

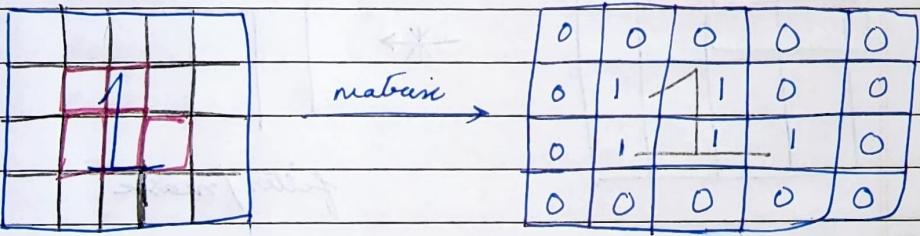
~~Book~~

Deep learning by Ian Goodfellow

18/09/22

Convolutional Neural Network

Image - basically a signal represented in the form of a matrix of discrete values



Image

0 → bg

1 → has parts of the
img of 1

why CNN instead of ANN?

1. Computationally expensive to use ANN for images

For images have many features, say 1024.

Say image is 1024×1024 neurons \Rightarrow if there are 1024 neurons in H_1 and 1000 in H_2 , then already $1024 \times 1024 \times 1000$ computations are being carried out. (flatten the image and pass it as an ip)

Grayscale image:

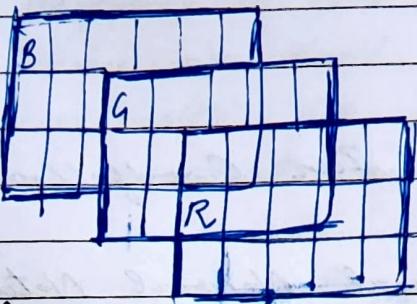
0 - black

255 - white

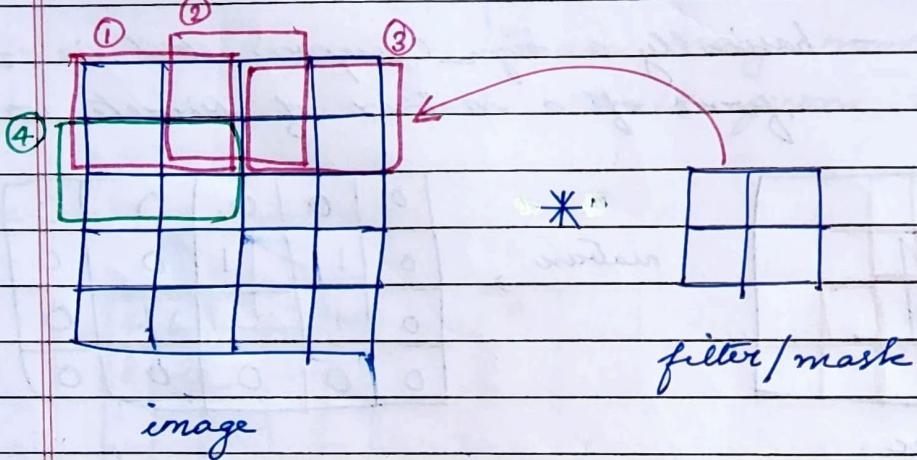
any value b/w 0 & 255 - gray

Color images :

3 channels typically (RGB).



Convolution :



Place the filter on the left end of the img
do some computation & move the filter 1 step
to the right (say stride is 1)

once a row is complete, move to the next
row's left end and repeat (assuming
stride is 1)

(Ex):

a	b	c
d	e	f
g	h	i

*

w	x
y	z

=

aw + bx	bw + cx
+dy + ez	+ey + fz
dw + ex	ew + fx
+gy + hz	+hy + iz

img * filter

element-wise
multiplication
(not matrix multiplication)

$$(3 \times 3) * (2 \times 2) \rightarrow (2 \times 2) (2 \times 2)$$

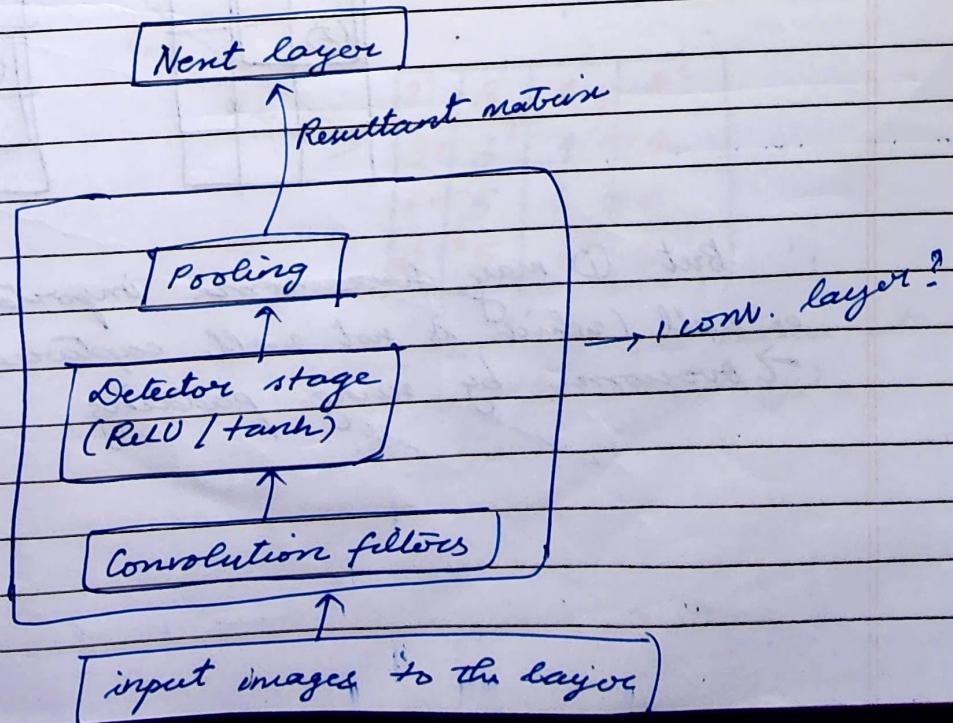
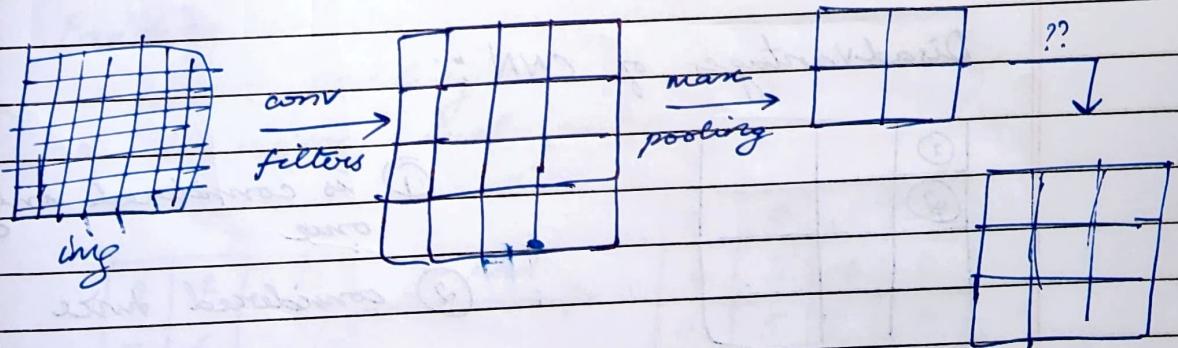
Typically, the felt image is convolved many times with the filters before activation fn is applied

Keras eg :

Dense(filters = 16, kernel = (3,3), input_shape = (240, 240), channels = 3).

