

Course - 4

Week 1 : Computer Vision

Computer Vision Problems

→ Image classification

→ Object detection

→ Neural style Transfer (Artwork)

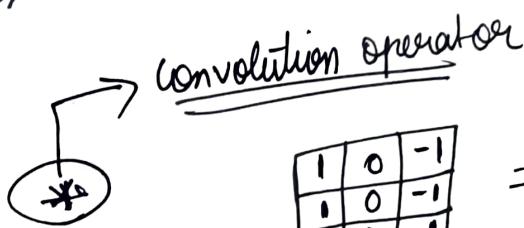
[Merge your face
with some filter]

Edge detection example

Task 1.: Given an image find its vertical and horizontal edges.

Vertical Edge detection

0	1	4
1	3	2
2	0	8
-1	0	7



$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

3×3

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

Suppose $6 \times 6 \times 1$
grayscale

• Called as filter
(or)
Kernel

4×4

Take a 3×3 in the

6×6 matrix

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	1	3	9

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

$$(3 \times 1) + (1 \times 2) + 0 \times 0 + 5 \times 0 + 7 \times 0 + (1 \times -1) + (8 \times -1) + (2 \times -1) = -5$$

-5	-4		

only do for other

for the next box shift the 3×3 by 1 to right
 $(0 \times 1) + (5 \times 1) + 7 \times 1 + (-1) \times (2) - 9 - 5$
 $+ 0 \times 1 + 0 \times 8 + 0 \times 2$

softmax.

keep sliding the 3×3 matrix over the 6×6 matrix

∴ final ans:

4×4

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	+3	-16

(4×4)

another image

} Vertical edge detector

python: $\otimes \rightarrow$ conv-forward

tensorflow \rightarrow tf.nn.conv2d

Keras \rightarrow conv2D

→ what asterisk means in different languages

Intuition of the filter

1	0	-1
0	-1	0
0	0	-1
0	-1	0

image
Dark on the left

Light on
the right

∴ Edge is deleted

More Edge Detection

(+ve edge detection)

light to
dark transition

(-ve edge detection)

dark to light transition

Smaller the pixel value,
darker the colour

e.g. $\frac{1}{0}$

Horizontal Edge filter

1	1	1
0	0	0
-1	-1	-1

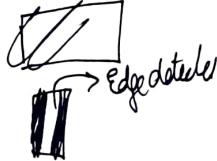
Vertical edge (light \rightarrow dark)

0	10	10	0	0	0
0	10	10	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0
1	0	0	0	0	0

6×6

$$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{bmatrix}$$

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0



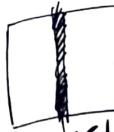
-ve edge (dark \rightarrow light)

0	0	0
0	0	0
1	0	-1
0	0	-1

6×6 .

$$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} = \begin{bmatrix} -30 & -30 \\ -30 & 30 & 0 \\ -30 & -30 \end{bmatrix}$$

-30	-30	
-30	30	0
-30	-30	
-30	-30	



Edge detected

Again, there are other filters too!!

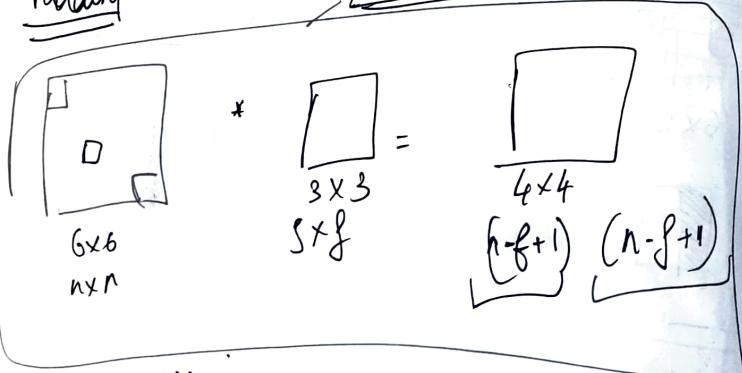
$$\text{Sobel filter} \quad \begin{matrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{matrix}$$

$$\text{Scharr filter} \quad \begin{matrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{matrix}$$

* Again you can just use DL to learn the filter matrix. GG!!

✓ the backprop will learn the 9 numbers in the filter !!

