

3rd March 2025:

Sequences:

MLP Fixed length data
CNN 3

Language Translation Length of i/p & o/p not fixed.

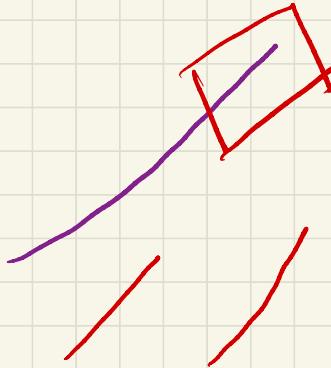
→ Images of Varying resolution.

medical records.

Varying length i/p : → fixed length problem.
"fixing the len sequence"

Collection of items in an order : Sequences.

Day	Temp
1	30
2	32
3	33
4	31
5	?

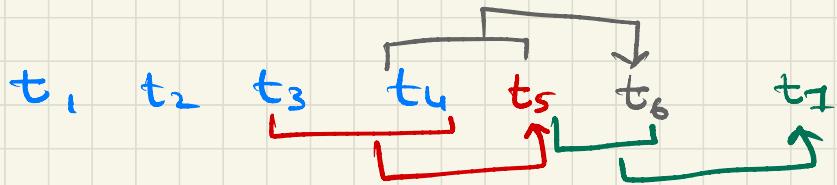


how much to go back?

taking Only the past 2 days.

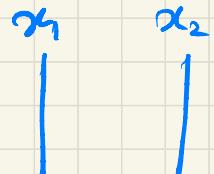
$$t_5 = b + w_1 t_4 + w_2 t_3$$

$$t_6 = b + w_1 t_5 + w_2 t_4$$

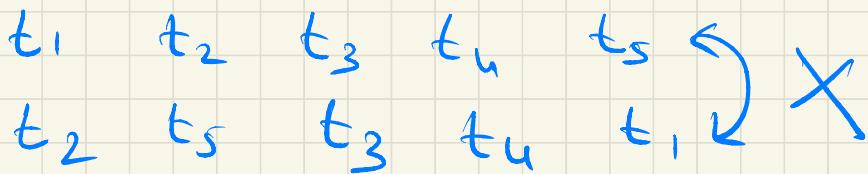


Auto Regressive models.

Self looks back in time.



height	weight		height	weight
h_1	s_1		h_1	s_1
h_2	s_2	Same	h_m	s_m
h_3	:		h_2	s_2
:	:		h_3	s_3
h_m	s_m		:	:



$t_0 \quad t_1 \quad t_2 \quad t_3 \quad t_4 \quad t_5 \quad t_6$

$\begin{matrix} X \\ t_0 & t_1 & t_2 \\ t_1 & t_2 & t_3 \\ t_2 & t_3 & t_4 \\ t_3 & t_4 & t_5 \end{matrix}$

$\begin{matrix} Y \\ t_3 \\ t_4 \\ t_5 \\ t_6 \end{matrix}$

fix "p"

$$AR(p) : x_t = b + \sum_{i=1}^p w_i x_{t-i}$$

$p=1 :$

$$AR(1) : x_t = b + w_1 x_{t-1}$$

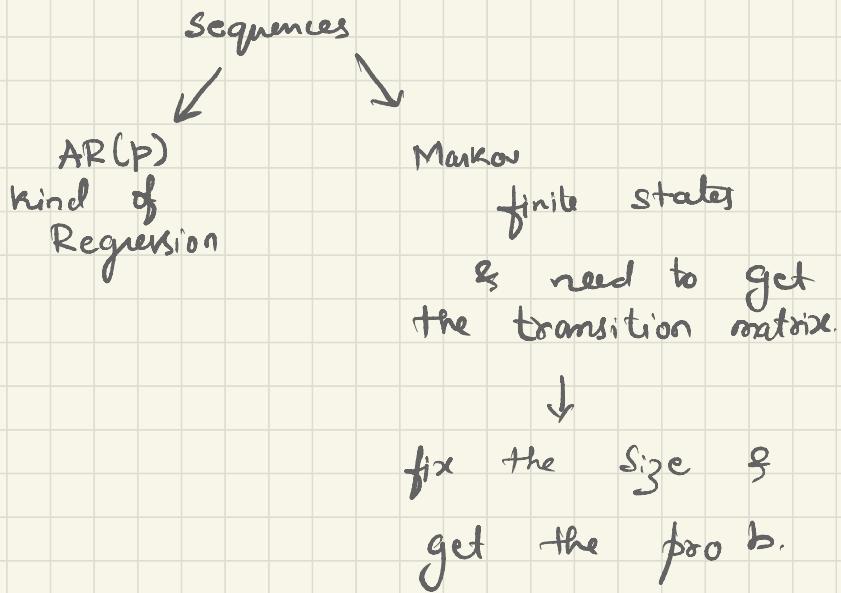
$$x_2 = b + w_1 x_1$$

$$x_3 = b + w_1 x_2$$

Markov:

O/p: next state

how many next states?



Text to Numbers?

I am enjoying my studies.

0 1 2 3 4

I love my country.

0 1 2 3

(242)

Oxford English Dict.

1 world

2

3

.

.

242 : country.

enjoy : 821

Coffee is great

car battery is dead, oh great

each word Some number based on
the index of constructed dictionary

I am from Mysore]
My name is Raghava]
I like eating Dosa.

- Tokenisation:



Vocabulary.

6th Mar 2025:

✓ fixed length i/p \rightarrow fixed length o/p
 \rightarrow variable length i/p \rightarrow fixed length o/p

Named Entity Representation: NER

XX Harry Potter won the Triwizard Cup, XX
| | O O | O

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6$

x_t : t^{th} feature

$x_{i,j}$: j^{th} feature in i^{th} example.

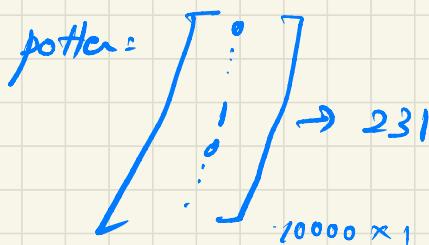
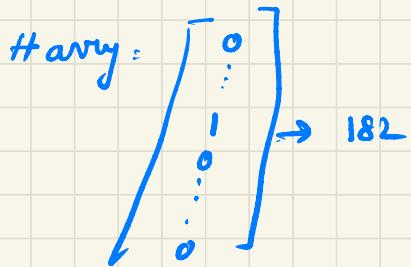
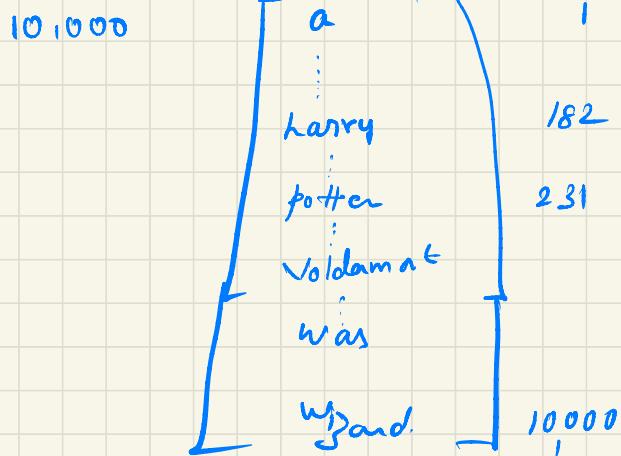
In this specific case of NER.