Pooling : done to reduce matrix size

- ① Pass over input img to conv layer (with a number of filters) after which we a
- ② apply ReLU
- ③ apply max Pooling



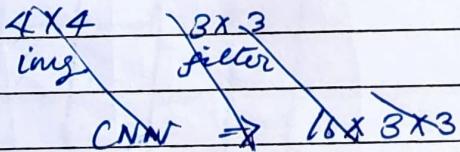
Date _____
Page _____

1st layer (convolution) detects edges / outline
 next detects little more complex features
 like 2D shapes (circles, squares, etc.)
 3rd layer detects shapes eyes, etc..
 and so on

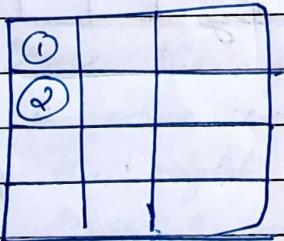
(each subsequent layer detects more & more complex features in the image).

feature learning via convolutions

- * Number of parameters in CNN is much less than that of ANN

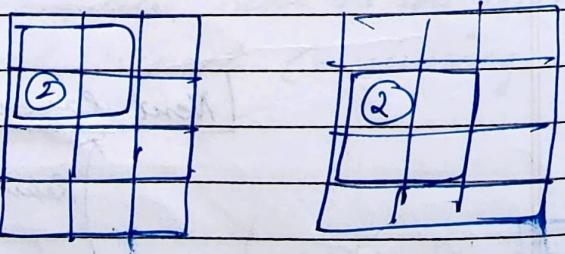


Disadvantages of CNN:

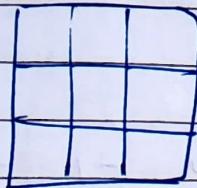


① is considered only one

② considered twice



But ① may have some important feature as well (which is not well captured in the ①).
 → overcome by using padding



img

Padding →

0	0	0	0	0
0	1			0
0				0
0				0
0	0	0	0	0

Now 0 is also considered since
→ 1

0	0
0	1

0	0
1	0

0	1
0	0

1	0
0	0

⇒ 4

considered

4 times
like all other
pixels.

Padding :

1. Zero padding
2. boundary level padding

3	4
5	6

0 pad →

0	0	0	0
0	3	4	0
0	5	6	0
0	0	0	0

boundary level

3?	3↑	4?	4?
3←3		4→4	
5←5		6→6	
5?	5↑	6?	6?

(Just copy the values at
the boundary).

13/09/22

Date _____
Page _____

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}_{6 \times 6} * \begin{bmatrix} +1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{bmatrix}_{3 \times 3}$$

$$= \begin{bmatrix} 3 & -1 & -3 & -1 \\ -3 & 1 & +1 & -3 \\ -3 & -3 & 0 & 1 \\ 3 & -2 & -2 & -1 \end{bmatrix}_{4 \times 4} \quad (\text{stride} = 1)$$

stride = 2

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}_{6 \times 6} * \begin{bmatrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{bmatrix}_{3 \times 3}$$

stride reduces pixel overlaps and the output volume
 ↓
 b/w img and filter

$$\begin{bmatrix} 3 & -3 \\ -3 & 0 \\ -3 & 0 \\ 3 & -2 \end{bmatrix}_{4 \times 2}$$

After Q, if we move the filter by a stride of 2, it goes beyond the image, so we skip it & move more to the next

row.

$$\begin{bmatrix} -1 & 1 & -1 \\ -1 & 1 & -1 \\ -1 & 1 & -1 \end{bmatrix} \rightarrow \text{detects vertical boundaries}$$

$$\begin{bmatrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{bmatrix} \rightarrow \text{detects a diagonal boundary}$$

To detect outline / edges

→ use kernel $\begin{bmatrix} 1 & -1 \end{bmatrix}$
 edge ↓
 ↓
 not an
 edge

Pooling :

Performs downsampling (reduction in image size)

Types - Max, min, average

Responsible to detect important information
 (helps in identifying different variants
 of similar images)

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

$\xrightarrow{\text{maxpooling}}$
 $\xrightarrow[2 \times 2]{\text{stride} = 2}$

6	8
3	4

Padding - performs upsampling

5 can be written in two different handwritings \Rightarrow Pooling helps in identifying that all are similar (belong to class 5).

Each conv layer uses many filters and stacks the ops depthwise so depth / channels inc.

but the dimension of each channel \downarrow

$$(\text{Ex}): \underline{160 \times 160 \times 3} \rightarrow \underline{80 \times 80 \times 9}$$

due to 3 filters?

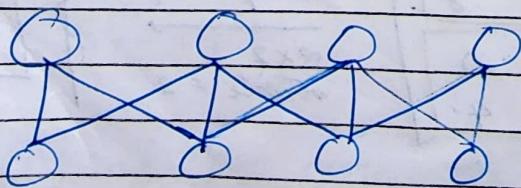
* CNN + fully-connected network

After using convolution, we flatten the matrix and pass it to a fully connected (why?) ~~or~~ network

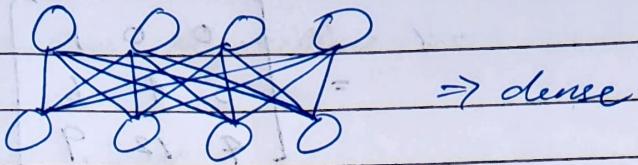
Activation \rightarrow generally softmax to predict the class with that the img belongs to

Sparse matrix

Sparse connectivity



\Rightarrow only some connections & not all dense



$$\text{Img size} = N \times N$$

kernel size = K x K

padding = p , stride = ~~s~~ S

$$\text{Output image size} = \frac{(N - K + 2P + 1)}{S} \times$$

$$\left(\frac{N-K+2P}{n} + 1 \right)$$

16/09/22

Padding = 1

$$stoxide = 1$$

$$\text{kernel} = 3 \times 3$$

$$\frac{N-K+2P}{S} + 1 = \frac{28-3+2}{1} + 1$$

$$= 98$$

(if you get fractions, take the floor value).

Transported Convolution:

✓ Upsampling technique

$$\begin{aligned}
 (\text{eq}): & \quad \left[\begin{array}{c|cc} 0 & 1 \\ \hline 2 & 3 \end{array} \right] \quad \left[\begin{array}{c|cc} 0 & 1 \\ \hline 2 & 3 \end{array} \right] \rightarrow \left[\begin{array}{c|cc} 0 & 0 & 1 \\ \hline 0 & 0 & 1 \end{array} \right] + \left[\begin{array}{c|cc} 0 & 1 \\ \hline 2 & 3 \end{array} \right] \\
 & \quad \text{img} \quad \text{kernel} \\
 & \quad \text{input} \quad \text{filter} \\
 & \quad \text{ip} \times 2 \quad \rightarrow \quad \left[\begin{array}{c|cc} 0 & 2 & 0 = 1 \\ \hline 4 & 6 & 0 = 9 \end{array} \right] + \left[\begin{array}{c|cc} 0 & 3 \\ \hline 6 & 9 \end{array} \right]
 \end{aligned}$$

$$= \begin{bmatrix} 0 & 0 & 1 \\ 0 & 4 & 6 \\ 4 & 12 & 9 \end{bmatrix}$$

1. Apply padding
2. Do convolution

(Eg): $ip = 3 \times 3$
filter = 3×3 .
reqd op size = 5×5 .

$$\frac{N - K + 2P}{S} + 1 = 5.$$

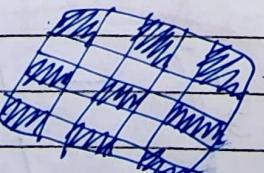
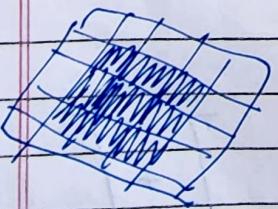
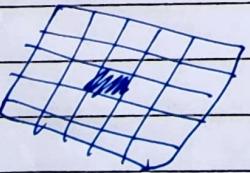
$$\frac{3 - 3 + 2P}{1} + 1 = 5. \quad \text{(padding)} \quad \text{1 = stride}$$

$$2P = 4.$$

$$P = 2$$

Dilated Convolution

Do padding in btw the pixels also (& not just the boundary).



