

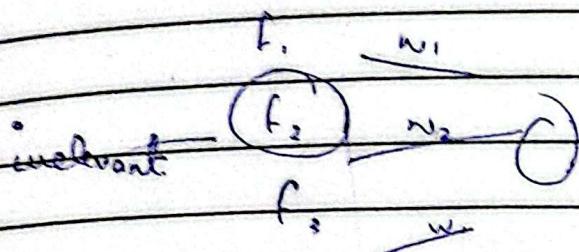


M T W T F S

Date _____

Feature Selection:

$f_1, f_2, f_3, \dots, \text{label}$



$$\chi^2 = \frac{\text{observed values}}{\text{expected value}} = \sum (O_i - E_i)^2$$

		E_i		values	
		Buy	Not buy	O_i	
Red	-10	20	30		$\frac{30 \times 25}{55} = \frac{30}{55}$
Blue	15	10	25		$\frac{25 \times 25}{55} = \frac{25}{55}$
	25	30	55		$\frac{11.36}{55} = \frac{11.36}{55}$
					$13.63 = 13.63$

$$\chi^2 = \frac{(10 - \frac{30 \times 25}{55})^2}{\frac{30 \times 25}{55}} + \frac{(20 - 16.36)^2}{16.36} + \frac{(15 - 11.36)^2}{11.36} + \frac{(10 - 13.63)^2}{13.63}$$

tagent
30 and 25

$$\chi^2 = \left[\frac{(10 - 13.63)^2}{13.63}, \frac{(20 - 16.36)^2}{16.36}, \frac{(15 - 11.36)^2}{11.36}, \frac{(10 - 13.63)^2}{13.63} \right]$$

$$= [0.97, 0.80, 1.16, 0.97]$$

largest difference

Partial distance
One relevant

selection criteria

~~Top k feature~~

Accept or reject a feature

f_1, f_2, f_3, f_4 Partial
 color decision columns

buy, not buy	... (red buy, red don't buy)
red)
i.	

Degree of freedom

[0.97, 0.80, 1.16, 0.97] ≈ 3.9

0 ... 9 - 1st deg
 10 ... 99 - 2nd

100 ... 99 - 3rd

Date _____

(M T W T F S)



day 1 2 3

20. 12...

10. 14...

5. 16...

1. 21...

80%. $3.9 < 12$ True

confident

I accept these facts with

80% confidence

 $3.9 < 21$

I accept with 99%

confidence

 $[0.97, 0.80, 20, 0.97] \text{ s } 23 < 21$

False

I reject these

facts w/

Date _____

M T W T F S



Individual Feature check

$$\begin{array}{r} \text{z-test} \\ \hline = 0.97 - 6 \\ \hline 0 \end{array}$$

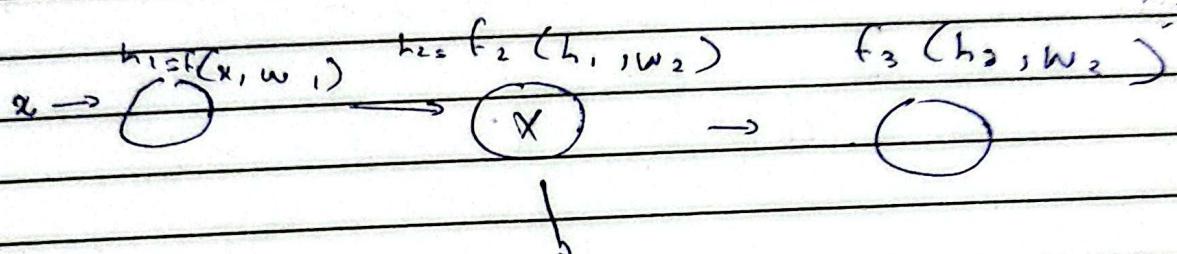
men

$$\begin{array}{r} x-4 \\ \hline 6 \end{array}$$



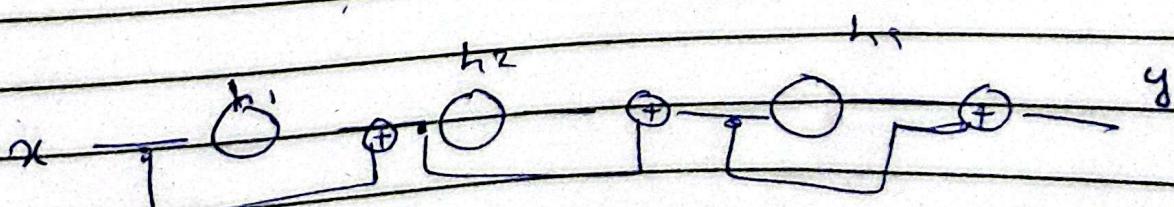
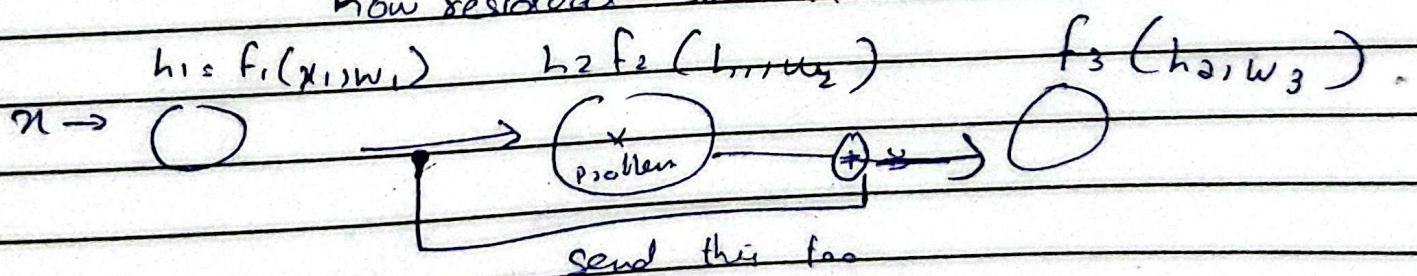
Residual Neural Network

CNN — Conv2D → MaxPooling → depth ↗
 ↓ vanishing gradient issue



$$\begin{aligned}
 y &= f_3(h_3, w_3) && \text{if problem is here it goes forward} \\
 &= f_3(f_2(h_1, w_2), w_3) \\
 &= f_3(f_2(f_1(x_1, w_1), w_2), w_3)
 \end{aligned}
 \quad \left. \begin{array}{l} \text{nested} \\ \text{behavior} \end{array} \right\}$$

how residual solve it

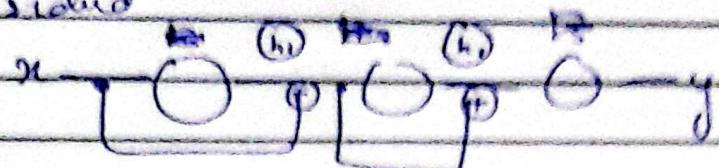


Here \rightarrow $x - \textcircled{O} \rightarrow \textcircled{O} - \textcircled{O} - y$
 $h_1 = f_1(x, w_1)$ $h_2 = f_2(h_1, w_2)$
 Date _____



(M/T/W/T/F/S)