$y_1 \quad y_2 \dots \quad y_6$

$y_{i,j}$: j^{th} output of i^{th} feature

Construct Vocabulary :

num of unique tokens (Size of vocabulary) is

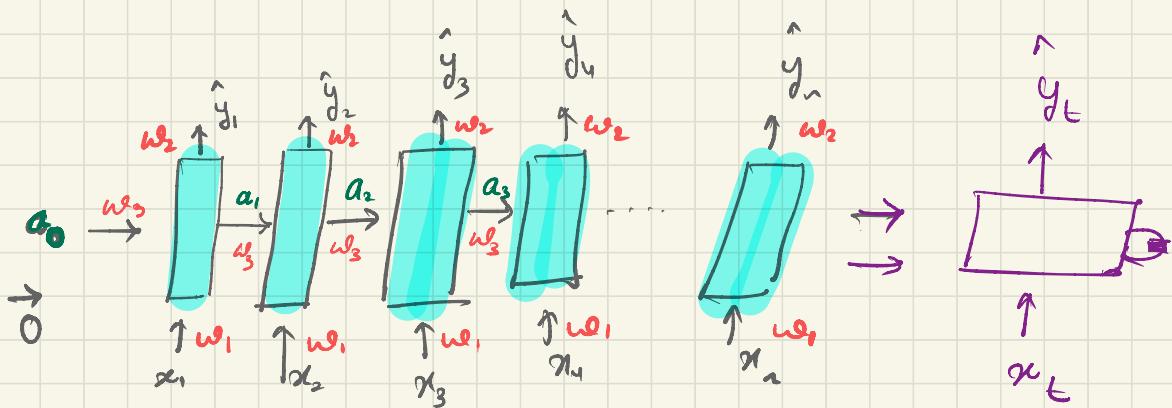


model to map from $X \rightarrow Y$.

word not in Vocabulary \rightarrow **Lunk**

→ we can perform zero padding which is inefficient.

RNNs: Recurrent Neural Networks:



input is given from left to right

$$a_0 = \vec{0}$$

$$a_1 = g(w_1 \cdot x_1 + w_3 \cdot a_0 + b_1)$$

$$\hat{y}_1 = g_2 (w_2 \cdot a_1 + b_2)$$

$$a_2 = g_1 (w_1 \cdot x_2 + w_3 \cdot a_1 + b_1)$$

$$\hat{y}_2 = g_2 (w_2 \cdot a_2 + b_2)$$

\$\tanh / \text{ReLU}\$: Activation.

$$a_t = g_1 (w_1 \cdot x_t + w_3 \cdot a_{t-1} + b_1)$$

$$\hat{y}_t = g_2 (w_2 \cdot a_t + b_2)$$

10th March 2025:

Problem with RNNs

→ Teddy Roosevelt was the president of USA.

1 1 0 0 0 0 1

→ Teddy Bears are for sale in Big bazaar.

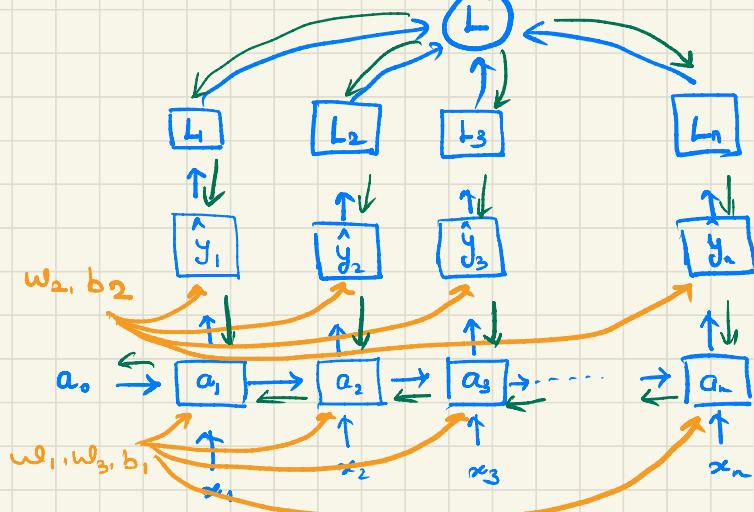
0 0 0 0 0 0 1

Note: to decide the NER we need to look toward for future words. Which is not possible in RNNs.

Bi-directional RNNs.

How to update the weights.

Computational graph



Back Propagation

Through Time

Binary cross entropy loss :

$$L_i(\hat{y}_i, y_i) = -y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$$

$$L(\hat{y}, y) = \sum_{i=1}^n L_i(\hat{y}_i, y_i)$$

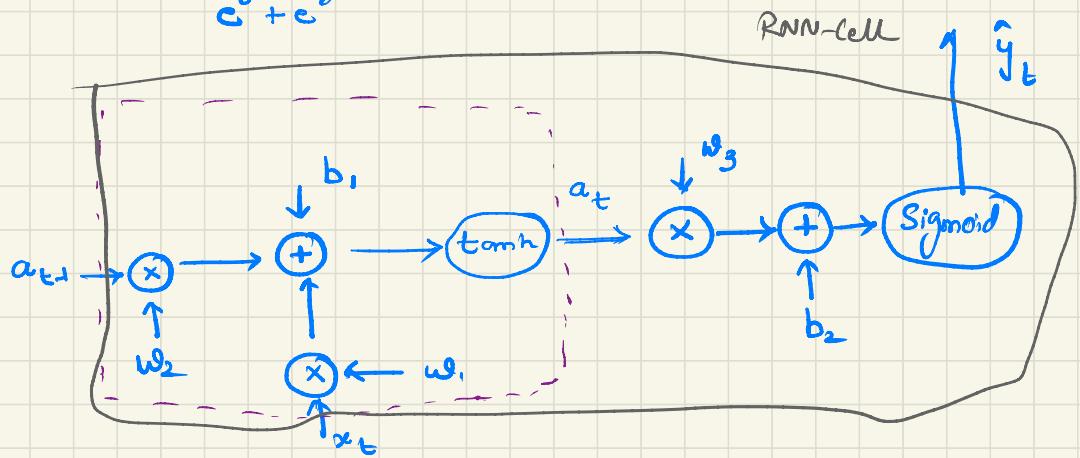
$$\begin{bmatrix} \hat{y}_1 \\ 1-\hat{y}_1 \end{bmatrix} \begin{bmatrix} \hat{y}_1 \\ 1-\hat{y}_1 \end{bmatrix}$$

11th March 2025:

$$a_t = \tanh(w_1 x_t + w_2 a_{t-1} + b_1)$$

$$\hat{y}_t = \text{Sigmoid}(w_3 a_t + b_2)$$

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



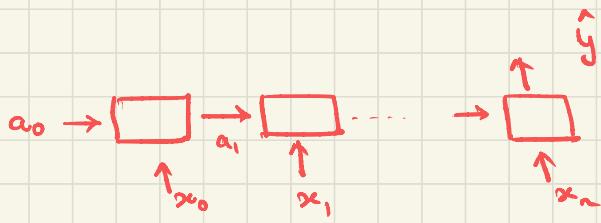
NER model we are looking is a Many to Many model.
& size of inputs & output were same.