3×3 padding

dilation = 1

(leave $1-1=0$

gaps b/w cells)

3×3 padding

dilation = 2

(leave $2-1=1$

gap

Convolution Algorithms

Stride Conv

Padding Conv

Transpose Conv

to Dilated Conv

Parameters of Conv layer

- ✓ How many filters
- ✓ footprint (size of the filter - $3 \times 3, 4 \times 4, \dots$)
- ✓ striding (Reqd or not, if yes how much)
- ✓ activation fn
- ✓ Padding

Op of conv layer \rightarrow feature map

- ✓ learning rate
- ✓ loss fn
- ✓ batch size
- ✓ initial weights
- ✓ loss fn

ImageNet \rightarrow dataset?

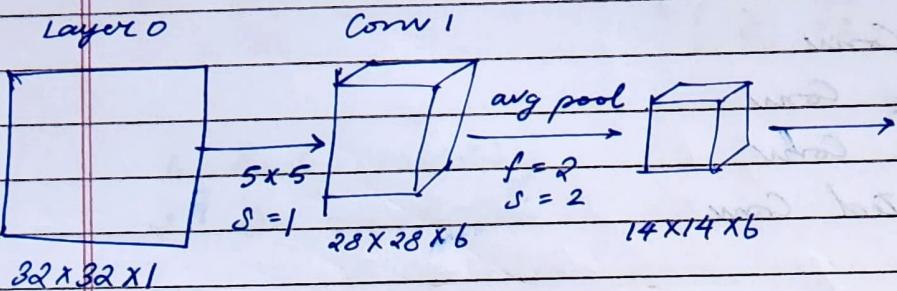
Initializing filter weights

Symmetrical filters \Rightarrow do the same job
 \Rightarrow Redundant

Approaches :

1. Random initialization
2. Glorot initialization
3. He initialization

LeNet 5



Pooling

✓ reduces size of ip to speed up processing

✓ gives some robustness to feature detection
 $512 \times 512 \Rightarrow$ lot of pixels (use pooling)

Conv layer

✓ convolve using filters

✓ Add bias

✓ Apply activation fn

Transposed Convolution :

was like conv. with fractional striding
 upsampling by just padding

Stride :

fast way to ↓ img size to speed up later blocks in the network

Why padding ?

- ✓ conv. reduces img size (use padding to make op size same as if or even ↑)
- ✓ pixels at edge are not used only in fewer convolutions than central pixels
- ✓ throws away a lot of information at the edges

Types of layers in CNN

- ✓ Conv
- ✓ Pooling
- ✓ FC

downsampling using strideed conv better than
conv + pooling

Deeper in the network

size of img ↓ { we use lower resolution to
channels ↑ } work with large collection
of features without
requiring large filters

~~Always~~ # channels in filter = # channels in input

23/09/22

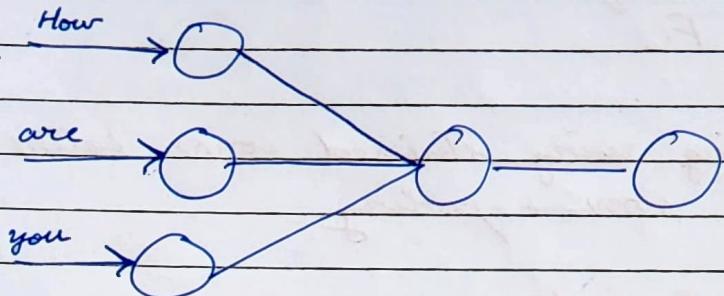
CLASSMATE

Date _____

Page _____

RNN

Sequence data Why RNN + not ANN?



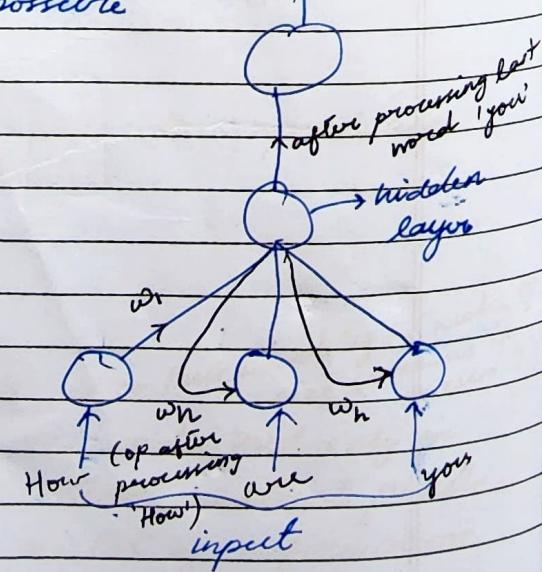
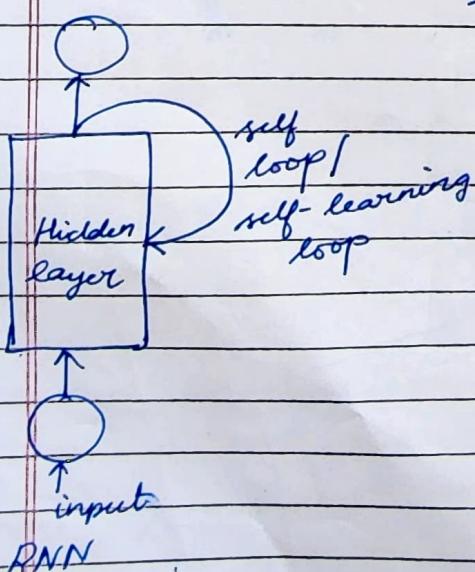
- ① Semantics are not preserved in ANN (when using text data).
(~~How~~ processed separately; ~~are~~ processed separately
and so on).

- ② Computationally expensive to use ANN
(sequence processing involves too many computations)

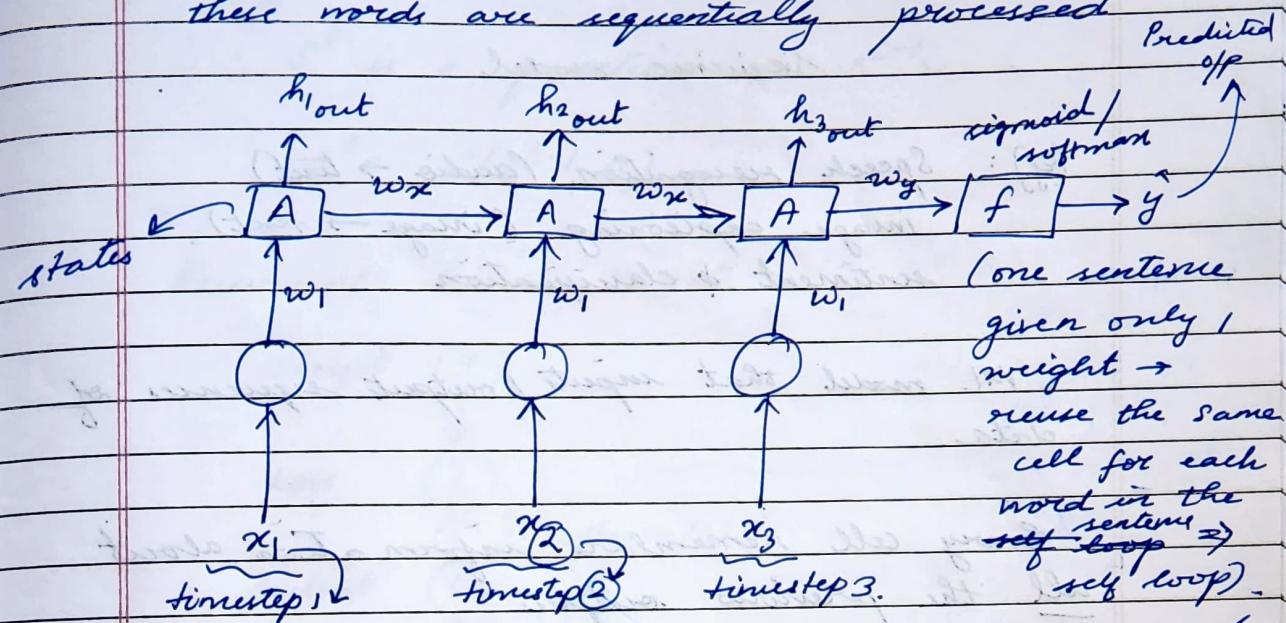
Encode $\rightarrow [1\ 0\ 0] \ [0\ 1\ 0] \ [0\ 0\ 1]$. (One-hot encoding).