Padding Normal



Downside of this is

→ You are shrinking the image ie $(6 \times 6 \rightarrow 4 \times 4)$
The edge pixel $(1,1)$ and $(6,6)$ are used only once, but the middle ones say $(3,3)$ are used twice or thrice atleast!,

it's as if you are discarding the edges of the image
corners

\therefore you pad the image

$$\begin{matrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{matrix} \quad \begin{matrix} 6 \times 6 \\ 6 \times 6 \\ 6 \times 6 \\ 6 \times 6 \\ 6 \times 6 \end{matrix} \quad \begin{matrix} 8 \times 8 \\ 8 \times 8 \\ 8 \times 8 \\ 8 \times 8 \\ 8 \times 8 \end{matrix}$$

$$p = \text{padding} = 1 \text{ (here)}$$

$$\begin{aligned} & (n+2p-f+1) \\ & = 8 + 2(1) - 3 + 1 \\ & \quad \boxed{6 \times 6} \end{aligned}$$

How much to pad?

Valid and Same convolutions

Valid: NO PADDING: $n \times n \times f \times f \rightarrow (n-f+1) \times (n-f+1)$

Same convolution: Pad so that output size is same as input size.
(called as zero padding)

$(n+2p-f+1) = n$

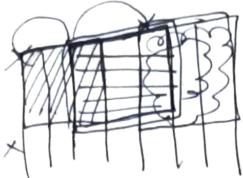
$$\Rightarrow \boxed{p = \frac{f-1}{2}}$$

f is usually odd

3×3
 5×5
 7×7

Strided convolution
You move the box while convolution with stride

$$\begin{matrix} & \times & 3 \times 3 & = & 4 \times 4 \\ 7 \times 7 & & & & \end{matrix}$$



If stride=2

~~2x2~~ move by 2 rows, 2 columns
while going down goes sideways.

$$n \times n \quad * \quad f \times f$$

padding p strides

$$\frac{n+2p-f}{s} + 1 \quad \left(\frac{n+2p-f}{s} \right)$$

if it isn't a
proper number,
take FLOOR

Actual convolution in math textbook

$$6 \times 6 \quad * \quad \begin{matrix} 3 & 4 & 5 \\ 1 & 0 & 2 \\ 9 & -1 & 7 \end{matrix}$$

filter 3x3

You just flip it ~~turn it left~~,
horizontally and vertically, before the
element wise multiplication

$$\begin{matrix} 7 & 9 & -1 \\ 2 & 0 & 1 \\ 5 & 4 & 3 \end{matrix}$$

- what we have doing is called cross-correlation rather than convolution

Doesn't matter

This is what is done in DCL

Convolution over Volumes (RGB images)

$6 \times 6 \times 3$ $3 \times 3 \times 3$ 4×4

height width no of channels

width height # channels

have to be equal

$$3 \times 3 \times 3$$

9 on R
9 on G
9 on B

get ~~*~~ by

$3 \times 3 \times 3$ filter, so on!

multiple filters (Both horizontal, vertical edge detection)

$6 \times 6 \times 6$ $* \quad \begin{matrix} 3 \times 3 \times 3 \\ \text{vertical detect} \end{matrix} \quad = \quad 4 \times 4 \quad \rightarrow \quad 4 \times 4 \times 2$

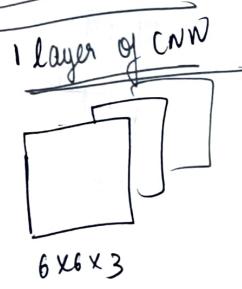
$* \quad \begin{matrix} 3 \times 3 \times 3 \\ \text{horizontal detect} \end{matrix} \quad = \quad 4 \times 4 \quad \rightarrow \quad 4 \times 4 \times 2$

Summary

$$n \times n \times \frac{n_c}{\text{no of channels}} * f \times f \times n_c = \text{no of channels}$$

$$= (n-f+1) \times (n-f+1) \times n_c$$

Each filter (bias)



$$a^{[0]}$$

$$\begin{aligned} z^{[1]} &= w^{[1]} a^{[0]} + b^{[1]} \\ a^{[1]} &= g(z^{[1]}) \end{aligned}$$

$$\boxed{\begin{aligned} z^{[1]} &= w^{[1]} a^{[0]} + b^{[1]} \\ a^{[1]} &= g(z^{[1]}) \end{aligned}}$$

1. You have 10 filters of $3 \times 3 \times 3$, how many parameters should it learn?

$$\begin{aligned} & \text{Diagram of a 3x3x3 filter} \\ & 3 \times 3 \times 3 + b \end{aligned}$$

$$= \underline{28 \text{ parameters per filter}}$$

$\therefore 280 \text{ parameters}$
for 10 filters

↳ response

* You have to learn 280 parameters irrespective of input image size! (\therefore No overlapping)

Notations

$f^{[c]}$ = filter size.