→ Movie Ratings

input: Review

Output: number of stars {1, 2, 3, 4, 5}

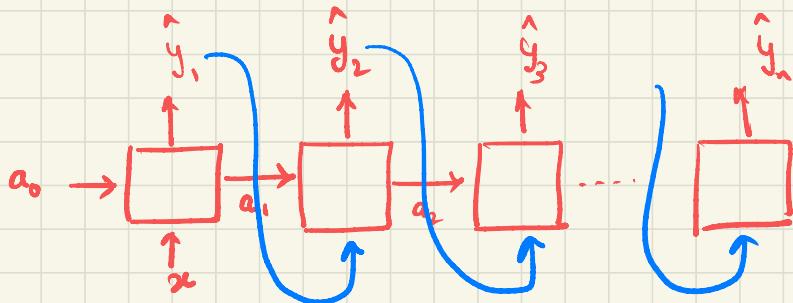
Many to One Mapping.



→ given the base node, play some Raga.



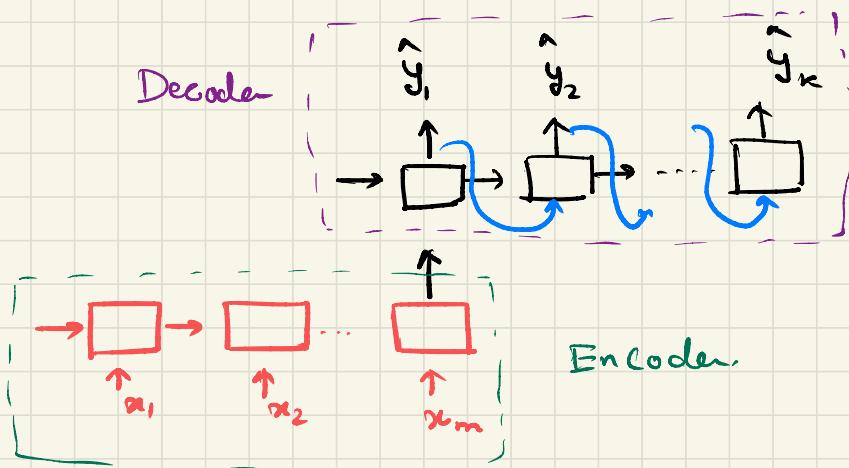
One to Many Problem.



Machine Translation :

Many to Many.

The length of
Input & Output
are diff



consider speech to Text

1: The apple & Pear Salad was good.

2: The apple & pear Salad was good.

Occurrence (S_2) > Occurrence (S_1)

I have total of "N" sentences in my Corpus.

$$\frac{o(S_2)}{N} > \frac{o(S_1)}{N}$$

$$\underline{IP[S_2]} > \underline{IP[S_1]}$$

I like eating Pongal

$x_1 \quad x_2$

$x_3 \quad x_4$

word Translation problem.

length of i/p & o/p is same.

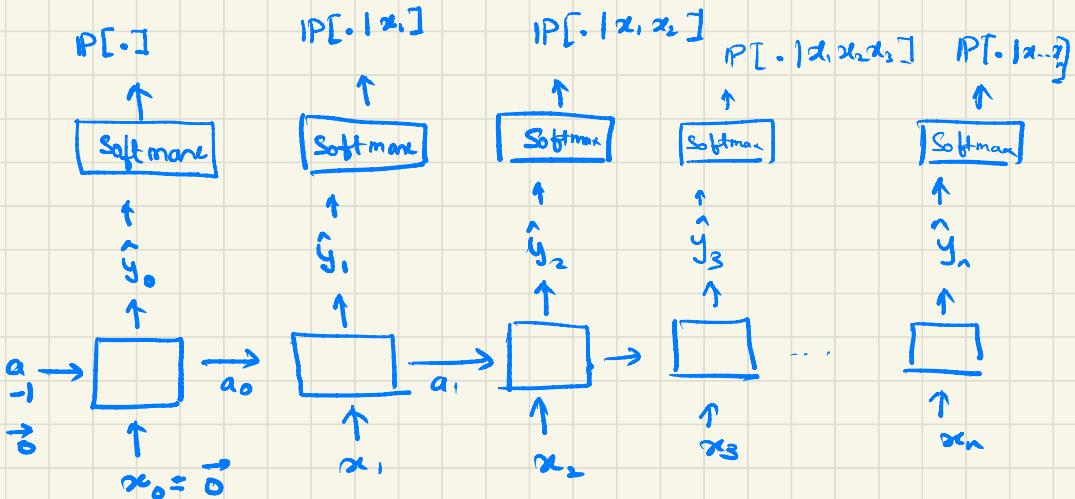
Output should choose one word out of all the words in the vocabulary.

\therefore need to use softmax.

Consider there are 10 words in vocabulary

then we should get a probability vector of

length 10.



[Eat, the life is going on, well, orange is good]

The —

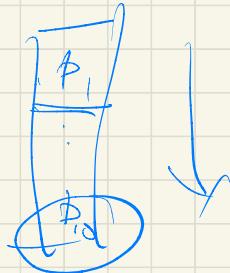
The orange —

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \rightarrow \begin{bmatrix} s \\ o \\ r \\ m \\ x \end{bmatrix} \rightarrow \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

$$p_1 = \frac{e^{q_1}}{e^{a_1} + e^{q_2} + e^{q_3}}$$

$$p_2 = \frac{e^{q_2}}{e^{a_1} + e^{q_2} + e^{q_3}}$$

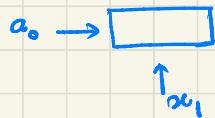
$$p_3 = \frac{e^{q_3}}{e^{a_1} + e^{q_2} + e^{q_3}}$$



12th March 2025!