- ③ No memory

- ④ Parameter sharing op
is not possible



RNN \rightarrow sentence is split into individual words and these words are sequentially processed



(previous output and current input are used to make current prediction)

$$h_{1\text{out}} = f(w_i x_1 + b) \quad f \text{ is activation fn}$$

$$h_{2\text{out}} = f(w_i h_{1\text{out}} + w_x x_2 + b).$$

$$h_{3\text{out}} = f(w_i h_{2\text{out}} + w_x x_3 + b).$$

$$\hat{y} = f(h_{3\text{out}} w_y + b)$$

f can be sigmoid / softmax

weights and activation fn's are same because we are processing the words in the same sentence and it performs the same task

sequential data :

Data that has to be processed

one after another.

(Eg) : videos, documents

$x_i \rightarrow$ word input at i^{th} timestep

Sequence model

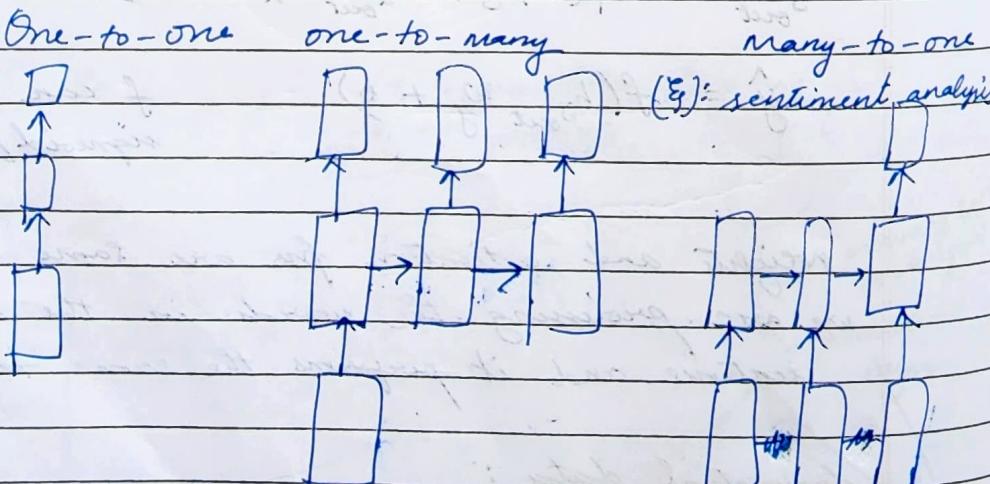
(Eg): Speech recognition (audio \rightarrow text)
 image captioning (image \rightarrow text).
 sentiment & classification

ML model that inputs / outputs sequences of data

Every cell remembers information about all the previous outputs

Type of RNN
 Language translation \rightarrow many-to-many
 Movie rating \rightarrow many-to-one
 generate poem from { one word } \rightarrow one-to-many

Product review (good or bad) \rightarrow many-to-one



(Eg): Music generation

sentiment analysis

language generation

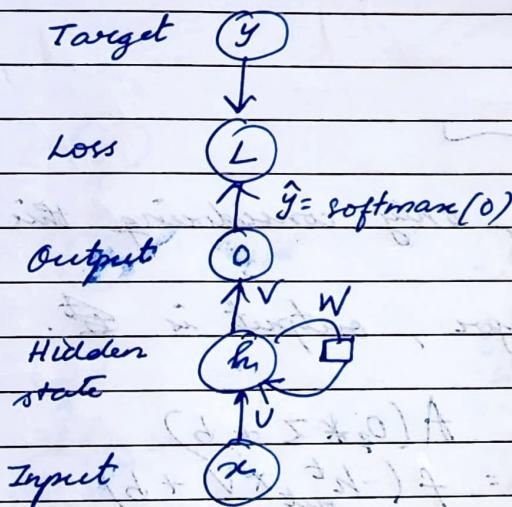
sequence input
single output
(Movie rating)

A system is recurrent if

s^t is a function of s^{t-1}

$$s^t = f(s^{t-1}, \theta)$$

Recurrent hidden units



for every time step

Same context \Rightarrow same weights, activation fn.

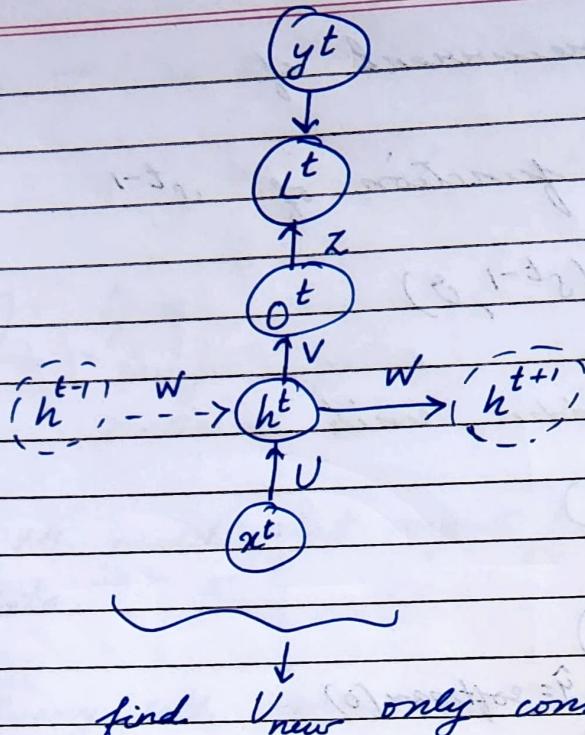
RNN loss

1. compute at individual time steps
2. overall loss after all time steps

Backpropagation happens through time

BPTT

\Rightarrow Backprop applied to RNN unfolded in time



find V_{new} only considering this layer

For that layer, output is \hat{y}^t .