$p^{[c]}$ = padding

$s^{[c]}$ = stride

$n_c^{[c]}$ = no of filters

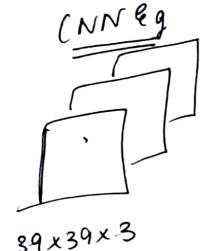
$$\text{filter} = f^{[c]} \times f^{[c]} \times n_c^{[c]}$$

$$\text{Activation: } a^{[c]} \rightarrow n_h^{[c]} \times n_w^{[c]} \times n_c^{[c]}$$

$$A^{[c]} \rightarrow m \times n_h^{[c]} \times n_w^{[c]} \times n_c^{[c]}$$

$$\text{Weight} = f^{[c]} \times f^{[c]} \times n_c^{[c]} \times \frac{n_c^{[c]}}{\text{no of filters}}$$

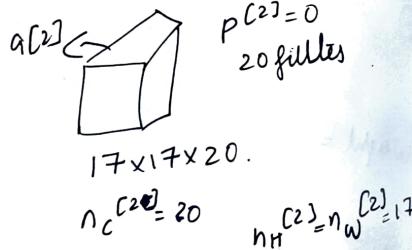
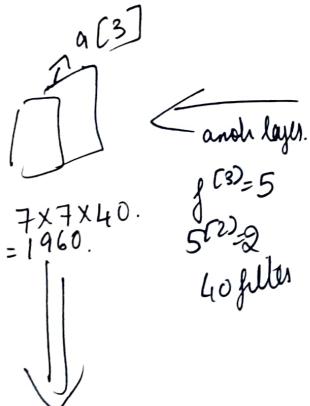
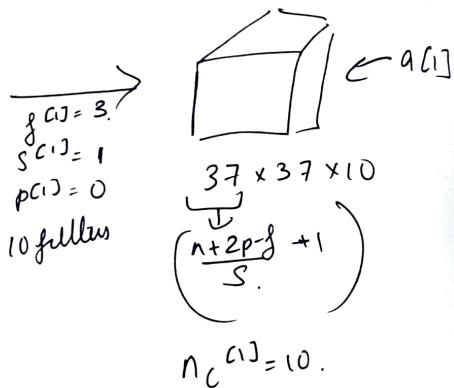
bias : $\star n_c^{[c]}$ dimension



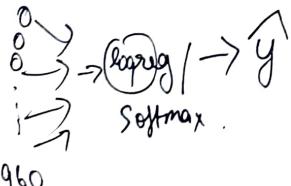
$39 \times 39 \times 3$

$$n_H^{[0]} = n_W^{[0]} = 39.$$

$$n_C^{[0]} = 3$$



Unroll it

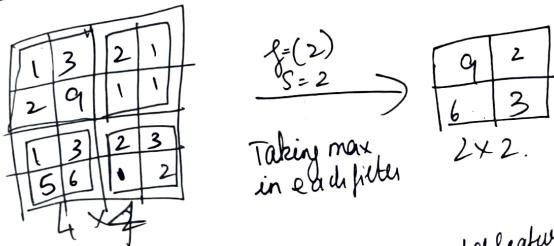


Types of layer in CNN

- Convolution (Conv) ← we have only seen this
- Pooling (Pool) ←
- Fully connected (FC) ←

Pooling layers

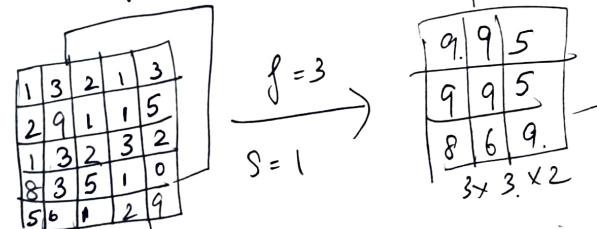
: Max pooling (Taking the max number from each filter)



Intuition : Think of 4×4 as some set of features in a NW.
If there is a large number, then it probably being it detected something!

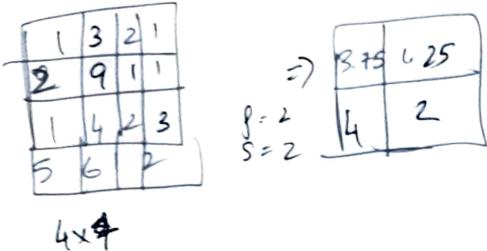
\Rightarrow It has no parameters to learn!

Another eg:



If there are multiple channels, you go channel by channel down one layer, then go to behind layer and so on

Average pooling (as the name suggests)
takes avg



max pooling is used much more frequently!

Hyperparameters

mostly $f=2$ ~~$B=2$~~ is used.