Size of Vocab: 5000

- S_1 : I like eating close
 S_2 : I am from myself
 S_3 : we are not in Mars
 S_4 : I like teaching
 S_5 : Once there was a lion



→ Size of Vocabulary : 5000

each word is a vector in 5000 dim $x_i \in \mathbb{R}^{5000}$

→ Consider the batch of size "m"

if we want to give the batch as input

then what will be the size 5000×5

input @ every time step :

5000×5

→ what will be the size considering all time stamps together.

max time stamp : 10

then @ max the whole input is $5000 \times 5 \times 10$

$x_1 \ x_2 \ \dots \ x_{10}$

$$a_t = (w_1 x_t + w_2 a_{t-1} + b_1)$$

$$x_t \in \mathbb{R}^{5000 \times 5}$$

$$\begin{array}{l} w_1 \in ? \\ \quad \quad \quad \textcircled{231} \quad \mathbb{R}^{10} \\ \quad \quad \quad ? \\ w_1 \in \mathbb{R} \end{array}$$
$$\underline{\quad \quad \quad \times 5000}$$

→ Size of w_1 & w_2 are dependent on size of a_t

→ this can be fixed by user : 231

that means for one word a_t generated is a vector of size 231

∴ for a batch 231×5

∴ for all the time stamps $231 \times 5 \times 10$

$$w_2 \cdot a_{t-1}$$

$$\downarrow \quad \downarrow$$

$$231 \times 231 \quad 231 \times 5$$

$$231 \times 5000$$

$$w_1 \in \mathbb{R}$$

$$231 \times 231$$

$$w_2 \in \mathbb{R}$$

$$a_t = w_1 * x_t + w_2 * a_{t-1} + b$$

$$= \underbrace{(231 * 5000) (5000 * 5)}_{(231 * 5)} + (231 * 231) (231 * 5)$$

$$a_t = (231 * 5) + (231 * 5) + \boxed{231 * 1}$$

$$\hat{y} = (w_3 * a_t + b_2)$$

$$a_t \in \mathbb{R}^{231 * 5} \quad w_3 \in \mathbb{R}^{a_L * 231} \quad b_2 \in \mathbb{R}^{(231 * 5)}$$

what should be the size of output

$$w_3 \in \mathbb{R}^{5000 \times 5}$$

$$b_2 \in \mathbb{R}^{5000 \times 1}$$

$$w_1 \in \mathbb{R}^{231 \times 5000}$$

$$w_2 \in \mathbb{R}^{231 \times 231}$$

$$w_3 \in \mathbb{R}^{5000 \times 231}$$

$$b \in \mathbb{R}^{231 \times 1}$$

$$b_2 \in \mathbb{R}^{5000 \times 1}$$

Vanishing gradients:

next word prediction.

The man, who was is

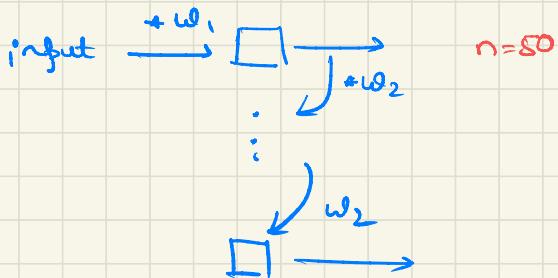
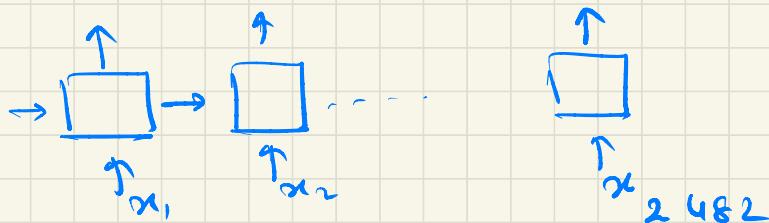
The men, who were are

→ long term dependencies.



200 layer MLP.

Vanishing gradients. (NaN)



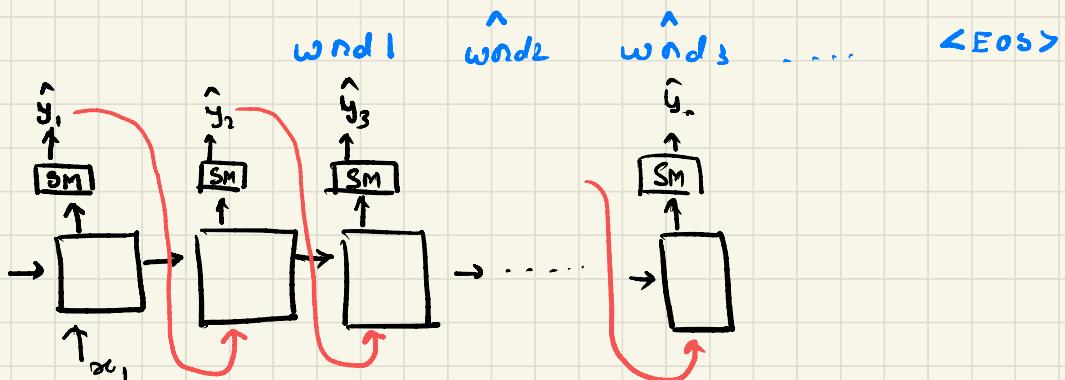
$n=80$

$$\text{input} * w_1 * w_2^n$$
$$\frac{\partial}{\partial w_1} \approx \text{input} * w_2^n$$
$$\text{if } w_2 = 2 \rightarrow \frac{\partial}{\partial w_1} = \text{input} * 2^{80}$$
$$w_2 = 0.5 \Rightarrow \frac{\partial}{\partial w_1} = \text{input} * (0.5)^{80}$$

13th March 2025.

Sequence generation:

Next word Prediction.



GRU : Gated Recurrent Unit.

What to remember. Problem: Sequence generation.

The cat is . . .

we need to remember that we have used singular earlier & generate the accordingly.

we introduce an idea of remember entity "Gt"

$C_t = a_t$ is same in GRU.

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

Previous
Content

$$w_c = w_c^1 C_{t-1} + w_c^2 x_t \\ = w_c \cdot [C_{t-1}, x_t]$$

$$f_u = \sigma(w_u[c_{t-1}, x_t] + b_u)$$

how much to remember from previous.

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

$$f_u = 1$$

↑

The Cat is

$$\begin{matrix} \downarrow & \downarrow & \downarrow \\ f_u = 0 & f_u = 0 & f_u = 0 \end{matrix}$$

→ worst value of c_t is pretty much retained so there is no problem of vanishing gradients.

Singular / Plural

Past / Present

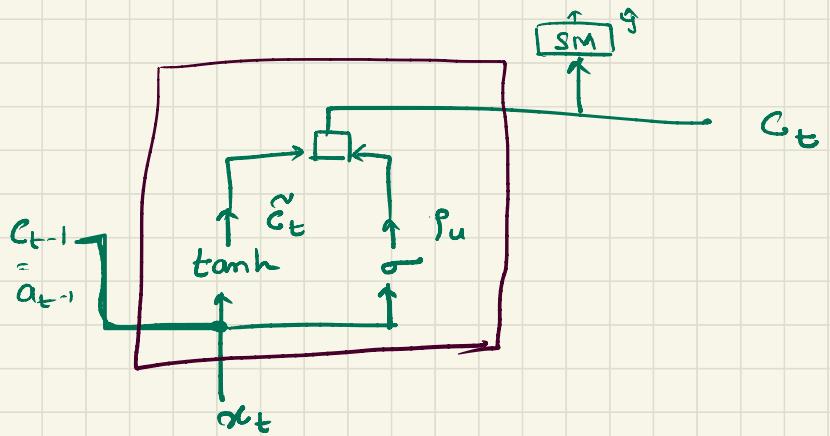
:

$$\left[\cdot \right]$$

c_t is a vector.

then f_u will also be a vector.

