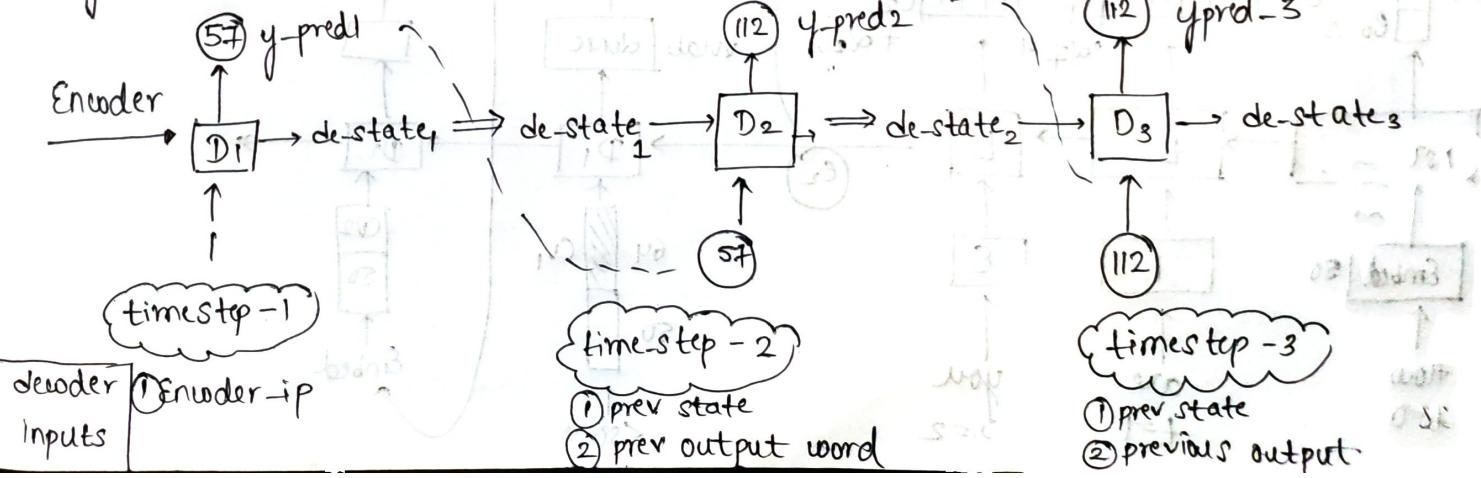
↳ (batch\_size, sentence\_length) → Embedding layer → (batch\_size, sentence\_len, output\_dim)

**Ex**



### Inference

↳ you don't have teacher forcing sentences



## Sq-2-Sq Attention

### Steps

T<sub>train</sub> :- Encoder

→ pass sentence to encoder : words → embedding layer

↳ pass it to LSTM (seq=True, st=True)

O/p: Ex i/p (None, 600, 60) → LSTM (64) → (None, 600, 64)

embeded

o/p. of  
each ts

(None, 64) → last h<sub>t</sub>  
(None, 64) → last c<sub>t</sub>

### Decoder

→ pass desired o/p sentence (+teacher forcing)

at t<sub>s</sub> (t)

① using  $h_{t-1}$  and  $h_e$  find context vector

$h_{t-1}$ : decoder hidden-state  
at prev t<sub>s</sub>

$h_e$ :  $(h_1, h_2, \dots, h_n)$   
encoder outputs

tx = input-length  
(call sentence)

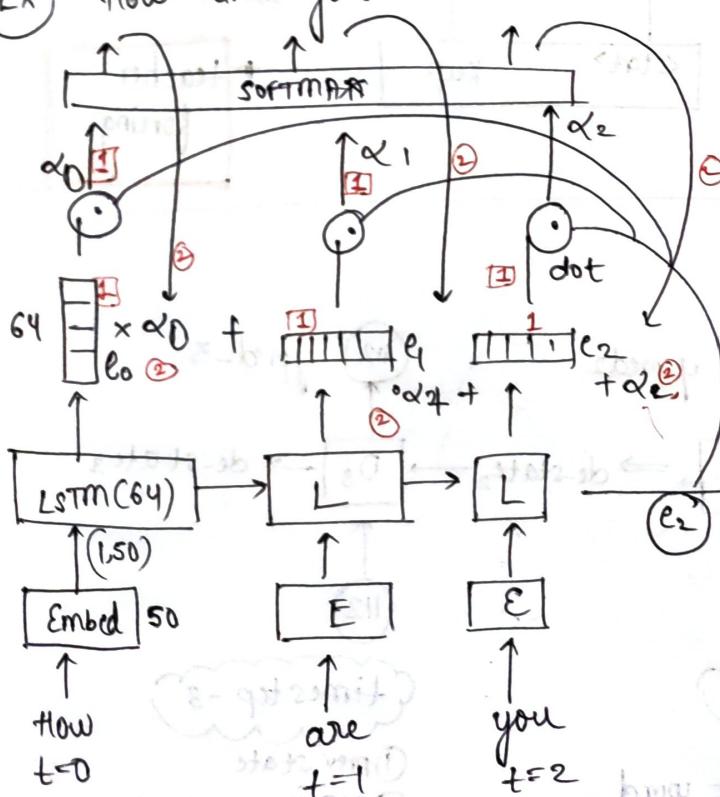
$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \cdot \begin{bmatrix} \alpha_{t1} \\ \alpha_{t2} \\ \alpha_{t3} \end{bmatrix} = \alpha_{t1} h_1 + \alpha_{t2} h_2 + \dots + \alpha_{tn} h_n$$

② concat [o/p sentence + context-vector] and pass on

the current LSTM decoder as i/p

Ex

How are you → ti como estas



$$CV_1 = e_0 \alpha_0 + e_1 \alpha_1 + e_2 \alpha_2$$

$$CV_2 = e_0 \alpha_0 + e_1 \alpha_1 + e_2 \alpha_2$$

$$\alpha_i = \text{dot}(h_{\text{decoder}}, E)$$

## How to find $\alpha_{ij}$

Encoder

inp  $\rightarrow (\text{None}, 600, 50) \rightarrow (\text{None}, 600, 64)$

$(600, 64) \rightarrow \text{all hidden states} \rightarrow h_e$   
 $(64, 1) \rightarrow \text{last h-s}$   
 $(64, 1) \rightarrow \text{last c-s}$

decoder  
 $LSTM = 64$

at time  $t$

① o/p from prev decoder step  $= (h_{t-1}) (64, 1)$   $\rightarrow$  what if decoder  $LSTM(128)$ ?