$$\hat{y}^t = A(o_t * z + b)$$

$$o_t^t = (h_{out}^t * V + b)$$

$$h_{out}^t = f(h_{in}^t)$$

$$h_{in}^t = (x^t * U + b)$$

$$\text{say } L = (y^t - \hat{y}^t)^2$$

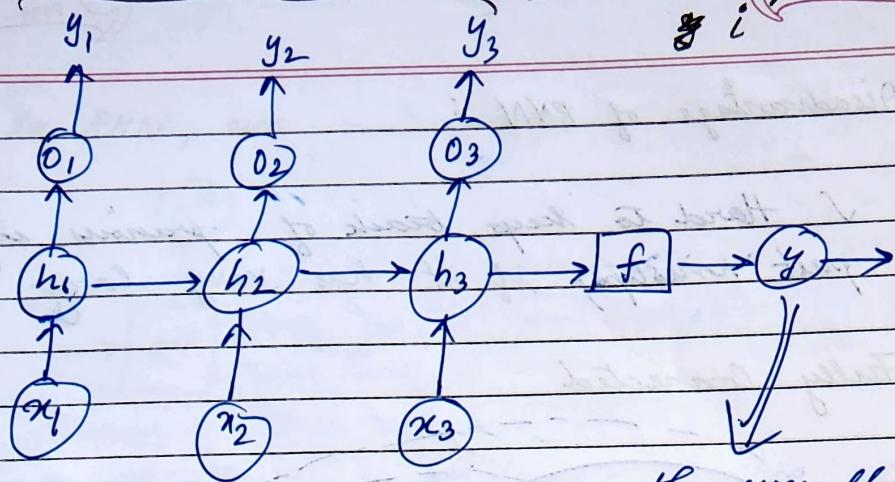
$$\frac{\partial L}{\partial U} = \frac{\partial L}{\partial \hat{y}^t} \times \frac{\partial \hat{y}^t}{\partial o^t} \times \frac{\partial o^t}{\partial h^t} \times \frac{\partial h^t}{\partial h_{out}^t} \times \frac{\partial h_{out}^t}{\partial h_{in}^t} \times \frac{\partial h_{in}^t}{\partial U}$$

Individual loss - useful to fill blanks in sentences (language generation?)

(Use output from layer 1, predict word for layer 2 and so on..)

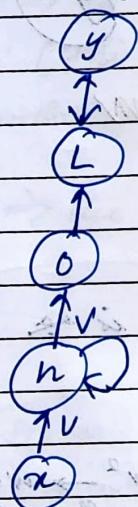
each y_i is used to compute individual loss at timestep i

classmate
Date _____
Page _____



the overall output
(used to compute
overall loss)

(Eq):



$$g^{(t)} = \text{softmax}(o^{(t)})$$

$$o^{(t)} = c + v h^{(t)}$$

$$h^{(t)} =$$

$$a^{(t)} = b + W h^{(t-1)} + U x^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

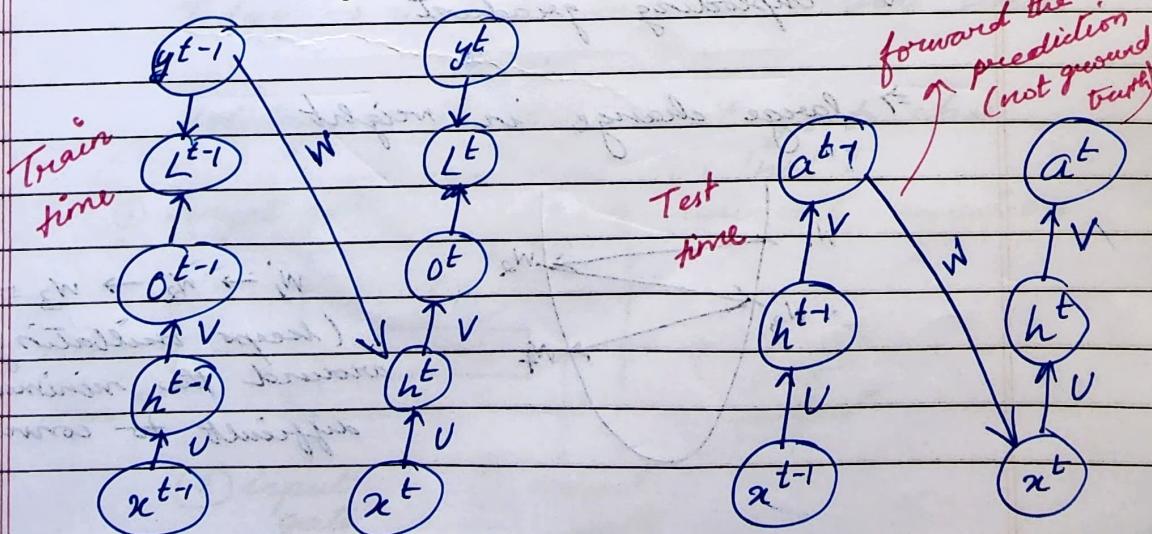
$$o^{(t)} = c + V h^{(t)}$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$

$a^{(t)}$ is $h^{(t)}$

$b, c \rightarrow$ bias.

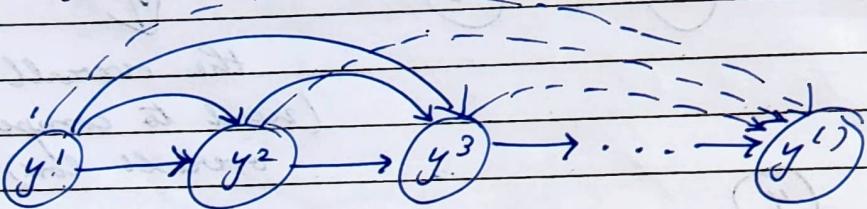
Teacher forcing



Disadvantage of RNN :

- ✓ Hard to keep track of previous info (at past timesteps) if it has very large inputs

Fully Connected



⇒ main issues with standard RNN

- ✓ Vanishing gradient

⇒ change in weight value is ~~is~~ very less

$$w^1 \rightarrow w^2 \rightarrow w^3 \\ (\text{very small update})$$

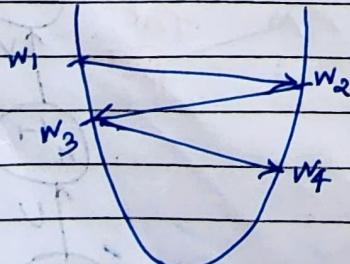
→ model does not learn even after many epochs.

- ✓ Exploding gradient

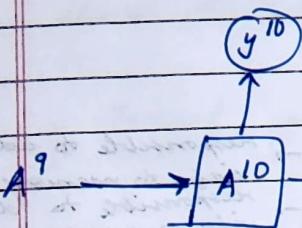
⇒ large change in weight

$$w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow w_4$$

(keeps oscillating around the minima difficult to converge).



In RNN, we



y^{10} resembles A^{10} more than the previous information (upto A^9).

→ so it does not remember past info well.
(This happens more & more as we go further through more & more timestamps).

→ Use LSTM

27/09/22 LSTM
(long-short term memory).

→ Stores whatever is important in a sentence in a memory and use it whenever needed

(Eg): Today is holiday, so students in class. are not

Keywords : Today, Holiday
(both are stored in memory)

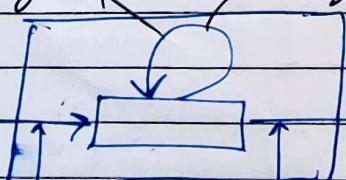
Today → used to indicate present tense
(are)

Holiday → so they won't go to class
(clue for 'not').