→ f_u will say which terms to update.
(bits)

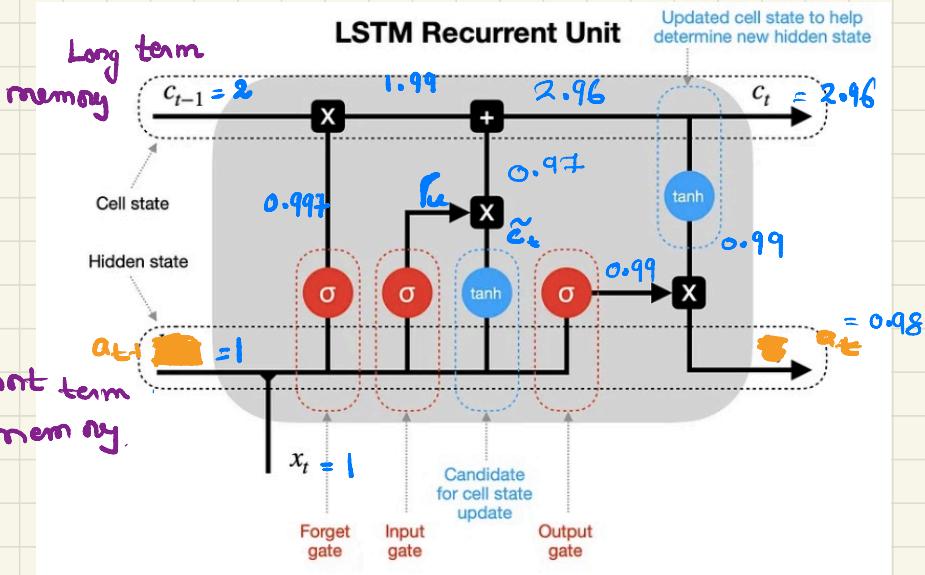


LSTM : Long Short term Memory

- Better than GRUs to get long term dependencies
- (1997)
- a_t & c_t are different

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

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



+ 2 separate paths to make predictions

* long term memory

* short term memory.

* upper path for long term its modified by multiplication & addition \Rightarrow No weights & Biases.

Forget Gate:

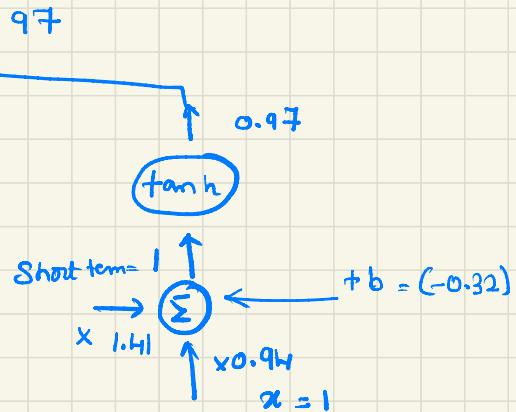
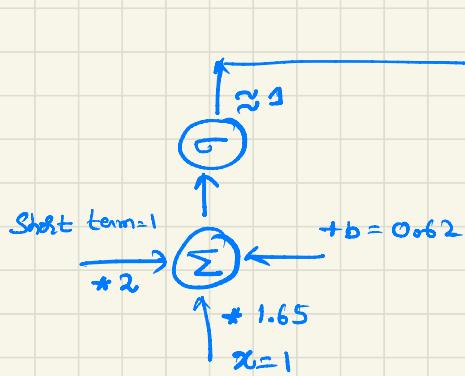
What % of long term memory do we need to remember

$$(1) \rightarrow 0.997$$

$$\begin{aligned} \text{short term} &= 1 \\ \xrightarrow{x \cdot 2.7} &\Sigma \\ &\leftarrow + b = 1.62 \\ &\times 1.6 \\ &x_t = 1 \end{aligned}$$

$$\begin{aligned} \sigma(1 + 2.7 + 1 + 1.6 + 1.62) \\ = \end{aligned}$$

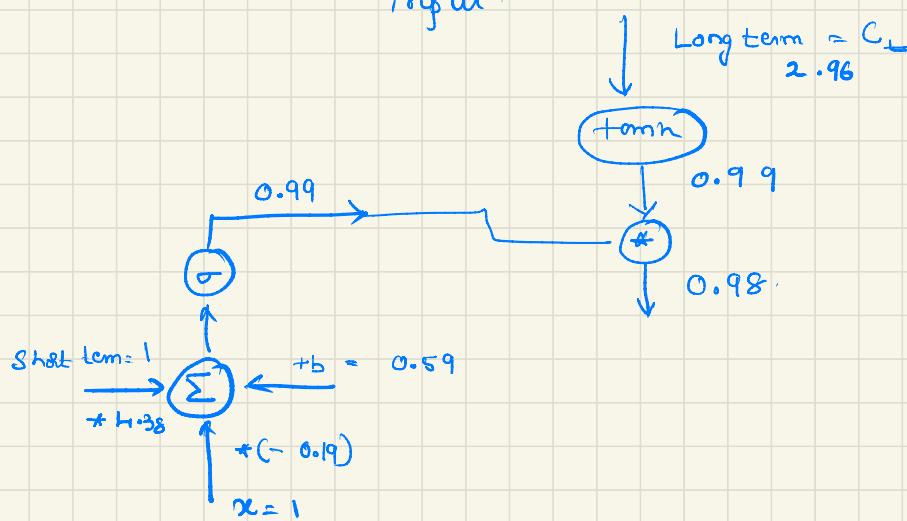
Input gate: Potential memory need to added to the long term based on short term & input.



Output gate:

Potential short term memory obtained by modulating.