Now  
 $h_{t-1}(h_{t-1}) = (600, 64) \cdot (64, 1)$   
 $= (600, 1)$

$\therefore$  you get 600  $\alpha$   $= [\alpha_0 \dots \alpha_{600}]$

$\rightarrow$  pass it to softmax so that

$$\sum \alpha_i = 1$$

\* use a matrix  $W$  of shape  
 $(\text{enc-units} \times \text{decoder-units})$

$$h_e(k) \cdot h_{t-1} = (600, 64)(64, 128)(128, 1)$$

$\rightarrow$   $600, 128$  in bottleneck  
 $\rightarrow$  learned by backprop.

Overall

inp  $\rightarrow$  encoder  $\rightarrow$  (calculate)  $\alpha_{ij}$   $\rightarrow$  decoder  $\rightarrow$  output

attention weight of  
each word

$\rightarrow$  to tpus to put no effect words to ppus

pixel distribution

soft max  $\rightarrow$  attention  $\rightarrow$  ppus



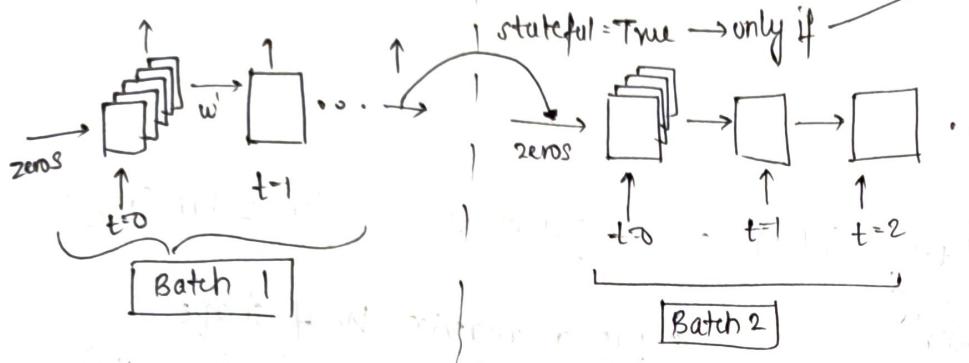
## char-RNN concepts

group-wise  
or continuity

### Stateful LSTM

→ used when batches are sequential

↳ `LSTM(5, stateful=True)`



Batches must be sequential

1, 2, 3, 4	5, 6, 7, 8
9, 10, 11, 12, 13	14, 15, 16, 17
:	:
Batch 1	Batch 2

Here stateful can be used

\* while training batch 1, initially weights are zeros and are gradually updated  
They are discarded at last step

+ when batch 2 is trained again initial weights are zeros

(Imp) ↳ Instead we can use the weights which are obtained from last unit of batch 1 ← stateful

[use] if dataset is proper,

we can learn longer RNN sequence

Ex: given ex dataset

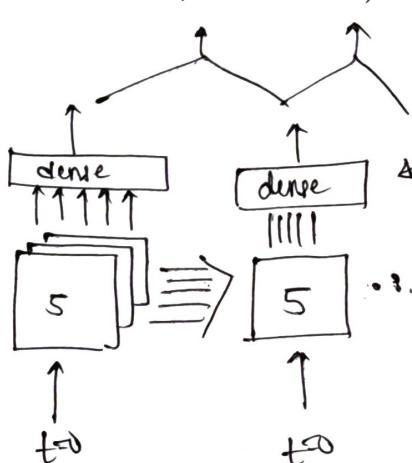
we would train sentence/sequence  
of len=4 (0-4-5-8)

→ \* with stateful we can learn longer  
ie (0,8)  
(9,17)

### TimeDistributed Layer

→ apply a dense layer on top of output of each LSTM timestep

↳ `LSTM(5, ret-seq=True)`



Time distributed  
dense layer

TimeDistributed (Dense(vocab-size))

Ex: `LSTM(64, ret-seq=True)`

$(\text{Batch}, 18) \rightarrow (\text{batch}, \text{time-steps}, 64)$

$(\text{Batch}, 18) \xrightarrow{\text{LSTM}} (\text{batch}, \text{ts}, 64)$

$\downarrow \text{TimeDist}(\text{dense}(16))$

$(\text{batch}, \text{ts}, 16)$

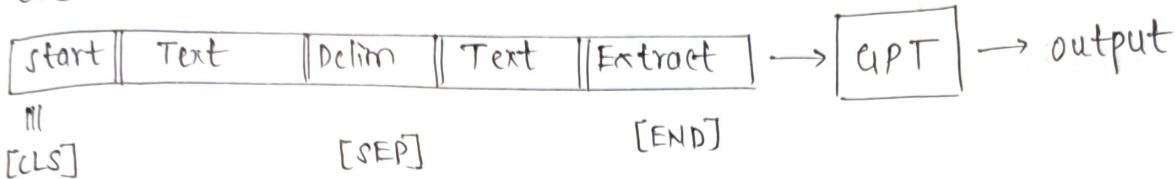
# Generative Pre Training

## EPT-1

Architecture: Decoder only transformer  
 ↳ 12 decoder cells (pg 56)

Training: Next word prediction  
 → give word ( $w_i$ ) predict ( $w_{i+1}$ )

Tokens



## EPT-2

Improvements on EPT-1

\* max input length = 1024 (EPT1 - 512)

\* vocab size = 50,257 (EPT1 - 40K)

\* Layer Norm is used @ input of each sub block

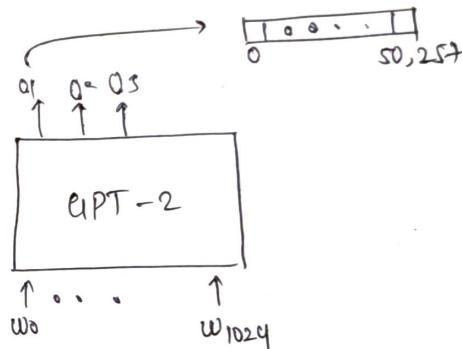
\* larger batch size : 512 (EPT1 - 64)

\* token embedding size =  $d_{model} = 768$

\* used 10x more data to train  
 +

Had 10x more parameters

\* Zero-shot learning



Params	layers	$d_{model}$	
117M	12	768	→ EPT2 small
345M	24	1024	
762M	36	1280	
1.542B	48	1600	→ EPT2 large

① Pretrain model (NWP)

② fine tune model by giving your dataset  
 ↳ No fine tune (Pretrain and directly predict)