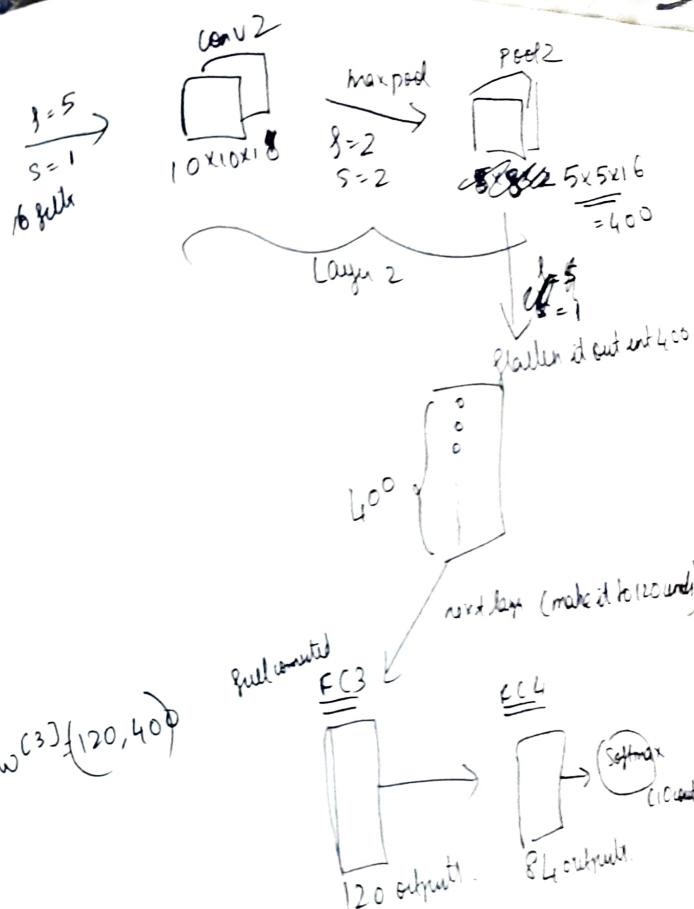
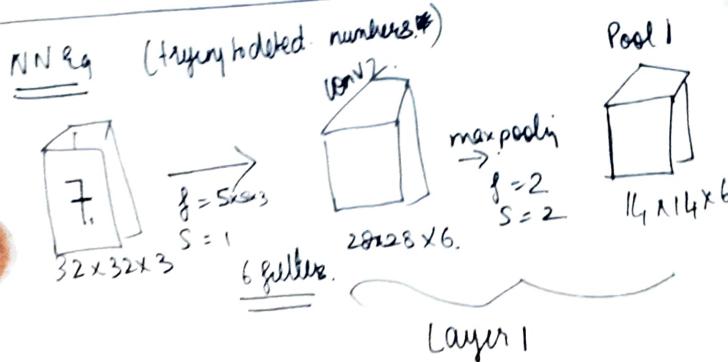
[Mostly no padding]

$$f=3 \quad s=2$$

→ max or avg pooling.

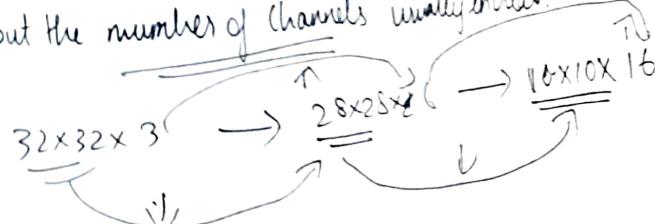
Quick handy trick
 $f=2 \quad s=2$ redundant
by a factor of 2
ie $28 \times 28 \xrightarrow{f=2} 14 \times 14$

Pooling: NO PARAMETERS TO learn!
it's a fixed function



As you go deeper into the NN, the (n_h, n_w) will decrease
height width

but the number of channels usually increase.



Note:

- The activation go down as you progress deeper in the NN.

$$32 \times 32 \times 3 = 3K$$

$$5 \times 5 \times 16 = 400$$

decreas

$$84 \times 1 = 84$$

$$10 \times 1 \rightarrow 10$$

• η shouldn't drop too quickly

- As for the hyperparameters (choose from some already worked literature)

Why convolution?

- There will too many parameters if you are fully connected network

• There is parameter sharing in convolution

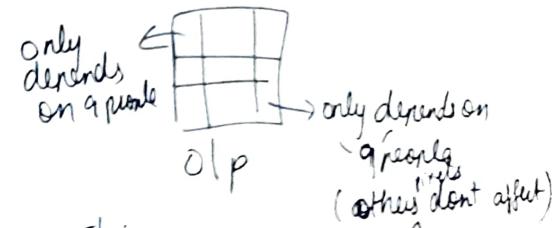
↳ a feature detector (such as a vertical edge detector) that's useful in one part of image is pretty useful in another part of image.

You can use the detector over and over.

Suppose fully connected then
 $100 \times 100 \times 3$ $50 \times 50 \times 3$

You need to learn ($50 \times 50 \times 3$, 100×3)
parameters.