In general



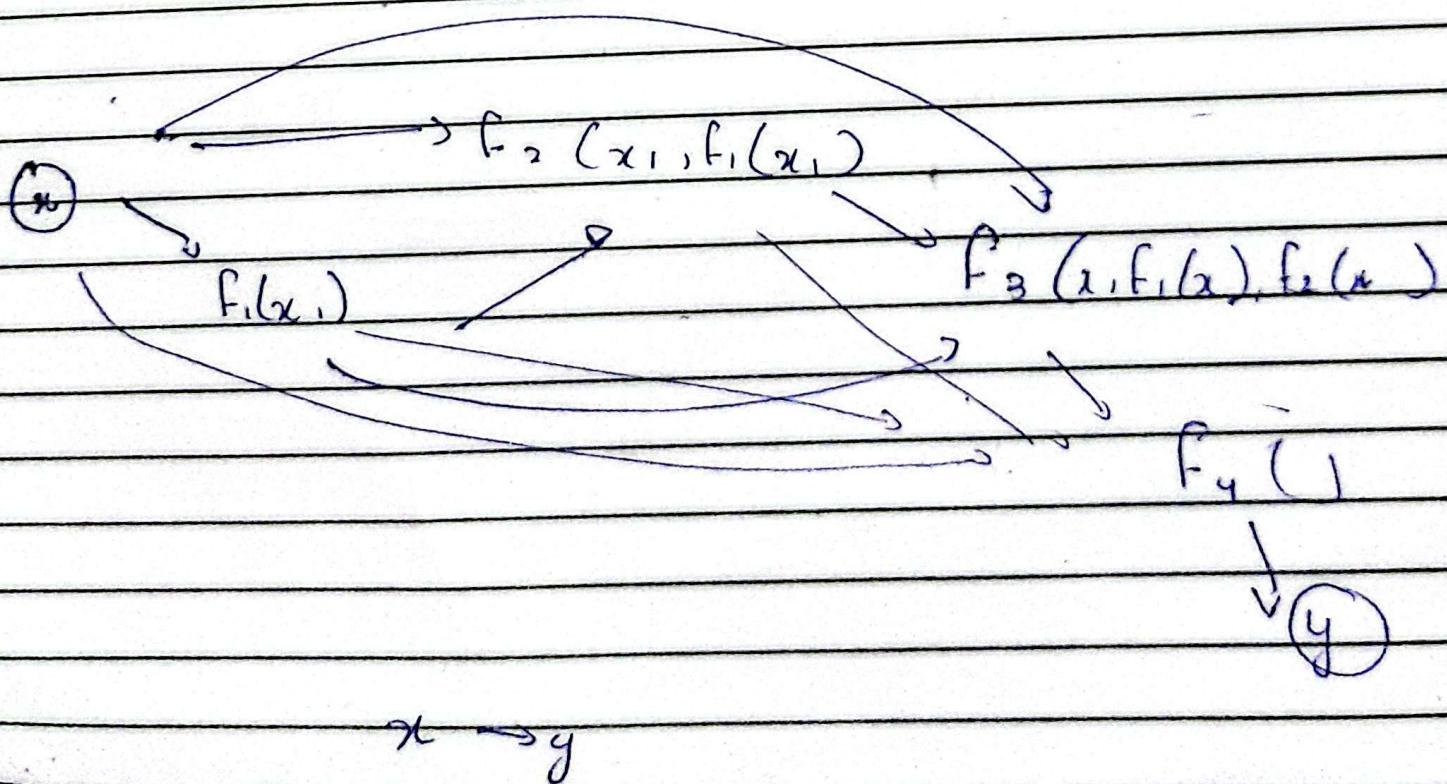
$$h_1 = f_1(x, w_1) + x$$

$$h_2 = f_2(h_1, w_2) + h_1$$

$$\cancel{h_3} y = f_3(h_2, w_3) + h_2$$

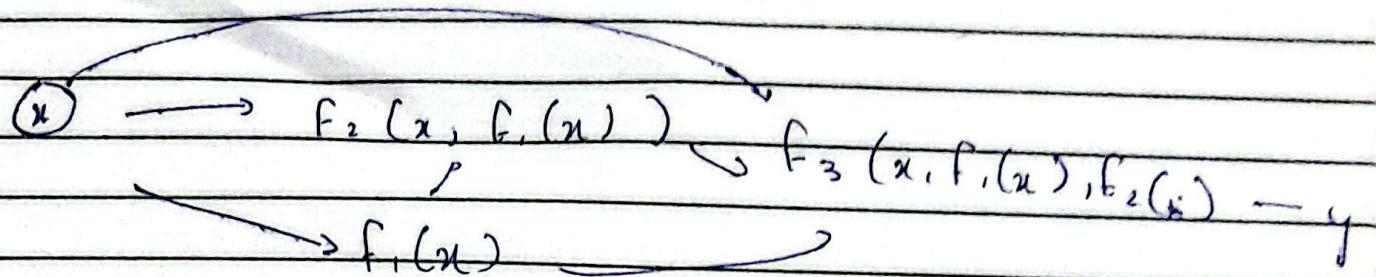
$$y = x + f_1(x, w_1) + f_2(x + f_1(x, w_1)) +$$

$$f_3(x + f_1(x, w_1) + f_2(x + f_1(x, w_1)))$$



Date _____

(M T W T F S)



Multiple parts

$$x \rightarrow y$$

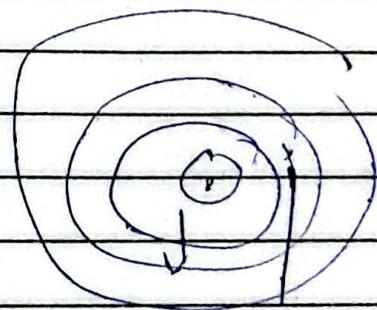
$$x \rightarrow f_2 \rightarrow f_3 \rightarrow y$$