③ make predictions

few-shot learning

① → ② → ③

## disadvantages

performs poorly on rare or specialized content

64

### GPT-3

→ 100x more params (175B)

→ trained on

[ 300B + human crowd data  
text tokens (2016-2019) ]

→ uses sparse attention

### disadvantages

- Cost and compute power
- Interpretability
- Biases in predictions

## LIME

# Local Interpretable Model Agnostic Explanations

## Intuition

- + Data =  $\{x, y\}$
- + classifier (model) =  $f(x)$

$\therefore x \rightarrow f(x) \rightarrow y$   
 $\downarrow$   
 may not be  
 interpretable (BERT S12)

## Algorithm

(Step 1): given  $x$ , calculate  $x'$

$$x \in \mathbb{R}^d \rightarrow f(x) \rightarrow y$$

$$x' \in \{0, 1\}^d$$

\*  $x'$  is interpretable

\*  $x'$  is a binary vector (each ele is 0/1 indicating presence or absence).

\* dimensions of  $x$  and  $x'$  may vary

(Step 2): Surrogate model  $g$

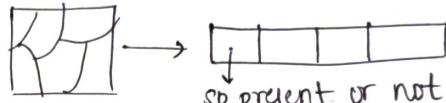
$g \in \mathcal{G}$  (models which are easily  
 interpretable ex: LR, DT)

$x \rightarrow x'$  sentence  $\rightarrow$  BERT  $\rightarrow x$  (S12 dim)

Text :- use BOW

$\rightarrow$  BOW  $\rightarrow x'$

Image :- divide image into super pixels



Tabular data  $\rightarrow$  ① Categories  $\rightarrow$  OHE

② Continuous  $\rightarrow$  binning values.

$\rightarrow \Omega(g) =$  complexity of  $g$   
 (depth, # non-zero weights)

optimization  $f \approx g$ 

$$\mathcal{J}(g) = \min_{g \in \mathcal{G}} L(f, g, \pi_x) + \Omega(g)$$

L: Linear models / DT

L: squared loss

$$\pi_x = \exp\left(-\frac{D(x, z)^2}{\sigma^2}\right) \leftarrow \text{Exponential kernel}$$

proximity:  $\pi_x(z) \rightarrow$  proximity b/w  $(x)$  and  $(z)$

$\pi_x \rightarrow$  locality of  $x$

local fidelity

$L(f, g, \pi_x) \rightarrow$  how best  $g(x')$  approximates  $f(x)$  in locality given by  $\pi_x$

$\downarrow$  better,  $f$  and  $g$  are similar  
 $\uparrow$  bad

$$\left. \begin{aligned} L(\hat{y}, y) &= \sum_{z \in \mathcal{Z}} \pi(z) (f(z) - g(z))^2 \\ z' &\in \mathcal{Z} \end{aligned} \right\} \quad \begin{aligned} \text{for each pt. } z &\text{ find } z' \\ \text{find locality of } \pi_q &\rightarrow \text{Kx} \end{aligned}$$

\* weighted avoiding col proximity

optimization :-  $\hat{w}(g) = \min(L(\hat{y}, y) + \frac{\lambda}{2} \|g\|_2^2)$  Lasso

K-Lasso :- Lasso regularized linear regression with only top-k features

Ex:- S-Lasso :- Lasso regularization keeps only 3 features and makes others 0.  
thus we get only 3 features used for interpretability

### Advantages

- \* works on image, text, tabular data
- \*  $x$  may not be same as  $x'$

### Drawbacks

- \* deciding neighborhood -  $\text{Thc}(z)$
- \* Hyper param =  $\sigma$  ↗  
 $D(x, z)$  may be issue with high-dim data
- \* finding  $(*)$  from  $(*)$  may not be straight forward.

### Unstable

Similar/closer points may give different interpretability.

### Lack of robustness

### Research paper algo

```

algo: sparse linear explanations using LIME
Ip : f(x), N, x → x', Thc(similarity kernel)
K (length of explanations)
algo :- Z ← { }
        for i ∈ {1, 2, ..., N} do
            z'_i ← sample-around(x')
            Z ← Z ∪ {z'_i, f(z'_i), Thc(z'_i) >
        end for
        w ← K-lasso(Z, K) with z' features
        f(z) as target
return w
    
```

Shapley Additive explanations

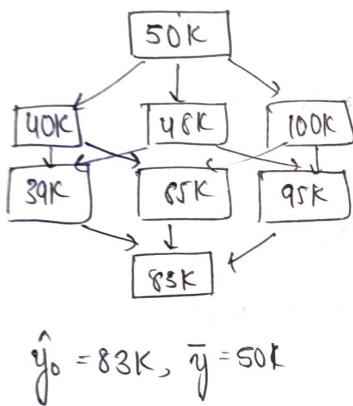
## SHAP

\* Model: Predict salary  
features: age, gender, job  
output:  $y_i$

$x_0$ : query point,  $\bar{y} = \frac{1}{n} \sum y_i$   
 $\hat{y}_0$ : pred value

① for each feature, train  
    (2<sup>nd</sup>) models

② pass  $x_0$  and find predictions for each model



## marginal contribution

$$MC(x_0)_{Age} = \bar{y} - \text{output}(Age)$$

$$MC(x_0)_{Age, Gender} = \text{output}(Age, Gender) - \text{output}(Age)$$

contribution of age when we have only age feature compared to model with {age, gender} feature

$$39 - 48 = -9K$$

③ Shapley values of each feature

$$\begin{aligned} SHAP(Age)(x_0) &= w_1 \cdot MC_{Age, \bar{y}} + \\ &\quad w_2 \cdot MC_{Age, [Age, Gender]} + \\ &\quad w_3 \cdot MC_{Age, [Age, Job]} + \\ &\quad w_4 \cdot MC_{[Age, Gender, Job]} \end{aligned}$$

→ No nodes at level = FC<sub>f</sub>

No of edges =  $f \cdot FC_f$

↳ no of features selected

$$SHAP(\text{feature}) = \sum_{x_0} \left[ \text{edge-weight} (\text{pred}(x) \setminus \text{feature}) - \text{pred}(x) \setminus (\text{feature} \cup \{f\}) \right]$$

$$\text{edge-wt} = \left( \frac{1}{|set|} FC_{|set|} \right)^{-1}$$

$\text{pred}(x) \setminus \text{feature}$  → model trained using feature  $= |set| - \text{feature}(f)$