① forget gate

→ self-relevant connection

memory
cell
input



memory
cell
output

② input
gate

③ output
gate

LSTM - 4 interacting layers



discovered in 1976

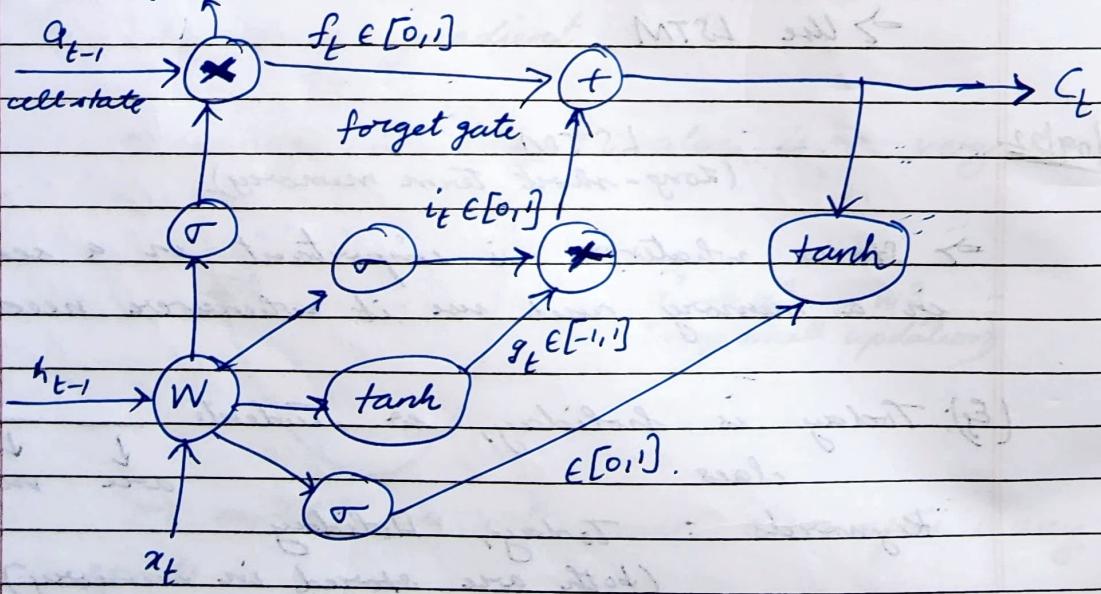
uses gating mechanism

(3 gates → input gate → responsible to add new info to memory
 output gate → responsible to decide what info from memory goes to forget gate - remember info if next important - important else forget

tanh - selects inputs

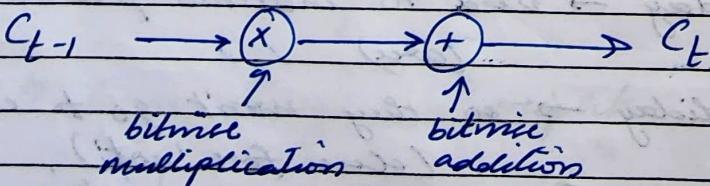
output from LSTM - 4x larger than

elementwise (or) bitwise multiplication that of RNN cell



State vector

cell state - represents memory of LSTM

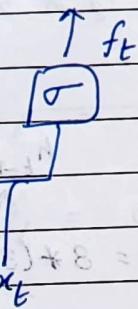


Gates

Have the ability to add or remove

1. Forget gate - controls what info to store
throw v is memory away from

remove
info which
is not important
from the
memory



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

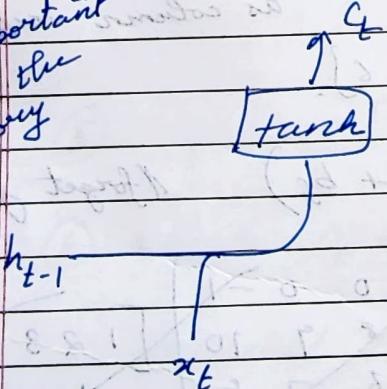
dimension is dimension of
 h_{t-1} + dimension of x_t

- $f_t = 0 \Rightarrow$ forget
 $f_t = 1 \Rightarrow$ remember completely

Remembering gate

2. Input gate → decides what info to put inside the memory

puts important
info into the
memory



(forget gate opp +
input → input
gate output)

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

$$\tilde{Q}_t = \tanh(W_i \cdot [h_{t-1}, x_t] + b_i)$$

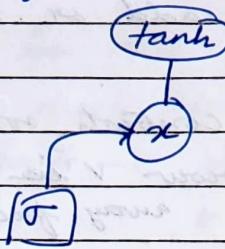
3. Memory update

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

~~C_t = C_{t-1} + i_t * \tilde{Q}_t~~

selection gate

Output gate - decides what to output from the memory



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

$$h_t = q_t * \tanh(c_t).$$

$$1. h_{t-1} = [1, 2, 3]$$

$$x_t = [4, 5, 6]$$

we have 3 units say and,

$$W_f = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & -1 \\ 5 & 6 & 7 & 8 & 9 & 10 \\ 3 & 4 & 5 & 6 & 7 & 8 \end{bmatrix}$$

$$h_{t-1}$$

$$\uparrow \rightarrow x_t$$

$$\text{so } W = 3 * (3+3) = 3 \times 6$$

matrix

$$b_f = [1, 2, 3].$$

(Consider all 1D vectors as column vectors)

$$[h_{t-1}, x_t] = [1, 2, 3, 4, 5, 6].$$

$$\textcircled{1} f_t = \sigma(W_f * [h_{t-1}, x_t] + b_f) \quad \text{forget gate}$$

$$= \sigma \left(\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & -1 \\ 5 & 6 & 7 & 8 & 9 & 10 \\ 3 & 4 & 5 & 6 & 7 & 8 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix} + \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \right)$$

$$W_f * [h_{t-1}, x_t] = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & -1 \\ 5 & 6 & 7 & 8 & 9 & 10 \\ 3 & 4 & 5 & 6 & 7 & 8 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix} = \begin{bmatrix} -6 \\ 175 \\ 133 \end{bmatrix}$$

$$W_f \cdot [h_{t-1}, x_t] + b_f = \begin{bmatrix} -6+1 \\ 175+2 \\ 133+3 \end{bmatrix} = \begin{bmatrix} -5 \\ 177 \\ 136 \end{bmatrix}$$

$$f_t = \sigma \left(\begin{bmatrix} -5 \\ 177 \\ 136 \end{bmatrix} \right)$$

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

$$\textcircled{4} \quad i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i). \quad // \text{input gate}$$

$$\textcircled{3} \quad \tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) -$$

$$\textcircled{4} \quad c_t = f_t * c_{t-1} + i_t * \tilde{c}_t \quad // \text{memory update}$$

elementwise multiplication

$$\left(\begin{bmatrix} a \\ b \end{bmatrix} * \begin{bmatrix} c \\ d \end{bmatrix} = \begin{bmatrix} ac \\ bd \end{bmatrix} \right) -$$

$$\textcircled{5} \quad o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad // \text{output gate}$$

$$\textcircled{6} \quad h_t = o_t * \tanh(c_t)$$

For next time step, pass o_t and h_t

For NN architecture, write the o/p layer
as

(1)

(2)

:

:

(25)

o/p
explicitly

16x16x3

3x3x3

5

→ # filters

16x16x3

3x3x3

consider

as filter

(implicitly 3rd dimension

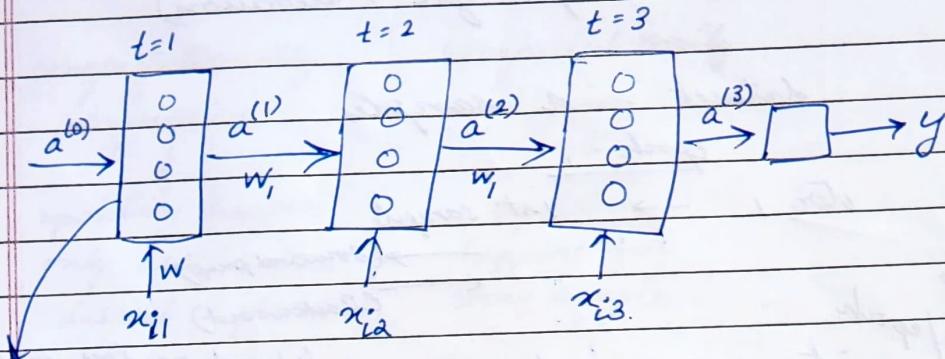
is 3 implicitly)

so answer is $n \times n \times 3$

07/10/22

Timeseries data input to RNN

$$x_i = \langle x_{i1}, x_{i2}, x_{i3} \rangle$$



summation
& activation

10/10/22

Optimization in DL

1. Gradient Descent

- ✓ calculate derivative (gradient)
- ✓ update weights using gradients during backpropagation.