$$x \rightarrow f_1 \rightarrow f_3 \rightarrow y$$

$$x \rightarrow f_1 \rightarrow f_2 \rightarrow f_3 \rightarrow y$$

Nested neural networks

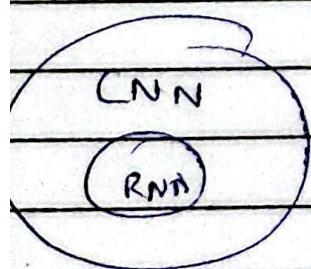
layers
in
neural
network



If some problem
path is going
in upward layers

You can
bypass
problematic
part

good part ;
also going



Res + ReLU

$$\text{ReLU} = \max(0, h_i)$$

solutions to
nonisotropic gradient

not a ~~form~~ perfect solution

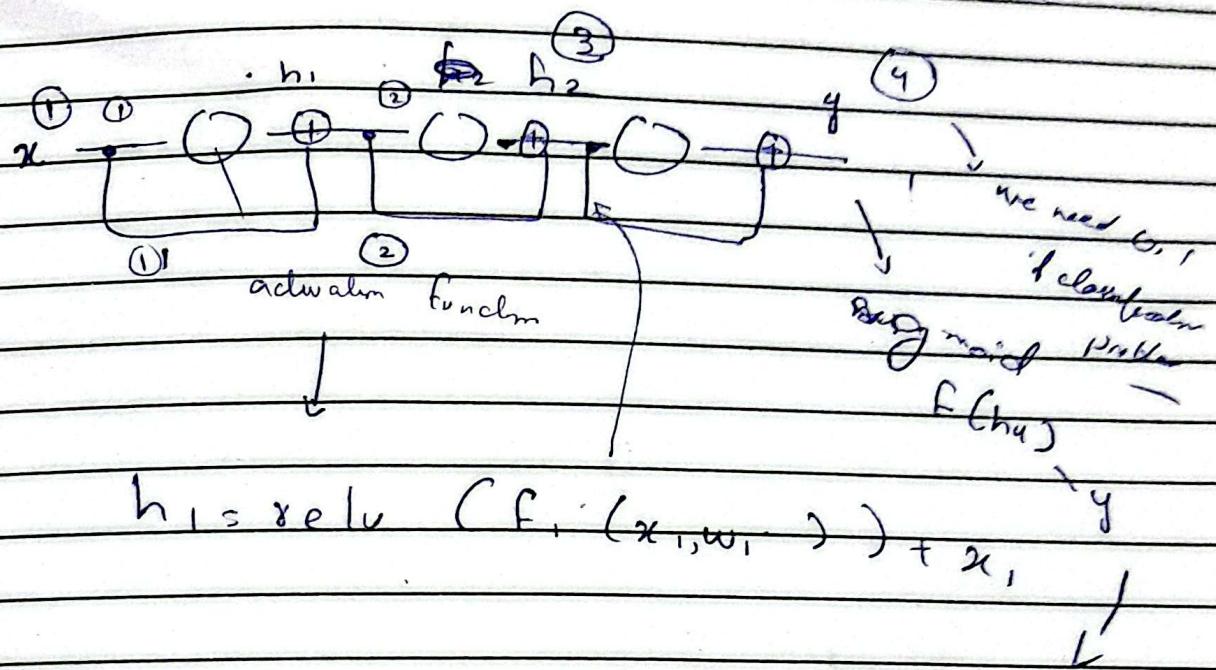
Whole solution

Res + Relv

+ Batch Normalization

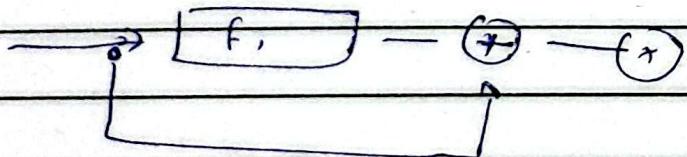
Date _____

MTWTFSS



We want to do

normalization on every step



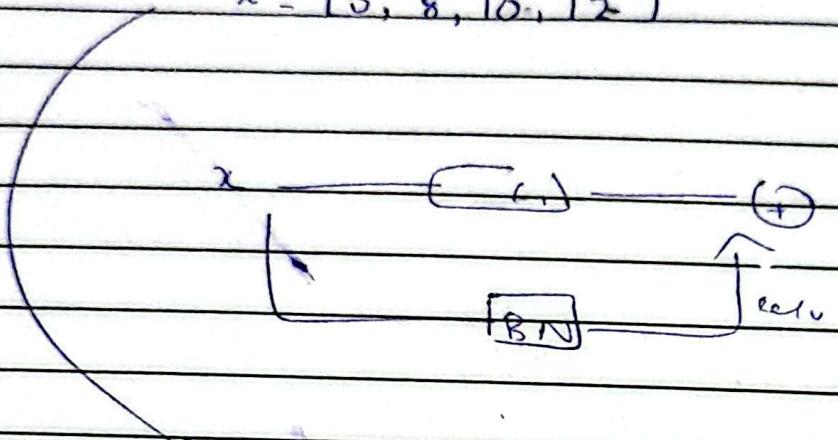
$$h_1 = (F_1(x) + x) \times \frac{1}{\sqrt{2}}$$

normalized

B

Batch Normalization

$$x = [5, 8, 10, 12]$$



$$\mu = 8.75$$

$$\sigma^2 = \frac{1}{4} [(5 - 8.75)^2 + (8 - 8.75)^2 + (10 - 8.75)^2 + (12 - 8.75)^2]$$

$$\sigma^2 = 6.68$$

$$\sigma = \sqrt{6.68}$$

$$\sigma = \frac{1}{\sqrt{N-1}}$$

$$x' = \left[\frac{5-\mu}{\sigma}, \dots \right]$$

~~zero problem~~ $\rightarrow \sigma = \sqrt{\sigma^2 + \text{epsilon}}$

$$x' = \left[\frac{5-8.75}{2.58}, \frac{8-8.75}{2.58}, \frac{10-8.75}{2.58}, \frac{12-8.75}{2.58} \right] = [2.58, 2.58, 2.58, 2.58]$$

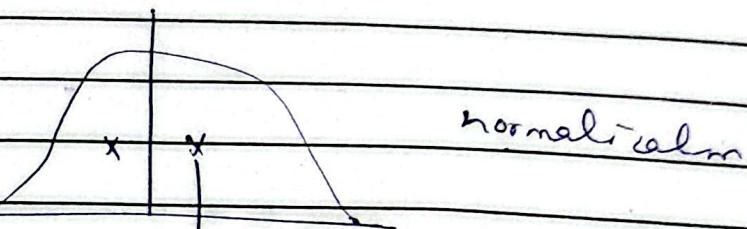
Date _____

(M T W T F S)



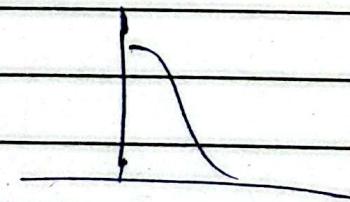
$$[- - + +]$$

$$[-1.45 \quad -0.29 \quad +0.48 \quad +1.25])$$

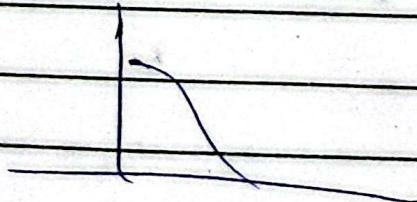


Symmetry

so we don't need other half



even with red,



positive
values

Date

(M|T|W|T|F|S)



$$x' = \frac{x - u}{\delta}$$

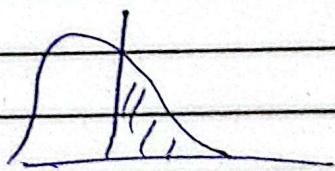
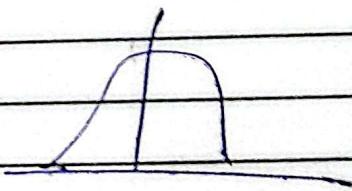
$$x'' = rx' + s$$

Batch normalization

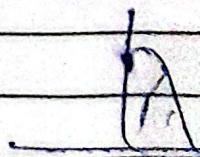
$$\text{let } r = 1$$

$$s = 0$$

$$[-1.45 \quad -0.29 \quad +0.48 \quad 1.25]$$



scale



a small portion



We update the value of λ and β

$$\frac{\partial E}{\partial \lambda}$$

$$\frac{\partial E}{\partial \beta}$$

so that

Vanishing Gradient

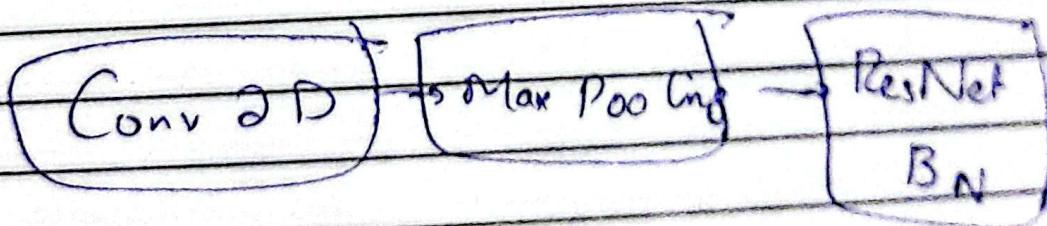
near zero

as θ

values after Iter. 9

push to other side

CNN





Addendums /

Regular \rightarrow timeseries

+ ext/ freeform

\hookrightarrow cat ate a mouse
 mouse ate a cat $\begin{matrix} \nearrow \text{from computer pol. of mode} \\ \searrow \text{sense} \end{matrix}$

cat \rightarrow mouse

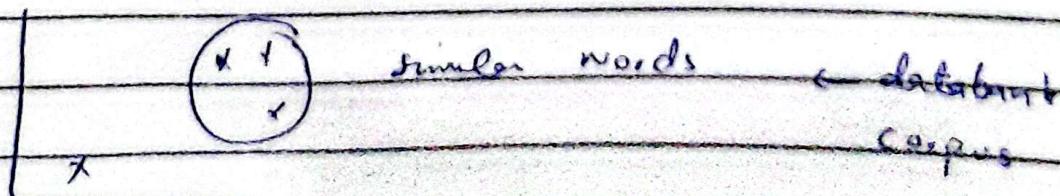
ambiguity is always there

position \rightarrow 1000	cat	cat ate
walks	ate	ate a
.	a	a mouse
.	mouse	
	1 gram	2 gram