## property

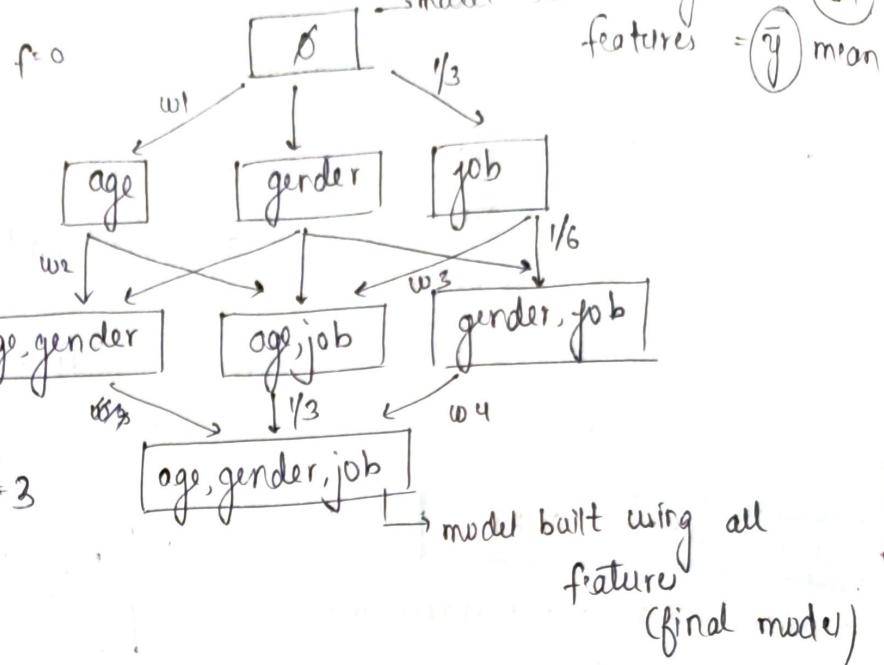
① at a level all  $w_i$ 's are same

② all weights at level are same

$$w_1 = w_2 + w_3 = w_4$$

$$SHAP(Age)[x_0] = -11.33K \$$$

∴ In finding  $x_0 \rightarrow \hat{y}_0$ , contribution of age is  $-11.33K$ .



model built using no features =  $\bar{y}$  mean (67)

# Understanding

$$y_0 \rightarrow \hat{y}_0$$

$$SHAP(\text{age}) = 81 = -11 \cdot 33k$$

$$j_{\text{ob}} = s_2 = -2.33 \text{ K}$$

gender = S3 = 46.60K

$\phi = 50 \text{ k\$}$

Sage, gender, job] = 83k

↳ our model predicted  $\hat{y}_0$  which was 33K away/different from  $\phi/\bar{y}$

$$\text{difference} = 33k = (-11.33k) + (-2.33k) + (46.66k)$$

↓  
wrtib of age

Paper — actual SHAP

\* convert  $x$  to interpretable repr  
 $x \rightarrow r^i \in \{0, 1\}^M$  (BoW)

\* train surrogate model  
 $g(x^*) \approx f(h_m(x))$

① sum of shapely values

$$= y_F - y_P \rightarrow \text{no feature}$$

all feature

Expectations or average value  
of RV

$$E(f(x)) = \sum_x f(x) \cdot p(x=x)$$

$$E(x|y=y_1) = \sum_x x \cdot p(x|y=y_1)$$

$$= \sum_{x \in X} x \cdot p(x=y|x=y_1)$$

## Additive feature attribution

$$g(z) = \phi_0 + \sum_{i=1}^M \phi_i z_i$$

$m = \# \text{ features}$

$$g(z) = \text{surrogate } f^n$$

$\phi_i$  = shapely value for  
feature i

$z^1 \rightarrow$  feature (i) is present  
or not

## KERNEL SHAPES

Kernel shape =  $\text{LIME} + \text{shapely} + \frac{\text{Shapely}}{\text{val}}$

(in UMF we fit interpretable model at last, here instead we find shapely values)

## problems

computationally NP-hard

∴ use approximations

$$\pi_k(z) = \frac{m-1}{m c_{|z|} |z| (m-1|z|)}.$$

$$M = \# \dim \text{in } x_1$$

$$|z_1| = \operatorname{len}(z_0)$$

SHAPELY KERNEL

## OBJECT DETECTIONS

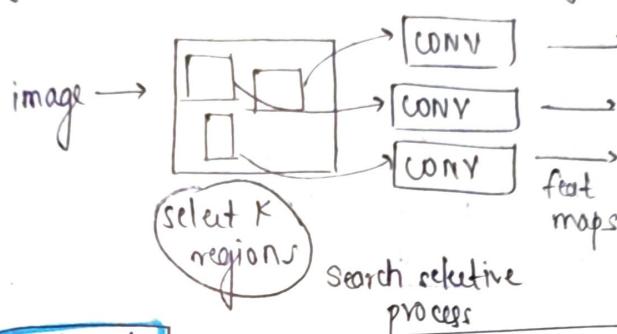
### Regional CNN

→ detect regions in image

69

### R-CNN

↳ slower and low FPS → CONV layer for each region

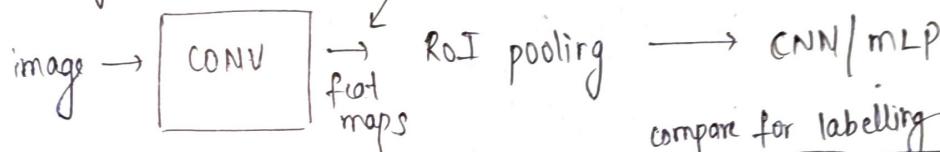


### RCNN

search selective → CONV on each regions

### Fast-RCNN

↳ apply CONV once and then choose feature maps



### Faster RCNN

image → pretrained model

K-pre-defined anchor boxes

feature maps

↓ 3x3 conv, 512 filters or 256

256d

conv 1x1

2K classifier

1x1 conv  
4K regres

work flow

Read img + boxes → get feature maps from a CNN

RPN

(region proposal net)

256d

ds 2K

4K

word

cls

obj

or not

regress

(x,y,w,h)