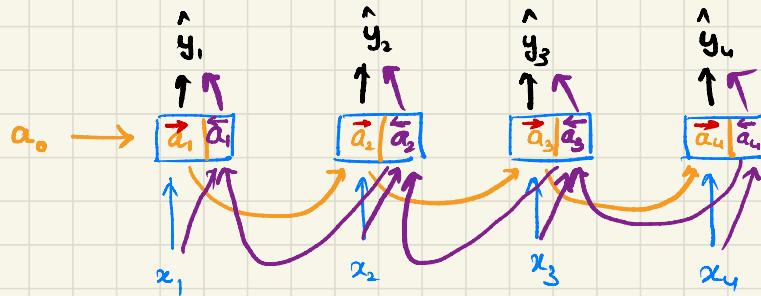
Long term
Short term
input



RNNs
GRUs
LSTMs

\Rightarrow Parsing i/p from left to Right

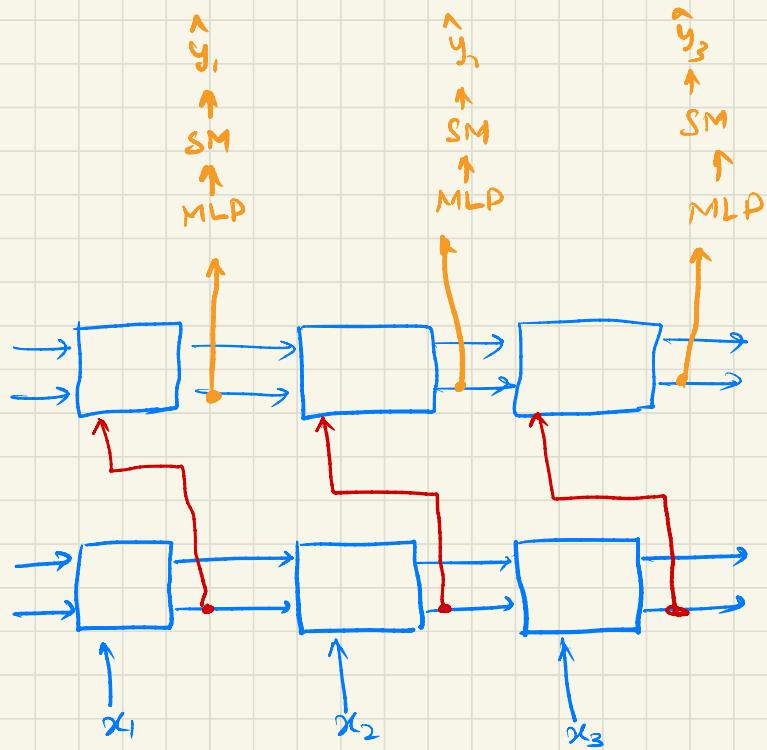
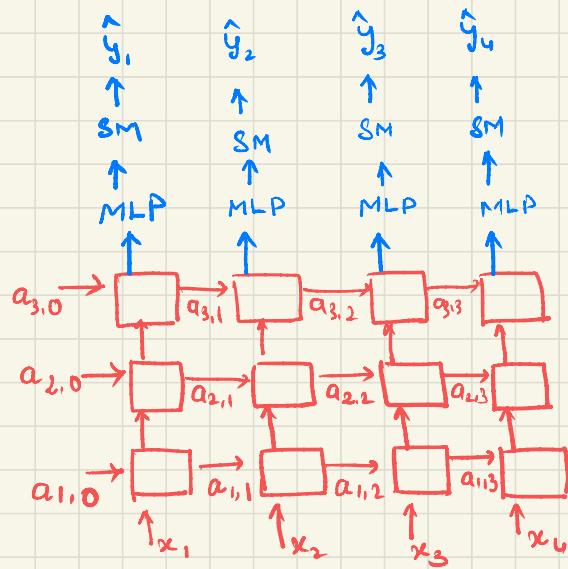
Bi-Directional RNNs

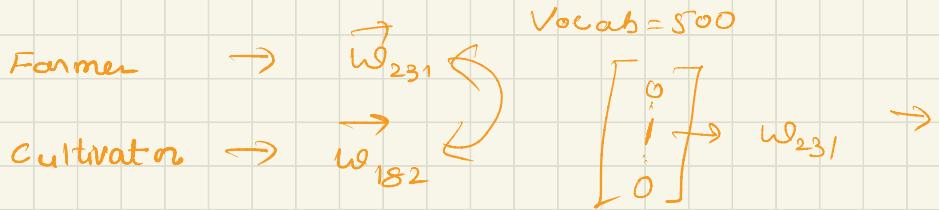


$$g(\omega_1 \vec{a}_i + \omega_2 \tilde{\vec{a}}_i + b)$$

bi-directional GRUs
bi-directional LSTMs

Deep RNNs :





Similarity

$$w_{231} \cdot w_{182} = 0$$

Word Embeddings:

Vocab Size: 500

Man
 w_{142}

Woman
 w_{420}

King
 w_{111}

Queen
 w_{341}

I want a glass of orange —

I want a glass of apple —

* relation b/w orange & apple is not easy to generalize.

* Dot product of orange & apple is zero.

\Rightarrow there is no relation, which is not true.

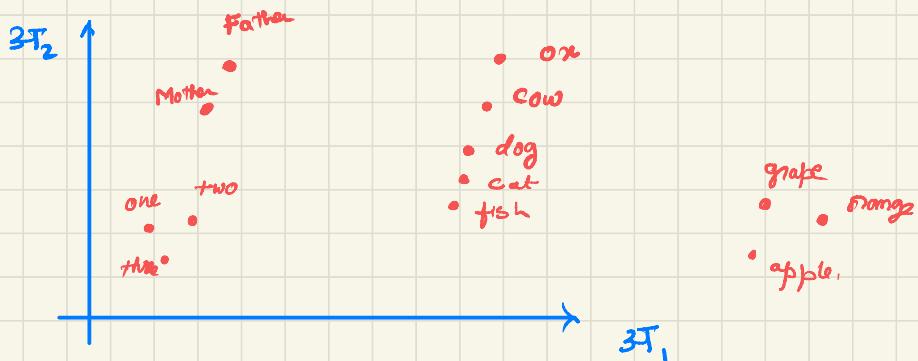
	Cow	Ox	Mother	Father	Apple	Orange
animal	0.99	0.99	0.01	0.01	0.01	0.01
Food	0.01	0.01	0.01	0.01	0.99	0.99
Relation	0.01	0.01	0.99	0.99	0.01	0.01
Gender (+1, -1)	+1	+1	-1	+1	0.01	0.01
.						
:						
:						
:						

c_i : feature vector of i^{th} word.

* we consider some 182 feature

* now the dot product b/w the feature vectors tell us how they are related.

* t-SNE Plot : 2D
Visualizing word embedding.



* huge Corpus of text

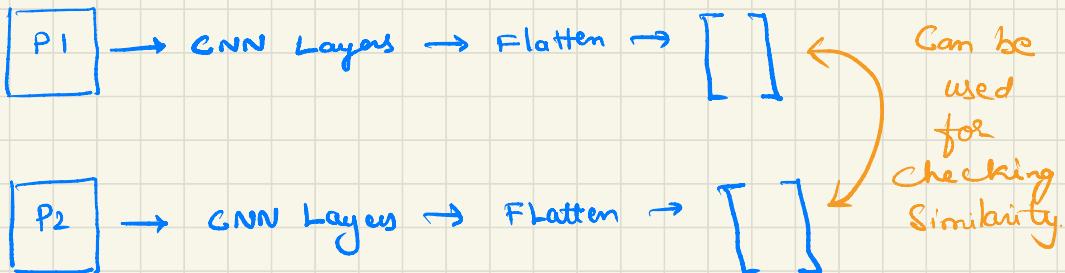
1 Billion words $\xrightarrow{\text{learn}}$ word embeddings

we can take these word embeddings which are learnt on huge Corpus and use for our task.