\rightarrow vectorized \rightarrow learning
Form also

position information is important

cosine \rightarrow similarity



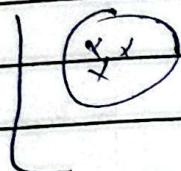
Date _____

MTWTFSS



Attention

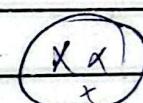
~~where cat attention is toward where - ate~~
~ mouse



nlp technique

d will never
learn

Attention



it looks

I love AI

(I) (O) (I) (I)
(O) (T) (I)

Love ↑ XAI

T

Date _____

MTWTF



vectors must be unigree and hex in size

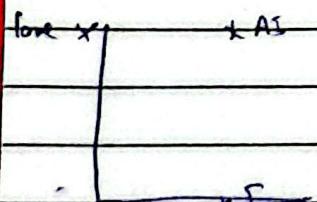
T have AS

$$\sigma(q; k_i) + \vec{v}_i$$

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

attention (Q, K, V)

(Query, key, value)



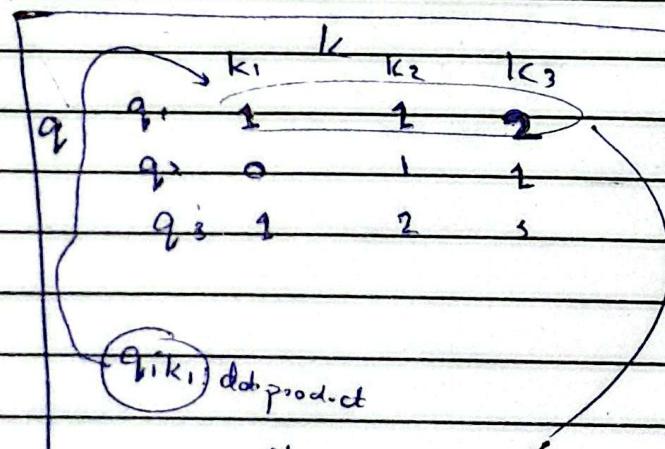
$$W_Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, W_K = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, W_V = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

For I

$$q_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$k_1 = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$v_1 = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$



For Lone

$$q_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$k_2 = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$v_2 = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\frac{e^1}{e^1 + e^1 + e^2} = \frac{2.71}{13.82}$$

$$\begin{bmatrix} 0.21 & 0.21 & 0.58 \\ 0.16 & 0.42 & 0.42 \\ 0.09 & 0.24 & 0.66 \end{bmatrix}$$

For A

$$q_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$k_3 = \begin{bmatrix} ? \\ 1 \end{bmatrix}$$

$$v_3 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Output for I

$$0.21 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0.21 \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 0.58 \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\begin{bmatrix} 0.21 + 0 + 0.58 \\ 0.21 + 0.21 + 1.16 \end{bmatrix}$$

$$\begin{bmatrix} 0.79 \\ 1.58 \end{bmatrix}$$



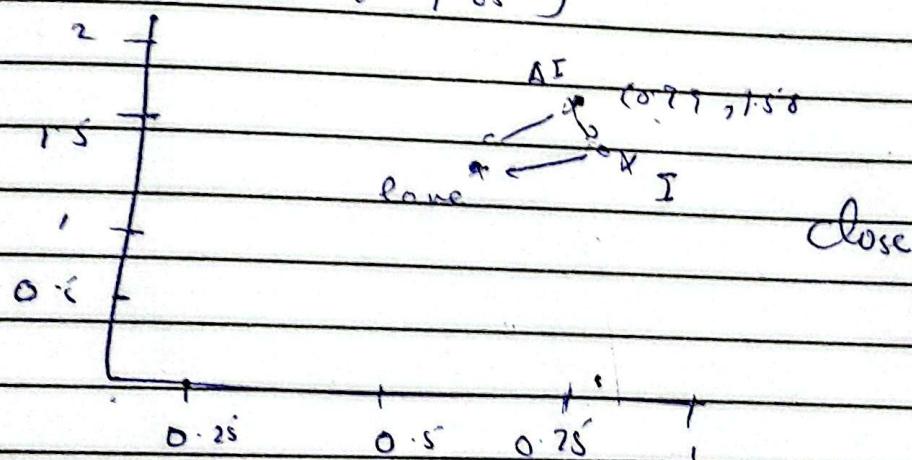
at at foot

$$0.16 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0.16 \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 0.16 \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$\left[\begin{smallmatrix} 0.88 \\ 1.42 \end{smallmatrix} \right]$

Output of AS

$$= \begin{bmatrix} 0.75 \\ 1.65 \end{bmatrix}$$



$\begin{matrix} & \begin{pmatrix} 1 & 2 \\ 1 & 1 \end{pmatrix} \\ \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} & \end{matrix}$

Σ

Weights
 startig value



Date _____

(M) (T) (W) (T) (F) (S)

Transformer \rightarrow Attention (Multi Head Multi) $x \leftarrow$ input mat $\in \mathbb{R}^{n \times d}$ $x \leftarrow ?$