### ROI pool vs ROI align

can deal with fractional stride by breaking each pixel into further subpixels

ROI = maxpool

↳ causes stride to be quantized  
say ROI Pool (7x7)

$$\text{stride} = \frac{19}{7} = 2 \dots - \boxed{\text{stride}=2}$$

$$\boxed{\text{ip} = 19 \times 19} \rightarrow \boxed{\text{ROI}} \rightarrow \boxed{\text{O/p } 7 \times 7}$$

we are ignoring decimal and hence some part of image

## RPN training

→ if feature map size is  $M \times N$

then classifier weights =  $M \times N \times (K \times 2)$  →  $p_c$  (foreground),  $p_b$  (background)

regressor weights =  $M \times N \times (K \times 4)$  →  $t_i = \langle t_x, t_y, t_w, t_h \rangle$

RPN Loss = log loss on classifier + distance loss for regressor  
 (dist b/w true and pred  $t_i$  ground  $t_i$  given by ground truth)

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i^{\text{CLS}} L(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i^{\text{reg}} p_i^* (L_{reg}(t_i, t_i^*))$$

$N_{cls}$  = no of anchors

$p_i^*$  = actual object present/not

$L_{cls}$  = log loss

$N_{reg}$  = no of anchor locations

$t_i^*$  = actual  $t_i$  of ground truth box

## Non-max suppression

→ we get  $N$  anchor boxes and assign probability to each anchor box

We assign +ve to areas having  $IOU > 0.7$  with any GTBB  
 -ve to areas having  $IOU < 0.3$  with any GTBB

we get final  $(N)$  bounding boxes called Regions of Interest

## OLS for Linear Regression

$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

Loss function =  $\frac{1}{N} \sum (y_i - \hat{y}_i)^2 = \sum (y_i - (\beta_0 + \beta_1 x_i))^2$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\beta_0 = c$$

$$\beta_1 = m$$

$$\frac{\sum y_i}{n} = \bar{y} \quad \frac{\sum x_i}{n} = \bar{x}$$

$$\sum y_i = n\bar{y} \quad \sum x_i = n\bar{x}$$

diff wrt  $\beta_0$

$$\frac{\partial L}{\partial \beta_0} = -2 \sum (y_i - \hat{y}_i) = 0$$

$$\sum (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

$$\sum (y_i - \bar{y} - \beta_1 \bar{x} - \beta_1 x_i) x_i = 0$$

$$\sum (y_i - \bar{y} - \beta_1 (x_i - \bar{x})) x_i = 0$$

$$\sum x_i (y_i - \bar{y}) - \beta_1 \sum x_i (x_i - \bar{x})$$

$$\sum x_i (y_i - \bar{y}) = \beta_1 \sum x_i (x_i - \bar{x})$$

diff wrt  $\beta_0$

$$\frac{\partial L}{\partial \beta_0} = -2 \sum (y_i - (\beta_0 + \beta_1 x_i)) = -2 \sum (y_i - \hat{y}) = 0$$

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\sum y_i - \sum \beta_0 - \sum \beta_1 x_i = n\bar{y} - n\beta_0 - \beta_1 \sum x_i$$

$$\rightarrow n\bar{y} - n\beta_0 - \beta_1 n\bar{x} = 0$$

$$n\beta_0 = n\bar{y} - \beta_1 n\bar{x} \quad (\text{divide by } n)$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

grad descent

$$\frac{\partial L}{\partial \beta_0} = -\frac{2}{N} \sum (y_i - \hat{y})$$

$$\frac{\partial L}{\partial \beta_1} = -\frac{2}{N} \sum x_i (y_i - \hat{y})$$

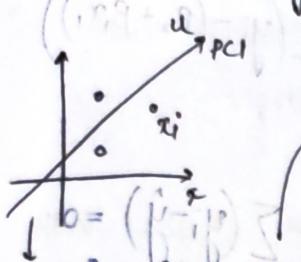
$$\beta_1 = \frac{\sum (y_i - \bar{y})}{\sum (x_i - \bar{x})}$$

multiply and  
divide by  
 $(x_i - \bar{x})$   
 $n$  and  $D$

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

## PCA

mathematical way



it can be proved that  
 $\max(x_i \cdot u) = \text{largest eigen vector}$

minimize fit and maximize variance

$$\text{var}(x_i) = \frac{1}{n-1} \sum (x_i \cdot u - \bar{u}^T \bar{x}_i)^2$$

proj $x_i$  of std

steps

standardize data → covariance mat → find eigen values  $\lambda_1, \lambda_2, \dots, \lambda_n$

$$X_{\text{std}} = \frac{X^T X}{n-1} = S + \text{eigen vectors } v_1, v_2, \dots, v_n$$

also,

$$v_1 \perp v_2 \perp v_3$$

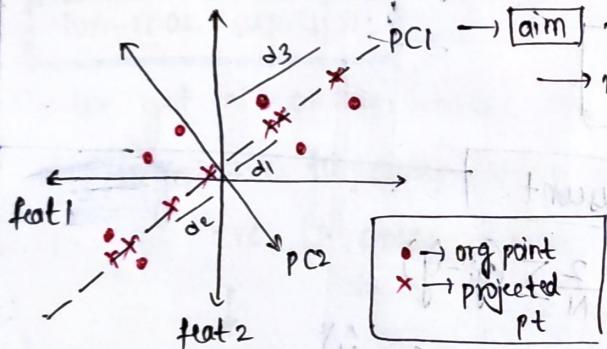
$$v_1 \cdot v_2 = 0$$

$$v_2 \cdot v_3 = 0$$

$$Sv_1 = \lambda_1 v_1$$

$$Sv_2 = \lambda_2 v_2$$

Graphically



standardize data ( $\mu=0$ , centered at  $(0,0)$ )

once PC1 is set,

fit another line fr to PC1  $\Rightarrow$  PC2

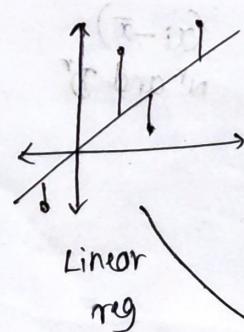
Now we can't fit another line fr to PC1, PC2  
 so for 2d we have 2 PC1

we can use % var / loading score to pick best one

Linear reg vs PCA

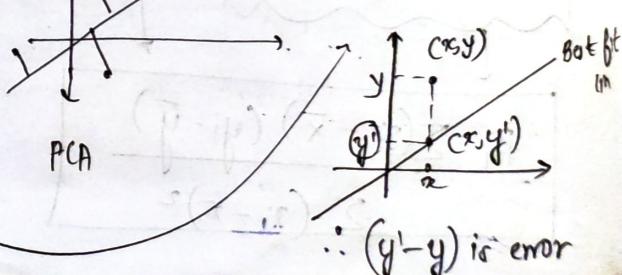
Linear reg → minimize dist b/w point and line