1 epoch - forward + backward prop through entire training set.

✓ considers whole dataset during 1 iteration
(processes entire dataset as a whole)
dataset - n samples

epoch 1

iter 1 = n samples
 → (forward prop on all n samples)
 ← (Backward prop together
 & for entire dataset at once)
 (loss)

2. Stochastic gradient descent (SGD) :

GD but for every individual sample
(considers 1 sample for 1 iteration)

dataset - n samples.

epoch - 1

iter 1 → 1st sample

→ (forward pass)

← (backward)

n samples
⇒ n iter / epoch
iter 2 → 2nd sample <sup>↳ based on loss incurred
for 1st sample
alone</sup>

→

←

:

iter n → nth sample

→

←

processes every individual sample and
updates weight

3. Mini batch SGD :

divide dataset into batches and process
1 batch during every iteration

Dataset - 100 samples

batch size = 10

$$\Rightarrow \text{#iters} = \frac{100}{10} = 10$$

1 epoch → 10 iters

epoch - 1

iter 1 → batch 1

(1 to 10 samples)

iter 2 → batch 2

(11 to 20 samples)

GD

SGD

mini-batch

More RAM

less RAM

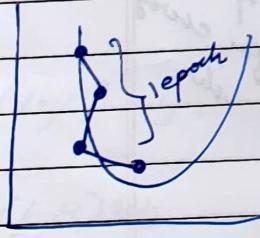
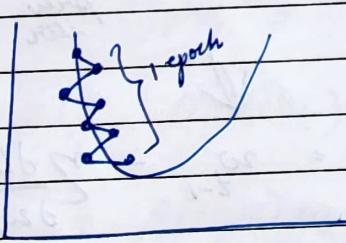
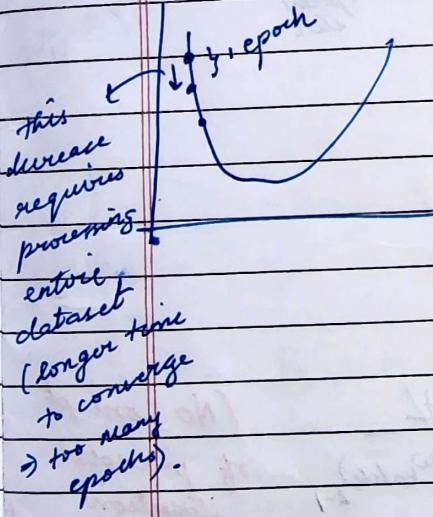
less RAM-
(> SGD, < GD ??)

computationally
expensive

computationally
expensive

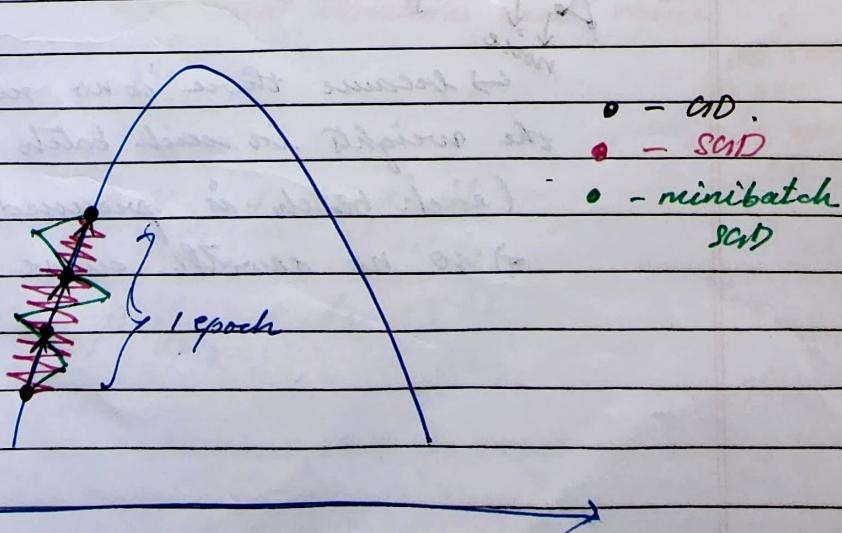
update happens
only after entire
dataset is processed

update
happens for
every sample



more zig-zag line
(100 samples
⇒ 100 updatios
so it won't be
uniform as
update depends
only on 1 sample).

100 samples
⇒ only 10
updatios



Drawback:

all the above methods, do not yield an optimum path to reach the peak
(the path followed is not smooth and keeps oscillating between some left and right extremes)



SGD update

$$(w_{\text{new}})_{\text{iter}} = (w_{\text{old}})_{\substack{\text{iter} \\ \text{prev}}} - \eta \frac{\partial L}{\partial (w_{\text{old}})_{\substack{\text{iter} \\ \text{prev}}}}$$

↓

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$

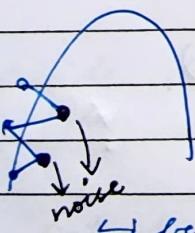
minibatch SGD
update

$$(w_{\text{batch}})_t = (w_{\text{batch}})_{t-1} - \eta \frac{\partial L}{\partial (w_{\text{batch}})_{t-1}}$$

(No concept
of previous
iteration
used in current
iteration)

$w_t \rightarrow$ current iter's
new value
 $w_{t-1} \rightarrow$ cur. iter's
old value

↳ because there is no relationship between
the weights in each batch
(each batch is processed independently)
⇒ so no smooth curve is observed



which is the
same as prev.
iter's old new
value?

f. SGD with momentum :

- ✓ Curves a smooth curve
- ✓ Considers previous weights
- ✓ Exponential weighted average is used

$$v_{t_1} = w_1$$

$$v_{t_2} = \beta v_{t_1} + (1-\beta)w_2 = \beta w_1 + (1-\beta)w_2$$

v_{t_1} \rightarrow velocity
 w_i \rightarrow acceleration importance given to previous weight (probability).

$$v_{t_3} = \beta v_{t_2} + (1-\beta)w_3.$$

$$= \beta(\beta v_{t_1} + (1-\beta)w_2) + (1-\beta)w_3.$$

$$= \beta^2 v_{t_1} + \beta(1-\beta)w_2 + (1-\beta)w_3.$$

$$= \beta^2 w_1 + \beta(1-\beta)w_2 + (1-\beta)w_3.$$

$v_{t_1}, v_{t_2}, v_{t_3} \rightarrow$ exponentially weighted averages.

SGD : $w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$

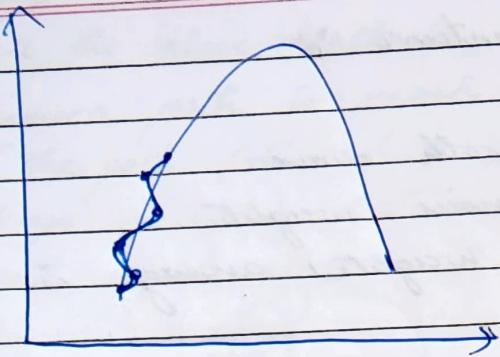
SGD with momentum : $w_t = w_{t-1} - \eta \frac{V_{dw_t}}{\partial w_t}$ derivative wrt current iterations old value
↳ current iterations new value ↓

$$\frac{V}{dw_1} = \frac{\partial V_{t_1}}{\partial w_t}$$

$$V_{dw_2} = \beta \frac{V_{dw_1}}{\partial w_t} + (1-\beta) \frac{\partial L}{\partial w_2}$$

$$V_{dw_3} = \beta \frac{V_{dw_2}}{\partial w_t} + (1-\beta) \frac{\partial L}{\partial w_3}$$

the same
as prev
iter's new
value??