- Sparsity of connections: In each layer, each o/p value depends only on a small no of inputs.

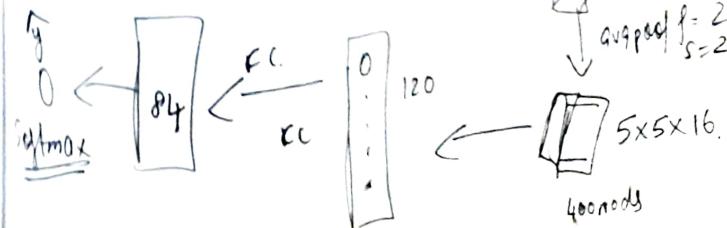
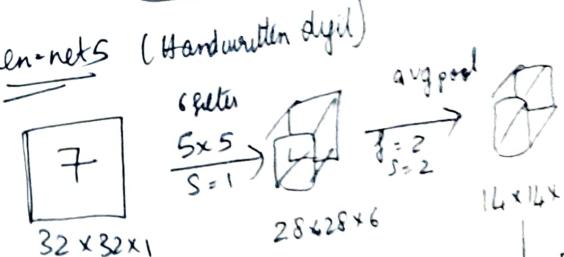


→ less prone to overfitting

- Translation variance (Even if a cat is shifted a few pixels to the left or right, it should still be a cat)
(more robust)

Week 2 : Classical Networks

Lennets (Handwritten digit)



Sequence Models

Course 5
Week 1

q: Speech recognition, music generation, Sentimental analysis, DNA sequence, Translation, video activity recog, Name entity recog.
Supervised learning

Notation (Name entity recog)

x : H.P and Herman graps anneta newsfull
 $x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(t)}$ $0 \ 0 \ 0 \ 0$
 y : 11 0 1 1 $y^{(1)} \ y^{(2)} \ \dots$ $y^{(t)}$

T =Temporal sequence

$$\bar{T}_x = q \quad \bar{T}_y = q.$$

\rightarrow tth element in ith ex:

$x^{(i) \leq t}$, $y^{(i) \leq t}$

$T_x^{(i)} \ y$ length of ith ex.

$T_y^{(i)} \ y$ l. of ith ex in OLP.

Vocabulary (Dictionary)

a	1
action	2.
Harry	4075
Patti	6830
zulu	10,000

one hot on x:
encoding

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{10,000} \ 4075$$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{10,000} \ 6830$$

$$\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{10,000} \ a$$

$x^{(t)}$ is a one hot vector.

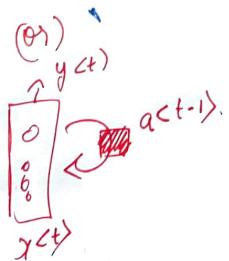
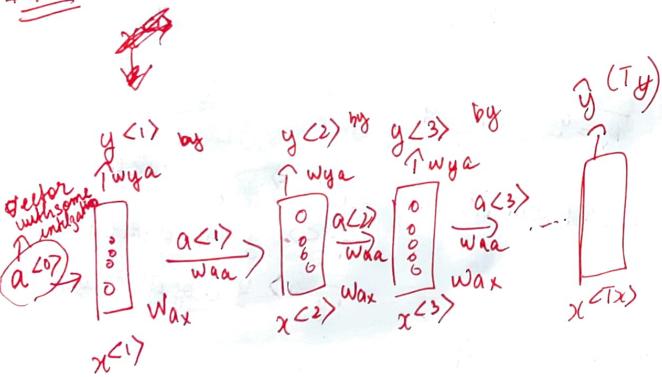
Why not standard neural.



Problems: I/P O/P lengths in diff. e.g.

Doesn't share features learned across different positions of text.

RNN



Waa, Wax, Way are the same throughout

all Paths:
You get learning ^{rate} α Input from previous $x^{(t-1)}$ (prev words)
not from the next forward set of words

$$a^{(0)} = 0$$

$$a^{(t)} = g(w_{aa}a^{(t-1)} + w_{ax}x^{(t)} + b_a)$$

$$y^{(t)} = g(w_{ya}a^{(t)} + b_y)$$

Some other activation function

mostly it with \tanh or relu
Mainly

Mostly sigmoid or softmax.