PCA → minimize orthogonal (shortest) distance



you can pick one line along z axis  
 but if you project 2d points, its variance = 0

in LR, we have to predict  $y$  given  $x$



$\therefore (y - \hat{y})$  is error

## XGBoost

### REGRESSION

dosage	age	effectiveness	resid
5	20	13	-1
10	21	14	-5
12	29	12	+10
8	30	11	-8

Bare model not accurate

$$\bar{y} = \frac{13+14+12+11}{4} = 12.5 \quad \text{assume } \sum \text{resid}$$

$$\text{Similarity score} = \frac{(\sum \text{resid})^2}{\text{SS}}$$

$$14 \text{ at node} = (\text{SS left} + \text{SS right}) - \text{SS cur}$$

child nodes

### Prediction

you start at root until leaf node is reached

$$\text{result} = \frac{(\sum \text{resid})}{\# \text{resid} + \lambda}$$

no square  
resid in root node

$$\text{Ex: } [1, -5] \text{ is leaf} \rightarrow \text{resid} = -4$$

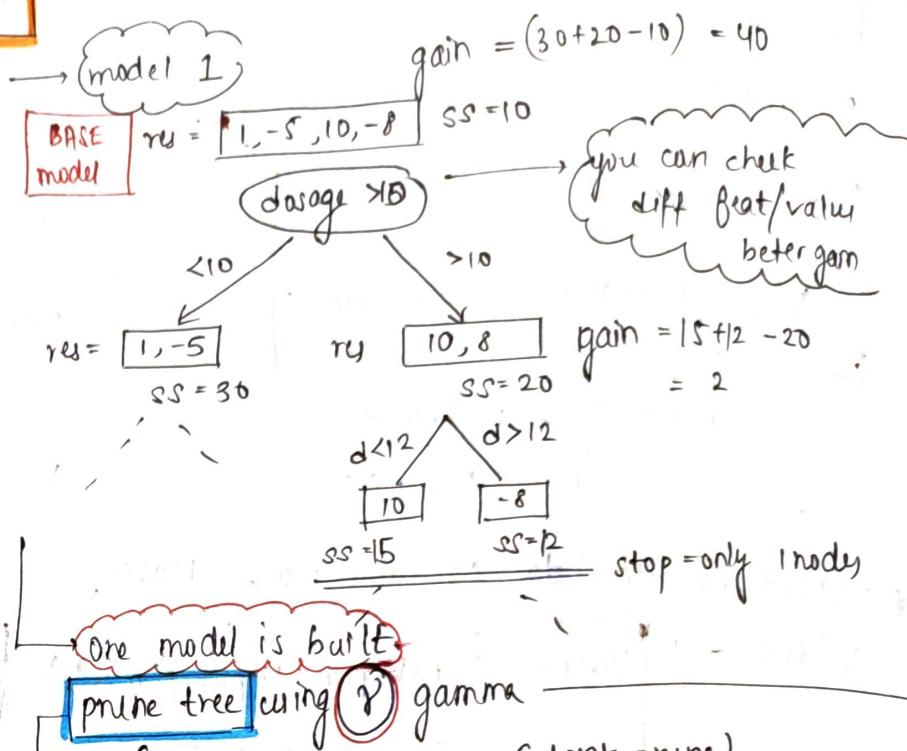
### Regularizations

① as  $\lambda \uparrow$ , SS ↓

as  $\alpha \downarrow$ , gamma ↓

as gain ↓, chance of pruning are here

∴ as  $\lambda \uparrow$  prevents overfit



if  $(\text{gain at node}) - \gamma \geq 0$  (don't prune)

if  $\text{gain} - \gamma > 0 \rightarrow \text{no prune}$   
 $\text{gain} - \gamma < 0 \rightarrow \text{prune}$

if  $\text{gain} > \gamma \rightarrow \text{keep}$   
 else prune

\* while pruning, always start from leaf node.  
 if you can't prune a leaf then you can prune its parent

$$\text{Ex: } \gamma = 5$$

branch gets removed

model 1 built model 2 which fits residuals of (model 1)

$$\text{pred resid}_1 \downarrow \text{resid}_1 = \text{resid}_1 - \text{pred}$$

gamma

?

if  $\gamma \uparrow \rightarrow$  more branches will get pruned

as  $\lambda \uparrow \rightarrow$  more pruning

→ smaller output value

## Classification

Base model ( $p = 0.5$ )

$$SS = \frac{(\sum res)^2}{\lambda + \sum (\text{prev-prob} (1 - \text{prev-prob}))}$$

$$\lambda + \sum (\text{prev-prob} (1 - \text{prev-prob}))$$

probability of previous fit

model (base = 0)

Model 0

Model -0

$$-0.5, 0.5, 0.5, -0.5 \quad SS = 0$$

degree > 10

$$SS = \frac{(-0.5)^2}{0.5(1-0.5) + \lambda}$$

\* xgboost has threshold for min no of leaves in tree →

$$\text{cover} = \left( \sum p_{xi} \cdot \text{prob} (1 - \text{prev-prob}) \right) + \lambda$$

$$\text{cover} = \text{dr of root} - \lambda$$

min-child-split

if at any node  
if tree < cover  
stop growth

∴ at a node, if  $\text{cover}_{\text{node}} > \text{min-cover}$

→ you can use **node** as root leaf  
if  $\text{cover}_{\text{node}} \leq \text{min-cover}$   
→ you can't use **node** as leaf

Ex: min-cover = 1

+ all leaf nodes

$\text{cover}(\text{node}) > 1$

∴ build tree model → leaf nodes must → prune if gain < ?  
with SS ①  $> \text{min-cover}$  ②

Output

$$res(\text{node}) = \frac{\sum res}{(\sum pp(1 - pp)) + \lambda}$$

we need probabilities

$$\log \left( \frac{p}{1-p} \right) \quad \log \text{odds}$$

$$p = res$$

## Overall

$\lambda > 0$ , reduces sensitivity of

tree to outliers

prevent overfit

Pick random subset of data to train models

HINT

→ just like bagging xgboost also uses random rows of data to train

→ increased speed + faster on large datasets

## XGBoost optimizations

### Loss function

#### model - 1

$$L = \left[ \sum_{i=1}^n L(y_i, \hat{y}_i) \right] + \gamma T + \frac{1}{2} \lambda O_{\text{val}}^2$$

$\lambda O_{\text{val}}^2$  → regularization term  
prevent  $O_{\text{val}}$  becoming too large

$\hat{y}_i$  → current out-put

$$= \text{Base model score} + \sum \text{residuals up to } n-1 + O_{\text{val}}$$

$O_{\text{val}}$  = output from current model/tree for point  $y_i$

$T$  - no of terminal nodes  
(used for post pruning)

$$L = \sum_{i=1}^n L(y_i, (\hat{y}_i)_{\text{old}} + O_{\text{val}}) + \gamma T + \frac{1}{2} \lambda O_{\text{val}}^2$$

### Approximate greedy

→ we decide node split based on gain and then just pick one with best gain and split. Once split we don't revisit node again which saves time.