matrix

s. 20

 $\in \mathbb{R}^{n \times d}$

i	love	ai	i
1	0	0	
0	1	0	$\downarrow d$
0	0	0	

 w_Q state $w_Q w_K w_V$

$$Q = x^T w_Q$$

 $\in \mathbb{R}^{n \times d}$

$$K = x^T w_K$$

 $\in \mathbb{R}^{n \times d}$

$$V = x^T w_V$$

 $\in \mathbb{R}^{n \times d}$

$$\text{softmax}(QK^T)$$

 \sqrt{d} $\in \mathbb{R}^{n \times n}$

A

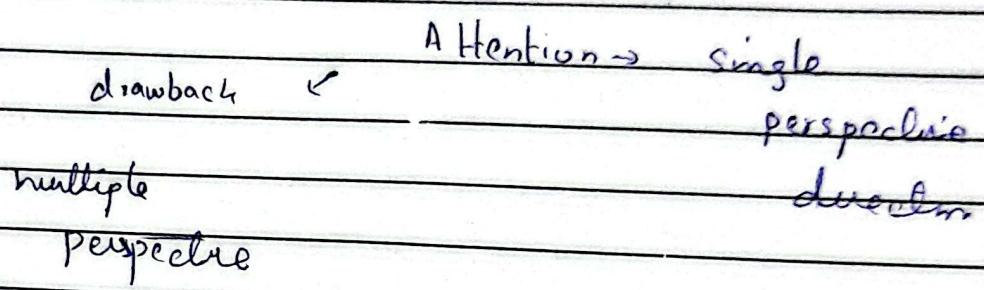
Date _____

MTWTFSS



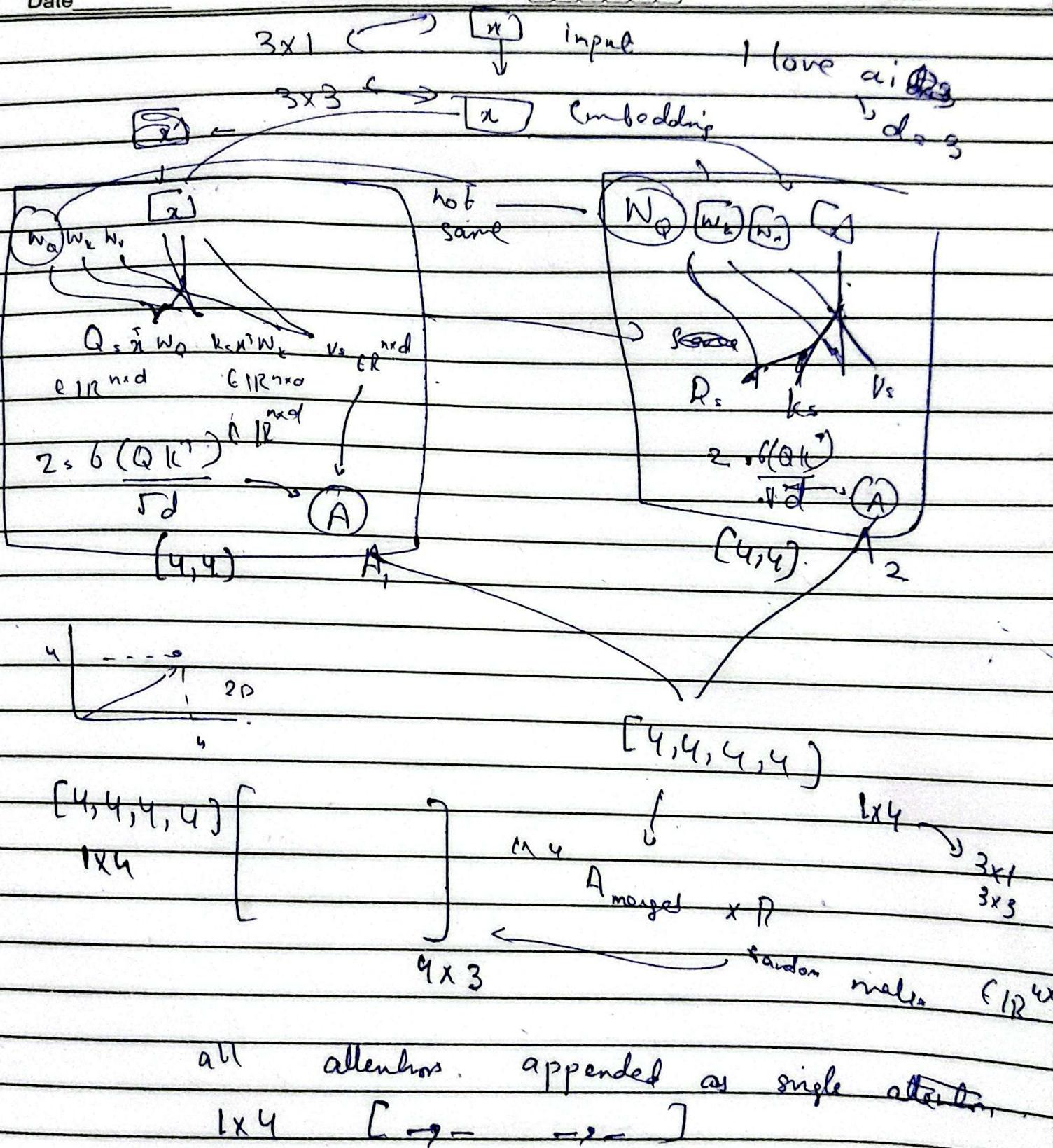
* Ambiguities

I saw the elephant with glasses.



Date _____

MTWTFSS



1 technique to keep back

to original dimension size.



DM

(MIT WIEF '18)

Positional Encoding

1 0 0 0

0 1 0 0

0 0 1 0

0 0 0 1

1) cat ate the mouse

2) mouse ate the cat

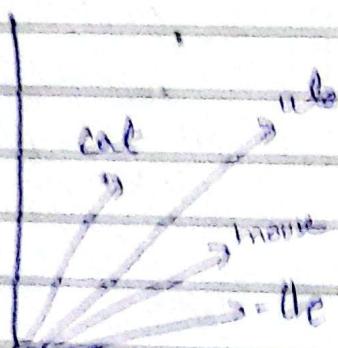
(1 0 0 0)

(0 1 0 0)

0 0 1 0

0 0 0 1

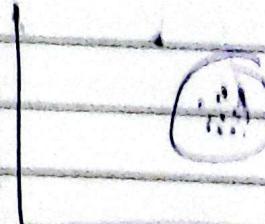
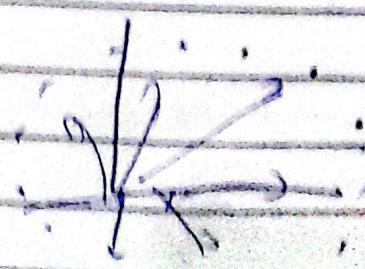
same



We can't

return our

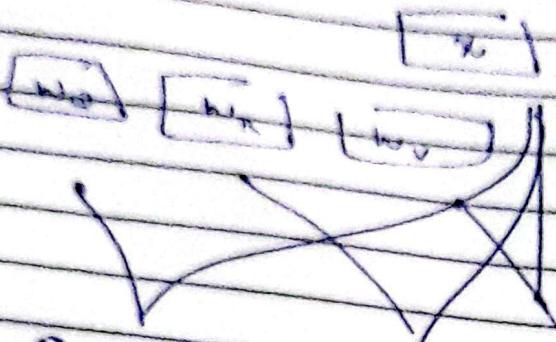
position.

Position
relationfind an
angle

Date _____



MTWTFSS



$$Q = x^T W_Q$$

G $\in \mathbb{R}^{n \times d}$

$$k = x^T W_k$$

C $\in \mathbb{R}^{n \times d}$

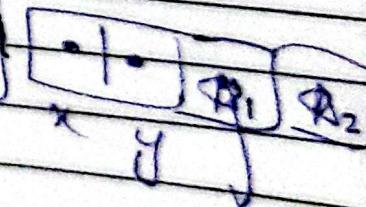
$$V = EIR^{n \times d}$$

$$Z = \frac{b(\oplus k^i)}{\sqrt{d}}$$

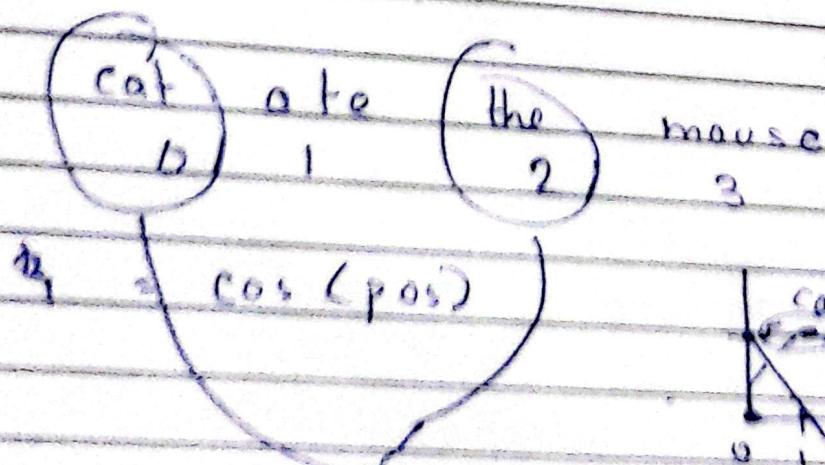
EIR $^{n \times m}$

A

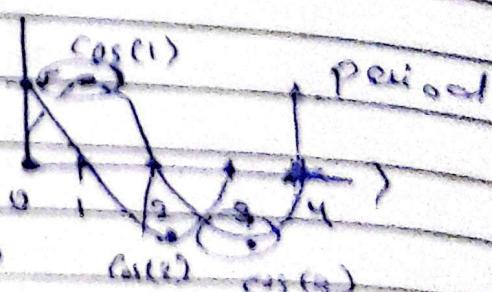
$$\delta = 2$$



positions



pos



$$f(x_2) = \sin(\text{pos}_2)$$

Date _____

MTWTF



cat ate the mouse
 pos 0 1 2 3

$$\begin{matrix} \cos(0) & \cos(1) \\ \sin(0) & \sin(1) \end{matrix}$$

cat

$$\boxed{\begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} \cos(0) & \sin(0) \end{bmatrix}}$$