"Transfer Learning"

* Fine tuning: train on your data for a very small time.

*



$$\text{Cow} = \begin{bmatrix} 0.99 \\ 0.01 \\ 0.01 \\ -1 \end{bmatrix}$$

$$\begin{array}{l} \xrightarrow{\text{Cow}} \\ = \begin{bmatrix} -0.2 \\ 0 \\ 0 \\ 0 \\ -2 \end{bmatrix} \end{array}$$

$$0x = \begin{bmatrix} 0.99 \\ 0.01 \\ 0.01 \\ +1 \end{bmatrix}$$

$$\text{Mother} = \begin{bmatrix} 0.01 \\ 0.01 \\ 0.99 \\ -1 \end{bmatrix}$$

$$\begin{array}{l} \xrightarrow{\text{Mother}} \xrightarrow{\text{Father}} \\ = \begin{bmatrix} 0 \\ 0 \\ 0 \\ -2 \end{bmatrix} \end{array}$$

$$\text{Father} = \begin{bmatrix} 0.01 \\ 0.01 \\ 0.99 \\ +1 \end{bmatrix}$$

$$e_{\text{cow}} - e_{\text{son}} \approx e_{\text{mother}} - e_{\text{father}}$$

$$\approx e_{\text{daughter}} - e_{\text{son}}$$

word1

word2

wordn

`<unk>`



Embedding Matrix.

$$E \cdot \theta_{182} = e_{182}$$

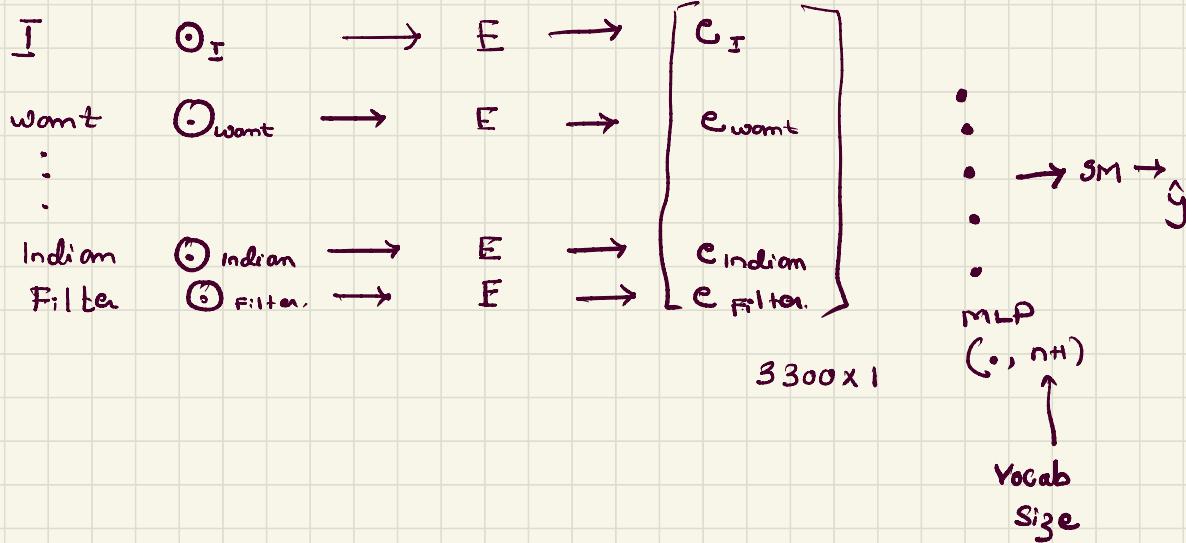
$$\begin{bmatrix} E \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} e_{182} \end{bmatrix}$$

Learning Word Embedding : Neural Language Models

: Next word prediction.

I want to drink a hot cup of south

Indian Filter



→ No need to take n vectors always
Fix a window previous n words

input $\in \mathbb{R}^{1200}$

if $\in \mathbb{R}^{1200}$ \vdots $\rightarrow SM \rightarrow \hat{y} \in \mathbb{R}^{n+1}$

→ Content → target

↳ words as i/p.

→ You can also modify the context

↳ words from right, 4 words from left.

→ Simple context : One word context.

$$\begin{bmatrix} 0.1 \\ 0.2 \\ 0.4 \\ 0.1 \\ 0.2 \end{bmatrix}$$

Word2Vec:

→ skip grams.

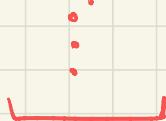
I want to drink a hot cup of South Indian
Filter Coffee, with Idli Sambhar.

→ Choose a word randomly as context.

→ Target is chosen in a window.

Context

Cup
Cup
Cup
.
:
:



↑
input

target

hot
drink
South
.
:
:



↑
o/p

model :

map Context \rightarrow target.
 "c" "t"

$O_c \rightarrow E \rightarrow E_c \rightarrow \text{MLP} \rightarrow \text{SM} \rightarrow \hat{y}$

$$\hat{y} = \text{IP}[t | c]$$

prob of target given the context.

this will be a prob vector of size (O+1)

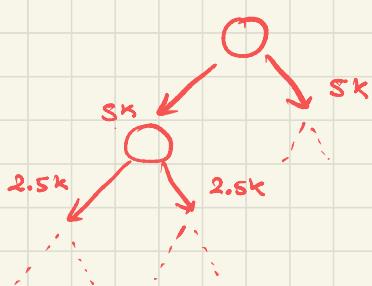
$$= \frac{\Theta_E^T \cdot e_c}{\sum_{j=1}^{n+1} \Theta_j^T \cdot e_c} \quad ||$$

$$L(\hat{y}, y) = - \sum_{i=1}^{n+1} y_i \log \hat{y}_i$$

CCE Loss.

Hierarchical Software

Vocab Size : 10000



binary classifier: Logistic regre

→ words of vocabulary is the leaf nodes.

it has a unique path from the root.

→ you will have around $\log_2 n$

→ The prob is Computed as a product of binary probabilities along the path from root to the word.

→ Skip gram $O(\text{size of Vocab})$

Hierarchical $O(\log_2(\text{Size of Vocab}))$

Glove algorithm:

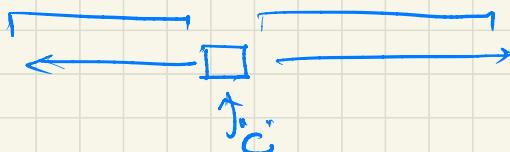
Global Vectors for word representations.

→ from the data come up with the co-occurrence matrix.

	is	the	apple	fruit	Computer
apple	5	20	100	50	5
fruit	3	15	50	100	1
Computer	1	10	5	1	100

→ window & then compute this.

x_{ij} : # times "i" appears in the context of "j"
 x_i : # times "i" appears in the context "C"



$$p_{ij} = \frac{x_{ij}}{x_i}$$

p_{ij} is greater in magnitude :
 $i \& j$ are strongly related.