### Parallel + weighted quantile sketch

→ if feature is continuous, calculating all thresholds and checking is slow

→ instead, we pick quantiles among data and then use these quantiles  
Generally, reboots were around 33 Quantile (test split 33 times)

→ sketch algo → they split data to various pcs and later combine to find approximate quantiles (histograms)

↳ in this each histogram region of quantiles will have some no of observations

+ weighting  
(use cover and prob)

↳ weighted Quantile sketcher.

### Cache-aware access

↳ try to store gradients and derivatives in cache as they are frequently accessed

→ if data is too large and can't fit in RAM  
compress it into chunks and save in HDD

→ when needed CPU reads this chunk parallelly → uncompress and use (although zip/unzip takes time but better than slow disk)

### XGBoost pros

Improved Grad Boost \* sparsity aware split

Regularization \* find

approx greedy algo \* Cache awareness

parallel learning \* pick random nodes

weighted quantile sketch like bagging.

\* block out-of-core computation

## Sparsity aware split

assume missing features

↳ initial residue can be found ( $y_i - \text{average}(\text{vals})$ )

dosage	effe	reso
10	9	3
21	2	4
19	7	-2
?	1	-6
?	8	-1

mean

split ? in new table

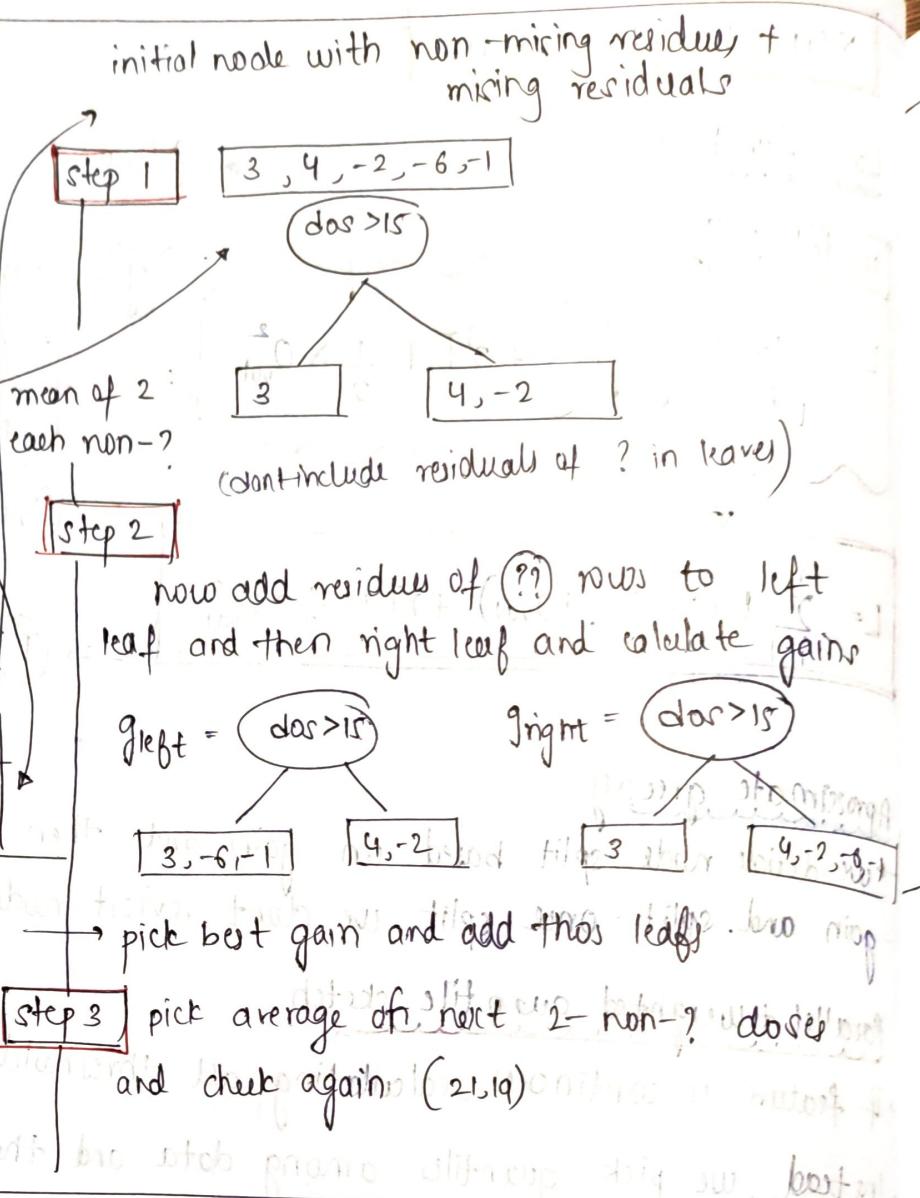
?	1	-6
?	8	-1

10	9	3
21	2	4
19	7	-2

as data is small we use avg of 2  
each non-? rows

else for large dataset,

we use quantiles



TSNE

t-distributed stochastic neighborhood embedding

Core idea: Instead of focusing to

preserve distances, focus on preserving neighbors.