Forward prop.

$$a^{(t)} = g(w_{aa}a^{(t-1)} + w_{ax}x^{(t)} + b_a)$$

$$\hat{y}^{(t)} = g(w_{ya}a^{(t)} + b_y)$$

Simplified version

$$a^{(t)} = g(w_a[a^{(t-1)}, x^{(t)}] + b_a) \quad \text{eg}$$

$$w_a = \begin{bmatrix} w_{aa} & w_{ax} \end{bmatrix}$$

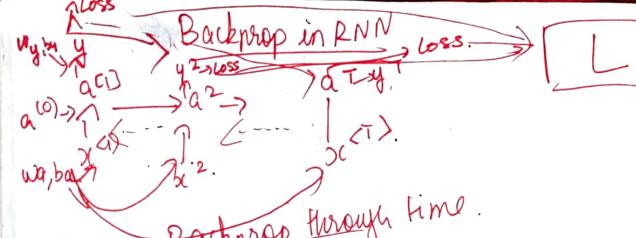
$$\begin{bmatrix} a^{(t-1)}, x^{(t)} \end{bmatrix} \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}^T \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}^T \cdot w_a = (100, 10, 000) \quad [a^{(t-1)}, x^{(t)}] = 10100$$

$$\begin{bmatrix} w_{aa} & w_{ax} \end{bmatrix} \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = w_{aa}a^{(t-1)} + w_{ax}x^{(t)}$$

Basically the same.

$$\hat{y}^{(t)} = \dots$$

3



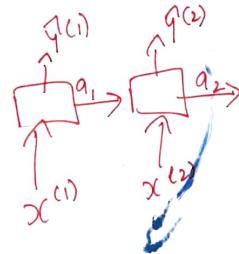
$$L^{(t)} = (y^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

$$L = \sum_{t=1}^{T_y} L^{(t)} (y^{(t)}, y^{(t)})$$

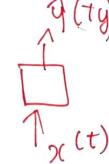
Sometimes in RNN length of I/p and O/p sequence aren't equal
e.g.: Machine Translation

Different types of RNN

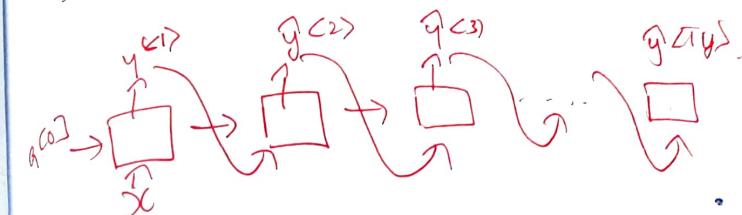
Until now $T_x = T_y$.



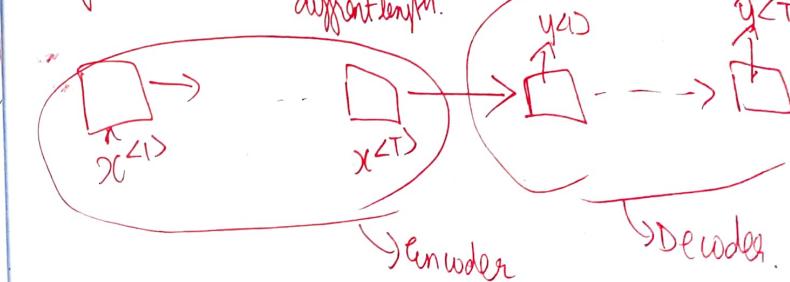
e.g.: Name entity recognition



Many to Many architecture
(Many I./O.ps and many O./I.ps)

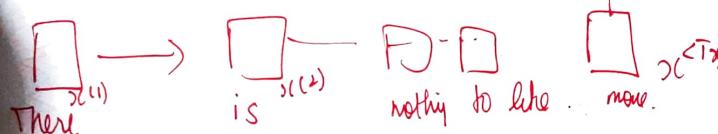


Many to Many (without $T_x = T_y$) (e.g.: Machine Translation)
different length.



In exploding gradient problem
You get Nan as O/P

Many to one architecture



Language Model and Sequence Generation

The example and pain Salad - y
 The " and near salad
 The language modelling will help us choose this

It will do a probability of either of these two sentences

$$P(\dots \text{ pain salad}) = 3.2 \times 10^{-3}$$

$$P(\dots \text{ near } \dots) = 5.7 \times 10^{-10} \text{ more probable.}$$

Language Model
 It calculates $P(\text{sentence}) = ?$

$$P(y^{(1)}, y^{(2)}, \dots)$$

↓ ↓
Sentence 1 Sentence 2

Training set: large CORPUS of English text body

Cats average 15 hours of sleep a day. $\langle \text{EOS} \rangle$ end of sentence extra token