ate

$$\boxed{\begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} \cos(1) & \sin(1) \end{bmatrix}}$$

$$\text{the } \boxed{\begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} \cos(2) & \sin(2) \end{bmatrix}}$$

mouse

$$\boxed{\begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} \cos(3) & \sin(3) \end{bmatrix}}$$

$$\cos(\text{pos})$$

$$\cos(\text{pos}) \rightarrow \cos(\text{pos})$$

$$\cos(\text{pos}_1)$$

$$\sin(\text{pos}) \rightarrow$$

$$\sin(\text{pos})$$

$$\sin(\text{pos}_2)$$

$$(n) (n \times 1)$$

$$1$$

$$\sin(\text{pos}_1) \quad \cos(\text{pos}_1)$$

$$(n \times d)$$

$$\vdots \quad \vdots$$

$$\sin(\text{pos}_i)$$

$$\cos(\text{pos}_i)$$

$$d = \frac{d}{\delta}$$

$$d = 8$$

$$i = 8, 4$$

$$i = 0, 1, 2, 3$$

$$d = 4$$

$$i = \frac{4}{2} = 2$$

$$i = 0, 1$$

Subject: _____

Sr. No.	Date	$d = 4$	$\frac{y_2}{2} \rightarrow 2(1)$ Topics
			$ = 0 \quad i = 1$
x	pos	cos sin	cos sin
cat	0	$P(0,0) P(0,1)$	$P(0,2) P(0,3)$
ate	1	$P(1,0) P(1,1)$	$P(1,2) P(1,3)$
fl	2	$P(2,0) P(2,1)$	$P(2,2) P(2,3)$
mouse	3	$P(3,0) P(3,1)$	$P(3,2) P(3,3)$
		<u>Unique position</u>	
$\cos \rightarrow$	$P(0,0), 2i$	$\cos \left(\frac{90^\circ}{100,000}, 2, 6d \right)$	
$\sin \rightarrow$	$P(0,0), 2it$	$\sin \left(\frac{90^\circ}{10,000}, 2, 1d \right)$	

$d = 4$

$$\frac{di}{d} \quad \frac{2}{4} : i \quad \frac{1}{2} : i$$

$$\begin{bmatrix} \cos\left(\frac{0}{100}\right) & \sin\left(\frac{0}{100}\right) & \cos\left(\frac{1}{100}\right) & \sin\left(\frac{1}{100}\right) \\ \cos\left(\frac{1}{100}\right) & \sin\left(\frac{1}{100}\right) & \cos\left(\frac{2}{100}\right) & \sin\left(\frac{2}{100}\right) \\ \cos\left(\frac{2}{100}\right) & \sin\left(\frac{2}{100}\right) & \cos\left(\frac{3}{100}\right) & \sin\left(\frac{3}{100}\right) \\ \cos\left(\frac{3}{100}\right) & \sin\left(\frac{3}{100}\right) & \cos\left(\frac{4}{100}\right) & \sin\left(\frac{4}{100}\right) \end{bmatrix}$$

$$1000^{\circ}/2 = 1 \\ 1000^{\circ}/2 = 100$$

If solve ambig angles related
to position in support of
altitudes

Date _____

MTWTF

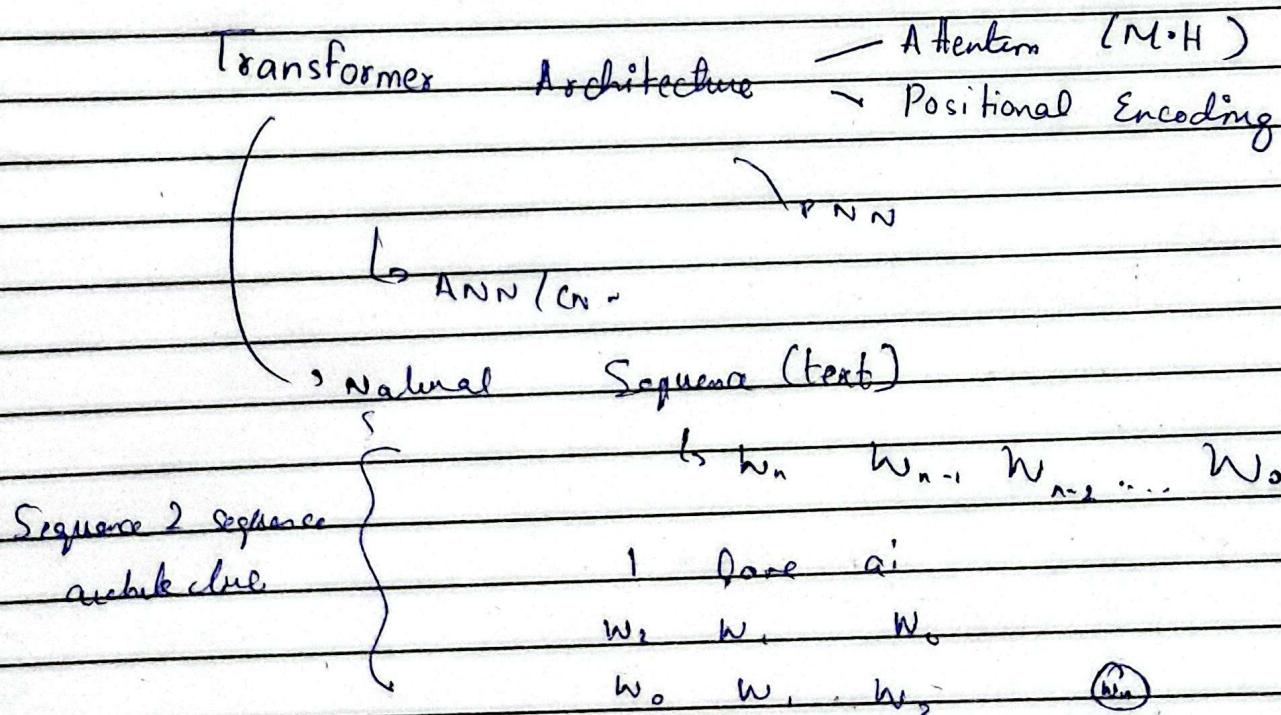


for periodic dep. calm increase
cycles

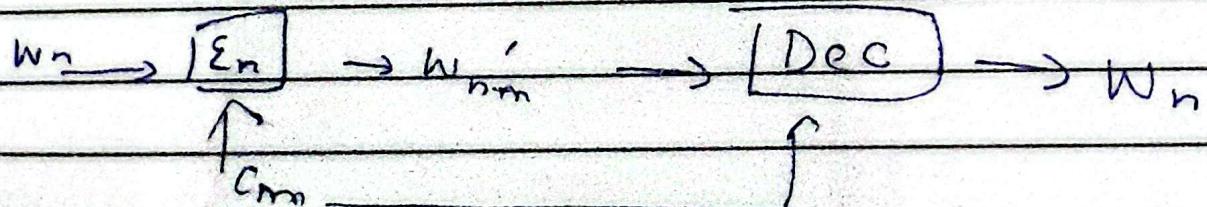
safe figure 10 multiple
decreas

log ~
scale

$10,000^i$



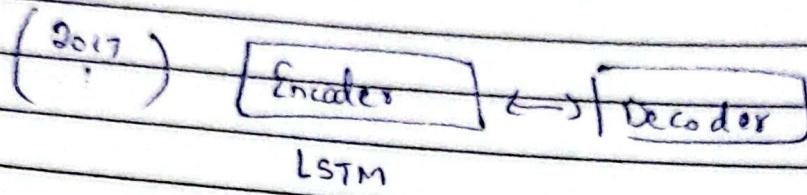
Encoder \leftrightarrow Decoder





Date

M T W T F S

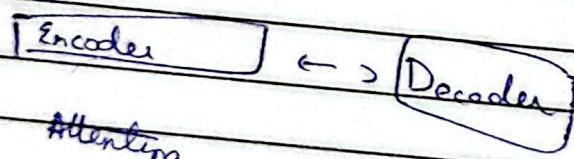


① small text length w.r.t threshold

↳ bad performance on large language

General test

Transformers



$$x' \xrightarrow{\text{EIR}^{n \times 1}} x \xrightarrow{\text{EIR}^{5 \times 2}} \text{fined}$$