① Stochastic nbrhood embedding

$$L = \sum_{ij \in \Omega} P_{ij} \log \frac{P_{ij}}{q_{ij}}$$

KL divergence b/w  
similarities

\* affinities  $\rightarrow 1$ ; if points are closer  
tends to 0; if far away  
if 2 points close in high-d, are projected  
far away in low-d, loss is greater

$P_{ij}$ : pairwise affinity b/w

$P_{ij}$  and  $p_j$  in high-dim space

$q_{ij}$ : pairwise affinities in

low-dim space

High-dim  
similarities

$$P_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

used such that  $\sum m = 1$

Perplexity → how many neighbors to preserve (around 30)

$$P = 2^H; H = -\sum_{j \neq i} P_{j|i} \log_2 P_{j|i}$$

$$w_{ij} = K(\|y_i - y_j\|^2)$$

$$P_{ij} = \frac{P_{ilj} + P_{lji}}{2n}$$

$$\text{loss fn} = \sum_{ij} P_{ij} \log q_{ij} + \sum_{ij} P_{ij} \log \frac{w_{ij}}{Z}$$

→ low-dim similarities

$$q_{ij} = \frac{K(\|y_i - y_j\|^2)}{Z} \quad \text{where, } K = K(d) = 1/(1+d^2) \quad [t\text{-distributed kernel for calculating similarities}]$$

$$Z = \sum_{k \neq l} K(\|y_k - y_l\|^2)$$

decision tree  
with 1 root, 2 leaf

steps

ADABoost

KEY POINTS

① Base learners are weak (stumps)

② Each BL is dependent on information from prev

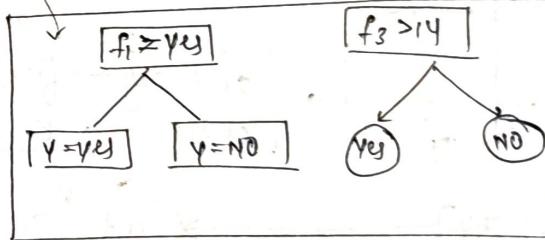
③ Each BL may have diff weight in final voting

	f <sub>1</sub>	f <sub>2</sub>	f <sub>3</sub>	y	wt
1	Y	0	-	Y	1/5
2	N	1	-	Y	1/5
3	Y	0	-	N	1/5
4	Y	2	-	N	1/5
5	N	2	-	N	1/5

→ fit a BL (stump)

→ create stump for f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>  
& check which gives more gini index

→ that stump is BL<sub>1</sub>



② update wt

→ predict ④ using BL<sub>1</sub>

ex say ④, ⑤ were incorrectly  $\Rightarrow TE = \frac{1}{8} + \frac{1}{8} = \frac{1}{4}$   
classified by BL<sub>1</sub>

Total Error =  $\sum$  weights of incorrectly classified samples

$$\text{amt-say} = \frac{1}{2} \log \left( \frac{1-TE}{TE} \right)$$

normalize after calculating wt<sub>i</sub>

so that  $\sum wt = 1$

③ new weights wt<sub>i</sub>

penalize incorrect more

normalize entire col

if sample

incorrectly classified  $\rightarrow wt = wt \times e^{-\text{amt-say}}$

correctly classified  $\rightarrow wt = wt \times e^{\text{amt-say}}$

$$\sum x_i \rightarrow x_i = \frac{x_i}{\sum(x)}$$

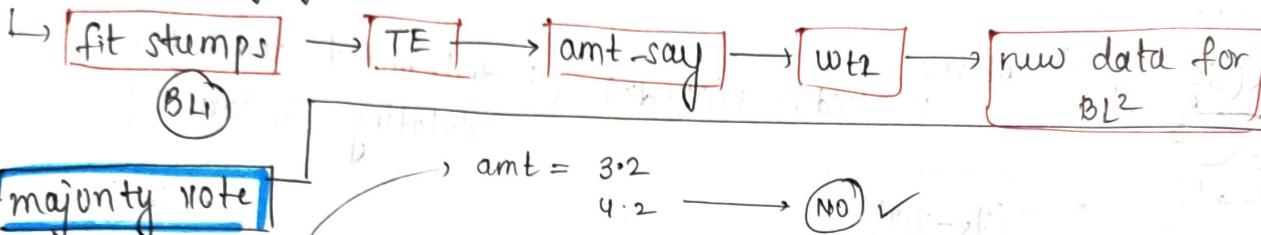
→ Next BL →

→ create new dataset with same size

→ say  $\text{len}(\text{dataset}) = 5$

\* pick 5 rows from datasets randomly, such that prob of each row getting picked is proportional to  $w_t$  value (proportional sampling)

\* some rows may get repeated

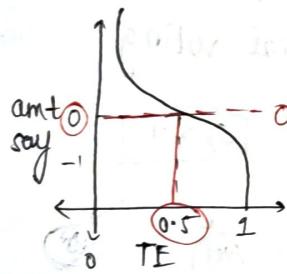


say 3 BL → Yes  
3 B → No → find  $\sum \text{amt-say} (\text{BL}_1)$  → which gives  $\sum \text{amt-say} (\text{BL}_2)$  more pick that result

### Observations

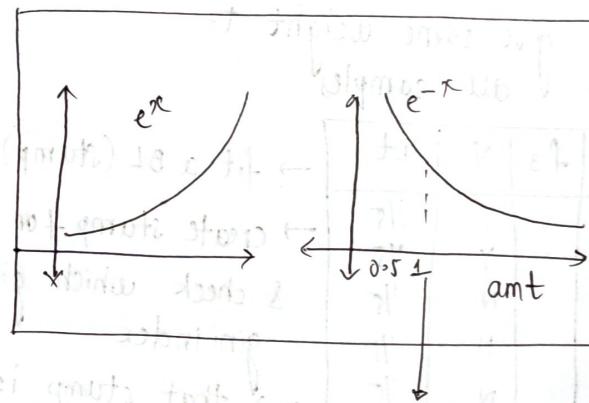
$$\text{AMTE} = \frac{1}{2} \log \left( \frac{1-x}{x} \right)$$

$$\text{amt-say} = \frac{1}{2} \log \left( \frac{1-\text{TE}}{\text{TE}} \right)$$



amt-say → fixed for BL  
change →  $e^{-x}$  or  $e^x$

error ↑ → amt-say ↓



AMT-SAY

\* If  $\text{TE} = 1$ , (worst)

amt-say < 0

(more error, less amt-say)

\* If  $\text{TE} = 0$  (best)

amt-say > 0

\* If  $\text{TE} = 0.5$  (random model)

amt-say = 0

$\text{TE}$  is small (better BL)

→ update INC less

→ reduce wrong msk

wt updates → +ve classified (current)

(wt,  $e^{-\text{amt-say}}$ )

→ more amt say,  
lesser is  $e^{-x}$

∴ if amt-say >>  
wt update very little  
(it may be < wt<sub>i-1</sub>)

-ve classified  
incorrect

wt ·  $e^{\text{amt-say}}$

→ more amt-say,  
more is  $e^x$   
∴ wt update is more

Idea

\* say  $\text{TE}$  is large,

→ then incorrect pts wt will be increased more  
as we want to predict them in next stump  
correct pts wt will reduce because they  
have already been correctly clif by BL<sub>1</sub>