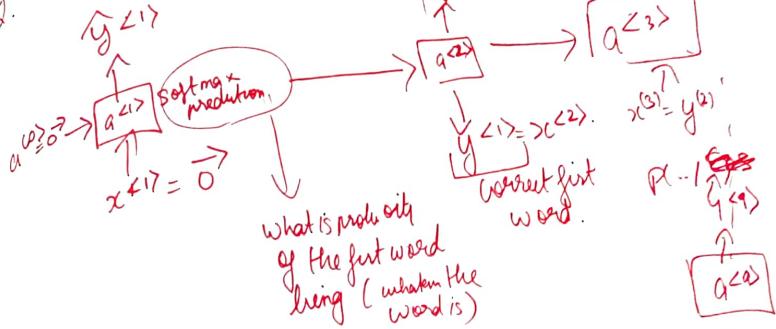
1. Tokenize. $y^{(1)} y^{(2)} y^{(3)} y^{(4)} \dots y^{(n)}$
 (one hot encoding)

Suppose: Some word is not there in your dictionary (large corpus)
 replace it with a $\langle \text{UNK} \rangle$ Token
 unknown.

$$x^{(t)} = y^{(t)}$$

Build RNN model.

2.



Each step will look at preceding words and calculate the probability of the next word occurring

3. Train RNN

$$L(\hat{y}^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

Softmax Loss func.

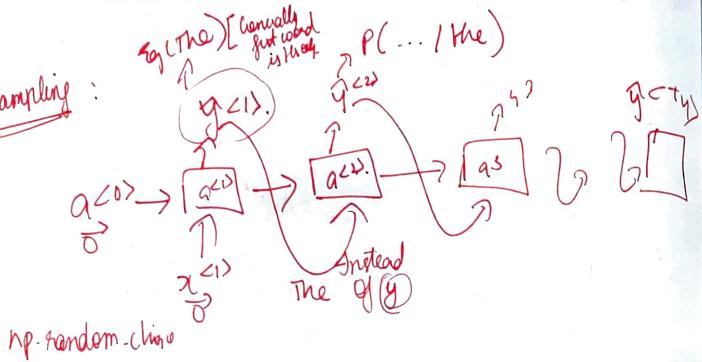
$$L = \sum L(q_t, y_t)$$

4. Tells us:

$$P(y^{(1)}, y^{(2)}, y^{(3)}) = \underbrace{P(y^{(1)})}_{\text{Prob of 1st word}} \times \underbrace{P(y^{(2)} | y^{(1)})}_{\text{Prob of 2nd word given 1st word}} \times \underbrace{P(y^{(3)} | y^{(2)}, y^{(1)})}_{\text{Prob of 3rd word given 1st and 2nd word}}$$

Sampling Novel Sequence (Random Sentence)

Sampling:

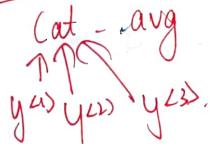


You can keep sampling till $\langle \text{EOS} \rangle$

(or)
Just stop when you want 20 or 30 words.

Here we're doing word modelling to build a sentence

You can also do character modelling to build a word



Adv of Char model: No need of $\langle \text{UNK} \rangle$

Disadv: → Not that effective
→ computational expensive

| You can use this
to generate a
Shakespeare poem!

Variants of recurrent unit & RNN

The cat was full
The cats were full.
long statement

• It's difficult for the RNN both during forward and back prop. to remember whether it saw

cat or cats!

Basic RNN, you can't remember so much!

• Highly influenced by local words
ie $y^{(2)}$ is highly influenced by $y^{(1)}$

$y^{(2)}$ is hardly influenced by $y^{(170)}$

• Can't capture Long-Range dependencies

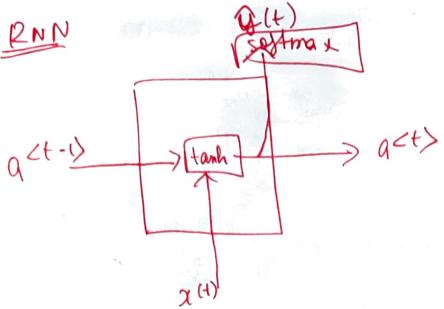
• Exploding gradients → the gradient clipping
 ↳ see if it is bigger than some threshold, if so rescale gradient vectors.

• Vanishing gradients → the GRU
 Gated Recurrent Unit

Gated Recurrent Unit



RNN

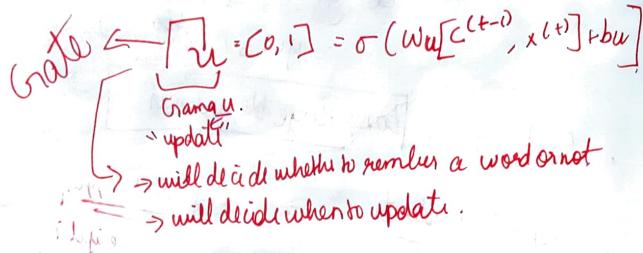


GRU (simplified)

C = memory cell.

$$c^{(t)} = a^{(t)}$$

$$\tilde{c}^{(t)} = \tanh(w_c [c^{(t-1)}, x^{(t)}] + b)$$



The cat was full.