$$n = 9$$
$$d = 5_{12}$$

~~1.1.1.1.0000 0000~~

h = 5

• • • (.) - [0 0 0 0 0]

embeddins

Date _____

(M T W T F S)



embeddings

I. Done AT

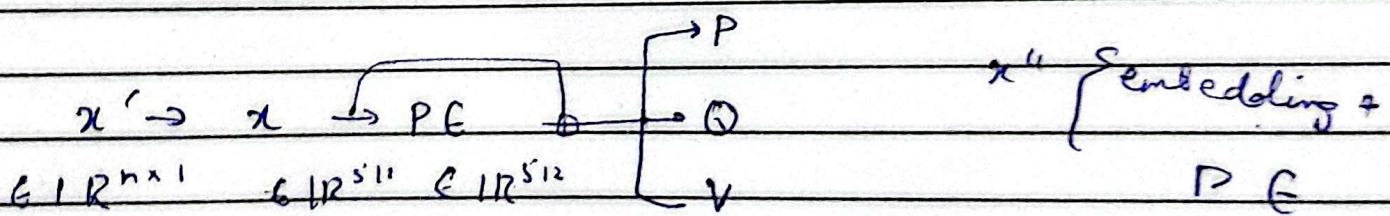
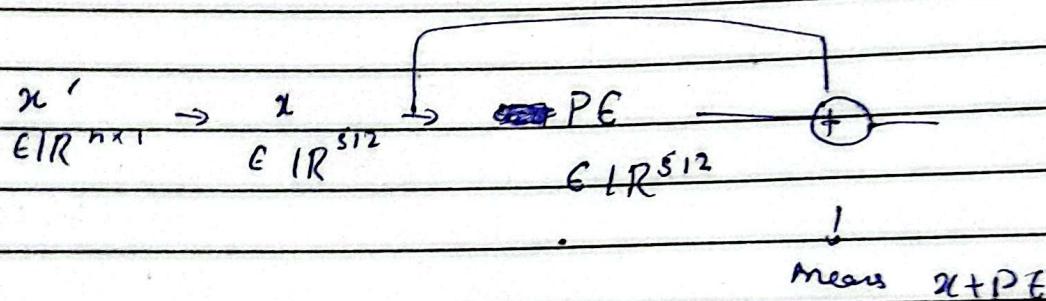
 $T \in \mathbb{R}^{n \times d}$

Done

AT

-

Skipgram ... -- one hot

 \hookrightarrow size
vector


$x'' \rightarrow$ Multi Head Attention

$W_p \cdot x'' = P$

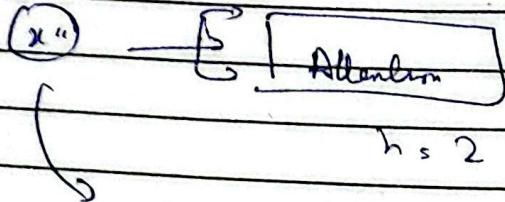
$$A = \left(\frac{\sqrt{d}}{\sqrt{d}} \right) V$$

Date _____

MTWTF



Castelli

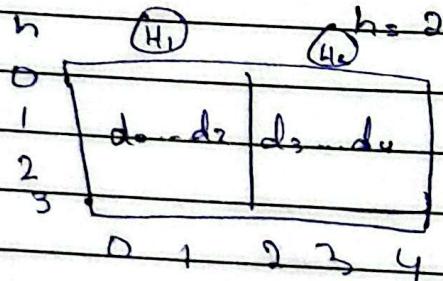


$$A = G \left(\frac{k Q^T}{\sqrt{d}} \right) V$$

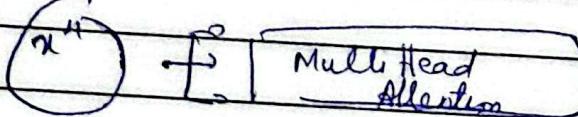
short sentences

ambiguities of sentence length

PAs



(1)

O(pn³)(O(h(n³+n³)) → O(n³))O(n⁴)

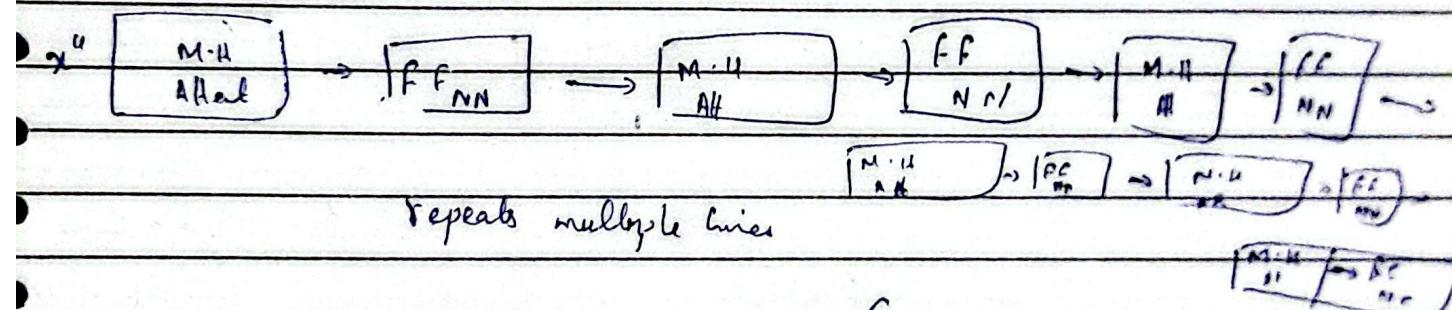
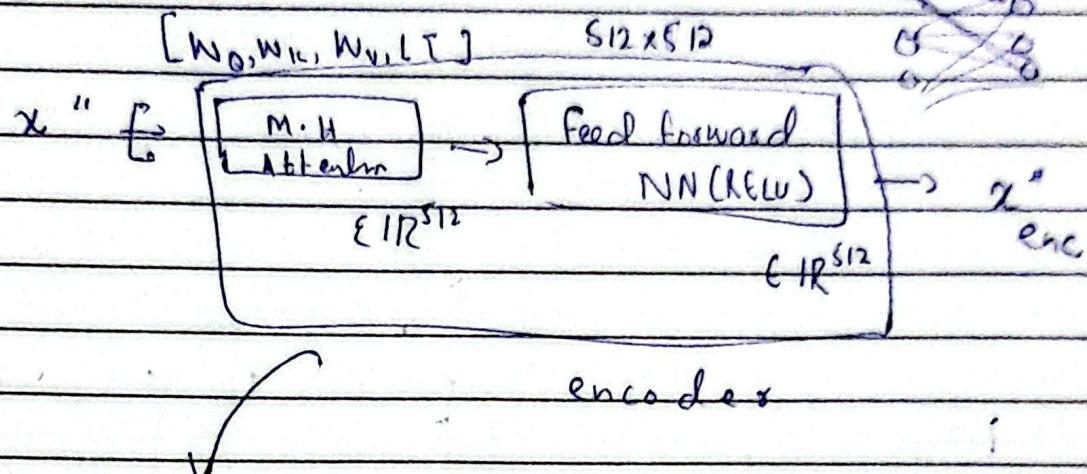
$$A = G \left(\frac{k Q^T}{\sqrt{d}} \right) V$$

bias, gender,
racialWe assume
 $p < n$ $p = 2, 3, 4, 8$ extreme case
if p is large
 $\approx n$ h should be
increased