(smooth curve).

curve depends on
β value11/10/22

5. Adagrad :

→ learning rate (α_t) is ~~so~~ dynamic
(Adaptive gradient descent).

→ gives better results than SGD with momentum

$$w_t = w_{t-1} - \gamma'_t \left(\frac{\partial L}{\partial w_{t-1}} \right)$$

$$\gamma'_t = \frac{\gamma}{\sqrt{\alpha_t + \epsilon}} \quad \text{where } \alpha_t = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_i} \right)^2$$

(cumulative sum)

$$\underbrace{\alpha_3}_{t=3} = \left(\frac{\partial L}{\partial w_1} \right)^2 + \left(\frac{\partial L}{\partial w_2} \right)^2 + \left(\frac{\partial L}{\partial w_3} \right)^2$$

γ → initial value we assign for the learning rate

Disadvantage :

May face vanishing gradient prob
($\gamma'_t \rightarrow$ might be very small if α_t is large)

→ very small updation, so model does not learn much).

→ provides better convergence than SGD with momentum

6. RMSProp :

→ Root mean square propagation

Change of w_t and use root mean square exponential averaging to compute it

$$w_t = w_{t-1} - \frac{\gamma}{\sqrt{V_t + \epsilon}} \left(\frac{\partial L}{\partial w_{t-1}} \right)$$

$$V_t = \beta V_{t-1} + (1-\beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

→ Improvement to Adagrad

→ Instead of computing cumulative sum, it considers exponential weighted average

7. AdaDelta :

→ Adaptive delta optimization

→ Completely removes η terms and added delta term

$$w_t = w_{t-1} - \frac{\sqrt{D_{t-1} + \epsilon} \times (\Delta w_t)^2}{\sqrt{V_t + \epsilon}}$$

where $\Delta w_t = w_t - w_{t-1}$

$$D_t = \beta D_{t-1} + (1-\beta)(\Delta w_t)^2$$

$$V_t = \beta_2 V_{t-1} + (1-\beta_2) \left(\frac{\partial L}{\partial w_t} \right)^2$$

→ improvement to Adagrad (just like RMSProp).

8. Adam

→ Adaptive moment estimation.

→ Combines momentum and RMSProp.

$$w_t = w_{t-1} - \frac{\gamma * v_{dw_t}}{\sqrt{s_{dw_t} + \epsilon}}$$

↑ momentum
RMSProp

$$v_{dw_t} = \beta_1 v_{dw_{t-1}} + (1-\beta_1) \frac{\partial L}{\partial w_t} \quad \text{Momentum}$$

$$s_{dw_t} = \beta_2 s_{dw_{t-1}} + (1-\beta_2) \left(\frac{\partial L}{\partial w_t} \right)^2 \quad \text{RMSProp}$$

After bias correction

$$\xrightarrow{\text{Adam formula}} w_t = w_{t-1} - \frac{\gamma * \frac{v_{dw_t}}{(1-\beta_1)^t}}{\sqrt{\frac{s_{dw_t}}{(1-\beta_2)^t} + \epsilon}}$$

→ This is done so that we consider the entire gradient $\frac{\partial L}{\partial w_t}$ (with raw weights).

without any exponential averaging

$$w_t = w_{t-1} - \frac{\gamma * v'_{dw_t}}{\sqrt{s'_{dw_t} + \epsilon}}$$

$$v'_{dw_t} = \frac{v_{dw_t}}{(1-\beta_1)^t}$$

$$s'_{dw_t} = \frac{s_{dw_t}}{(1-\beta_2)^t}$$

Epoch - 1

$$\textcircled{1} \quad w_t$$

$$w_{t \text{ new}} = w_{t-1}$$

$$\textcircled{2} \quad w_t$$

$$w_{t \text{ new}} = w_{t-1} - \eta, \nabla_{dw_t}$$

4

18/10/22

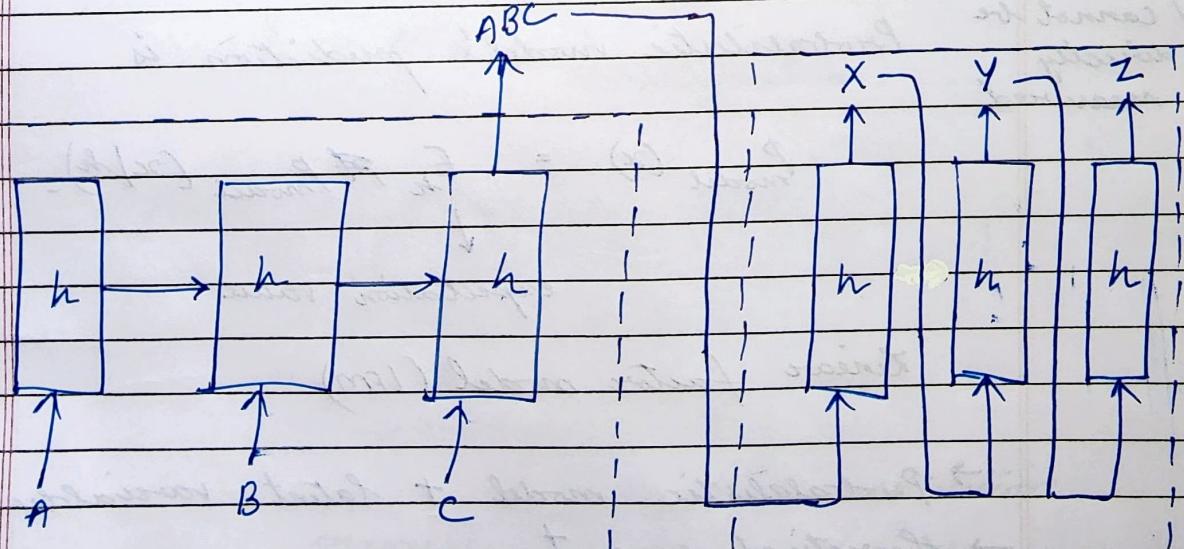
Encoder Decoder RNN

Text data

ip \rightarrow strings of lowercase, uppercase letters, digits, punctuation, whitespace, ...

Input - ABC
output - XYZ

Translation of a sentence



Encoder
(processes
the input)

composing

Decoder
(does the translation)

RNN - ? , 1 activation fn only
LSTM - has 4 gates , > 1 activation fns (sigmoid, tanh, ...)

Unit 3

Qns asked in miss feedback form:

- ✓ Ambience
✓ Arrangement
✓ Cleanliness
✓ Quality of food
✓ Quantity of food
✓ Variety in menu
- 6 features
- only 2 features
- } service
} Food
↓
can be grouped as above

→ Feature reduction

Techniques - PCA, LDA, ICA

Latent variables → new features obtained by say reducing or feature engineering input features.

✓ inferred from the observed variables

✓ cannot be directly measured Probabilistic model's prediction is

$$P_{\text{model}}(x) = \mathbb{E}_h p_{\text{model}}(x|h)$$

↓
expectation value

Linear factor model (LFM)

- Probabilistic model + latent variables
- theoretical concept
- linear decoder that generates x by adding some noise to a linear transformation of h

Data generation process

- sample h from a distribution

$$h \sim p(h)$$

$$p(h) = \prod_i p(h_i) \quad (\text{processing each letter of a word as in encoder-decoder})$$

$p(ABC) = p(A) \cdot p(B) \cdot p(C)$

to sample

$$x = Wh + b + \text{noise}$$

Factor analysis

$$h \sim N(h; 0, I)$$

\hookrightarrow identity matrix.

Read Deep Learning book

✓ Factor analysis

✓ PCA

\rightarrow probabilistic PCA

(to remove reconstruction

errors)

same as PCA as $\sigma \rightarrow 0$.

✓ ICA

\rightarrow prior $p(h)$ is fixed

\rightarrow deterministically generates

$$x = Wh$$

applications