$c^{(t)} = 1$
1 is singular
0 is plural

E! OH! it is singular; probably it is "WAS".

After this, you can forget the value.

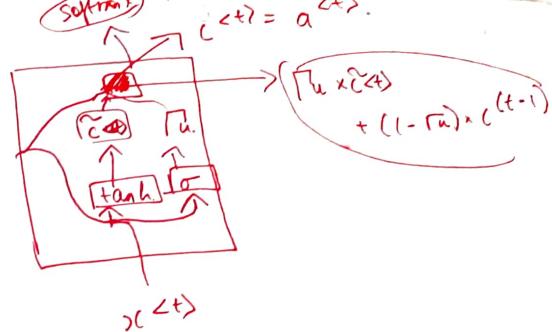
$$c^{(t)} = \underbrace{\Gamma_u \tilde{c}^{(t)}}_{\text{don't forget it}} + (1 - \Gamma_u) * c^{(t-1)}$$

The cat was full.
 $\Gamma_u = 1$ $\Gamma_u = 0$ $\Gamma_u = 0$
when $\Gamma_u = 0$.

Just ~~not~~ the previous value of $c^{(t)}$.

GRU simplified

$$c^{(t-1)} = a^{(t-1)}$$



How does it solve vanishing gradient?

$$\Gamma_u = \sigma(w_u [c^{(t-1)}, x^{(t)}] + b_u)$$

If it is a large negative value since it's sigmoid

$$\Gamma_u = 0$$

Then in the formula

$$c^{(t)} = \Gamma_u \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

$c^{(t)} = c^{(t-1)}$: value of $c^{(t)}$ is maintained

- Solves vanishing gradient problem.
- Can we it over long range dependency.

full GRU

BTW $\tilde{C}^{(t)}$ is called Candidate Memory Value

$$\tilde{C}^{(t)} = \tanh(W_c [\tilde{r}_c * C^{(t-1)}, x^{(t)}] + b_c)$$

relevance
gate, how relevant is $C^{(t-1)}$ for calculating
the next $C^{(t)}$

$$\Gamma_u = \sigma(W_u [C^{(t-1)}, x^{(t)}] + b_u)$$

$$\Gamma_g = \sigma(W_g [C^{(t-1)}, x^{(t)}] + b_g)$$

$$C^{(t)} = \Gamma_u * \tilde{C}^{(t)} + (1 - \Gamma_u) * C^{(t-1)}$$

Mutiple ways of deriving (with Γ_g , without Γ_g , ...)
But researchers said this is more robust

LSTM (Longshort term memory)

forget

~~$$\tilde{C}^{(t)} = \tanh(W_c [\tilde{r}_c * C^{(t-1)}, x^{(t)}] + b_c)$$~~

~~$$F_u = \sigma(W_u [C^{(t-1)}, x^{(t)}] + b_u)$$~~

LSTM Equations

$$\tilde{C}^{(t)} = \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c)$$

update gate $\Gamma_u = \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u)$.

forget gate $\Gamma_f = \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f)$

$$O = \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

$$C^{(t)} = \Gamma_u * \tilde{C}^{(t)} + \Gamma_f * C^{(t-1)}$$

$$a^{(t)} = \Gamma_0 * \text{tanh}(C^{(t)})$$

$a^{(t)} = \text{softmax}(a^{(t)})$

3 gates!!

LSTM and GRU are good at generating values.

peephole connection
just one equation right!

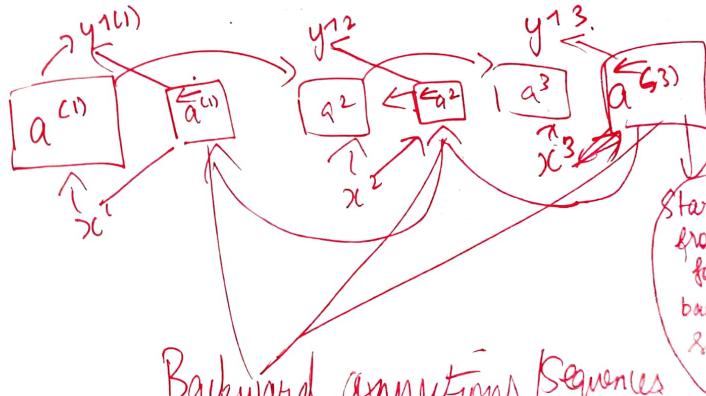
$$\Gamma_0 = \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

Adv of GRU: Simpler, easy to compute big network

But LSTM is much better!

BiDirectional RNN

Motivating eg: Teddy Roosevelt You need future information here!
In Name Entity



Batchwise annotations / Sequences