Date _____

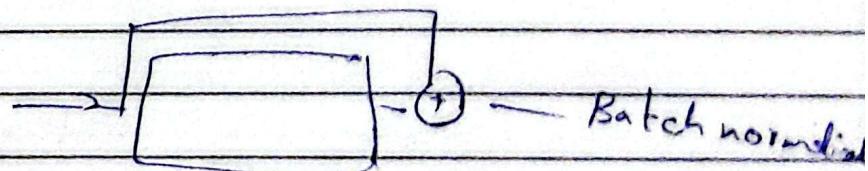
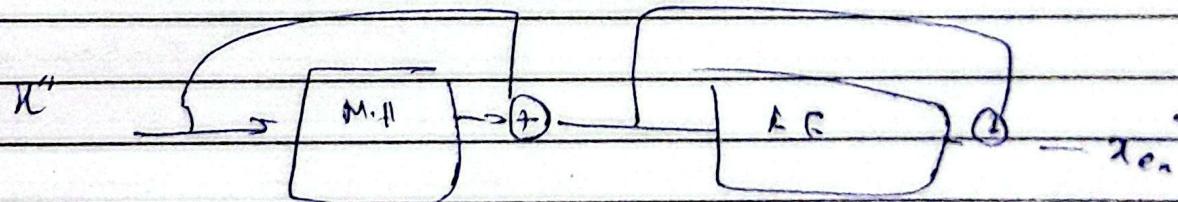
M T W T F S

3/3



6

Vanishing gradient problem if depth increases



Date _____



(M T W T F S)

Q Here we have

layer Normalization

$f_1 \ f_2 \ y_i$
- - -
- - -
- - -
- - -

	f_1	f_2	f_3	...
d_1				
d_2				
	x_0	x_1	x_2	x_3

δ

$x'_0 \dots x'_n \rightarrow 0 = 0$

$6 = 1$

Batch

Normalizat_n
for each row

→ layer normalization for each column.

1 2 3

4 5 6

Feature types

no difference

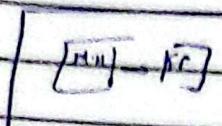
f_1, f_2 etc

branch

batch + layer batch

Encoder

Decoder

 $\in M^{64 \times 64}$

④

ART

⑧ $\rightarrow P_{\text{Encoding}} \rightarrow \text{Masked Multi Head Attn.} \rightarrow \text{Feed Forward} \rightarrow$

h

Glimpses

output

Decoder



I love ai and fast $\rightarrow ?$
 $\rightarrow \rightarrow \rightarrow \rightarrow \rightarrow$

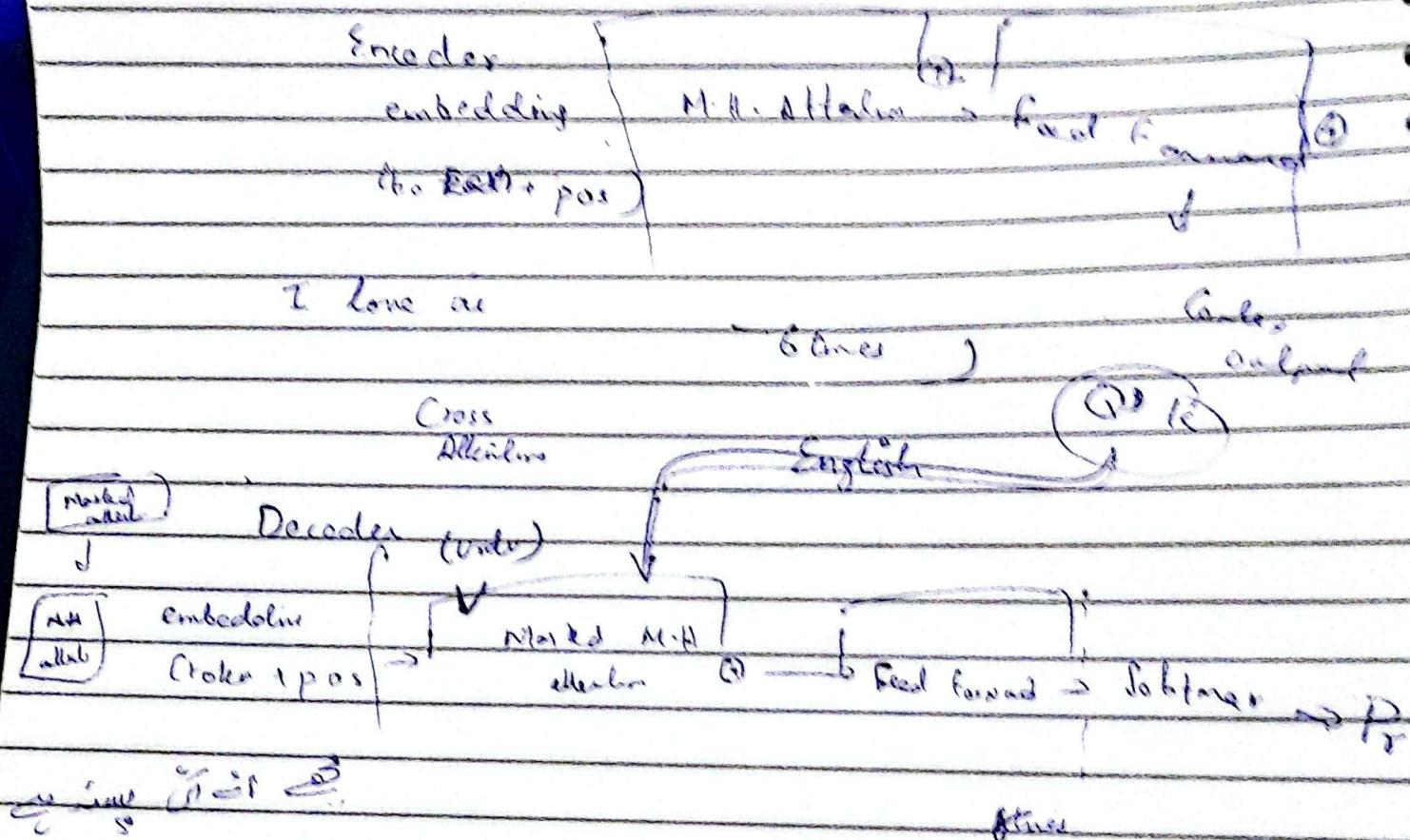
$$x'' = [i \text{ love } \underset{\text{mask}}{\text{fast}}] + [o \text{ o-o-o}]$$

$[-, -, \frac{\text{negative}}{\text{large value}}]$

BERT / pos enc TA / GPJ

Date

(M) (T) (W) (T) (F) (S)

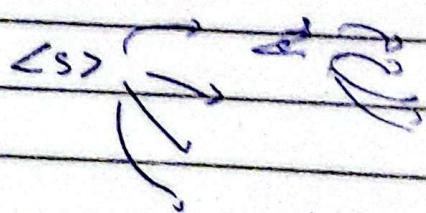


	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
CS	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
SS	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	0.1	0.2

Probability distribution over vocabulary.

Modeling select one word
Hide all others needed

= attention (Q, k, v)



MTWTF(S)

Generalism

after named model

just do encoding