* When we have target appearing ± 10 words from the target

$$x_{ij} = x_{ji}$$

but x_i may not be equal to x_j

$\therefore p_{ij}$ may not be equal to p_{ji}

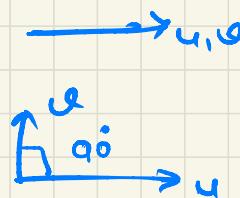
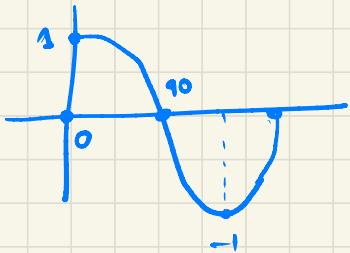
$$* \sum_{i=1}^{10k} \sum_{j=1}^{10k} \left(\underbrace{\Theta_i^\top e_j + b_i + b_j}_{\text{↑}} - \log(x_{ij}) \right)^2$$

how is the
similarity b/w
 $i \& j$

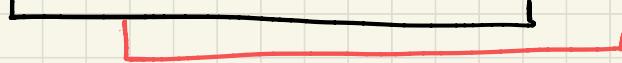


$$\text{Sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

$$\|u\|_2 = \sqrt{u_1^2 + u_2^2 + \dots + u_k^2}$$



index [\overline{s}_0 , \overline{s}_1 , \overline{s}_2 , \overline{s}_3 , \overline{s}_4 , \overline{s}_5 , ...]
word [\overline{st}_0 , \overline{st}_1 , \overline{st}_2 , \overline{st}_3 , \overline{st}_4 , \overline{st}_5 , \overline{st}_6 , ...]



1/10 [\overline{s}_0 , \overline{s}_1 , \overline{s}_2 , \overline{s}_3 , \overline{s}_4], [\overline{st}_1 , \overline{st}_2 , \overline{st}_3 , \overline{st}_4 , \overline{st}_5], ...

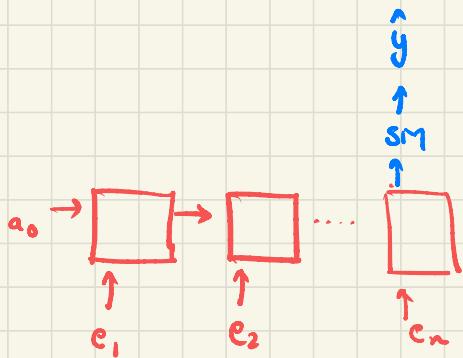
0/10 [[\overline{st}_5], [\overline{s}_6], ...]

Sentiment Analysis:

i/p : tweets

o/p :	Very offensive	- 0
	Offensive	- 1
	normal	- 2
	good	- 3
	very good.	- 4

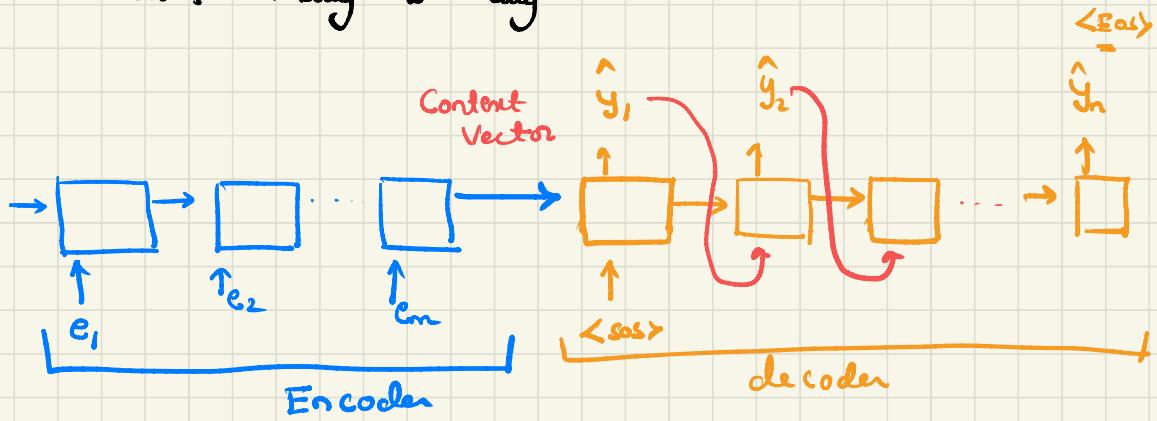
Arch : many to one.



Google Translate: Sequence to Sequence Models.

French to English.

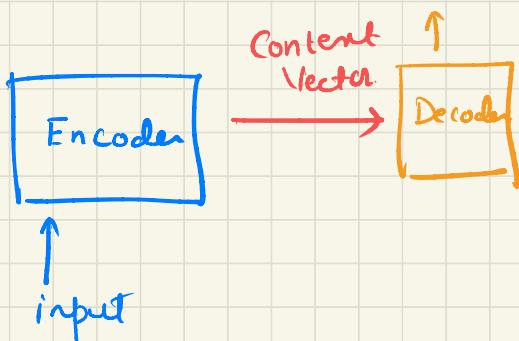
Arch : Many to Many



$\langle \text{sos} \rangle$: Start of Sentence.

$\langle \text{EOS} \rangle$: End of Sentence

Output :



BLEU : Bilingual Evaluation Understudy

i/p : (मैं) एक बड़ी बाला हूँ

Machine Translation : The cat is on the mat
(Candidate)

(Reference) Human Translation :

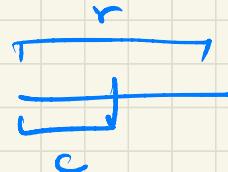
There is a cat on the mat

Candidate : [the, cat, is, on, the, mat]

Reference : [There, is, a, cat, on, the, mat]

Candidate unigrams :

	Candidate	Reference
the	2	1
cat	1	1
is	1	1
on	1	1
mat	1	1



$$\phi_1 = \frac{\text{matched unigrams}}{\text{Total candidate unigrams}} = \frac{5}{6} \approx 0.83$$

Candidate Bi-grams.

	Candidate	Reference
the cat	1	0
cat is	1	0
is on	1	0
on the	1	1
the mat	1	1

$$P_2 = \frac{\text{Matched bigrams}}{\text{Total Candidate bigrams}} = \frac{2}{5} = 0.4$$

* Decide on how many n-grams you need to take
 normally people go till P_4

$$\frac{1}{4} \log P_1 + \frac{1}{4} \log P_2 + \frac{1}{4} \log P_3 + \frac{1}{4} \log P_4$$

in our case we are stopping C bi gram

$$\frac{1}{2} \log 0.83 + \frac{1}{2} \log 0.4$$

$$= 3T$$

Brevity Penalty (BP)

Candidate length (c): 6

reference length (r): 7

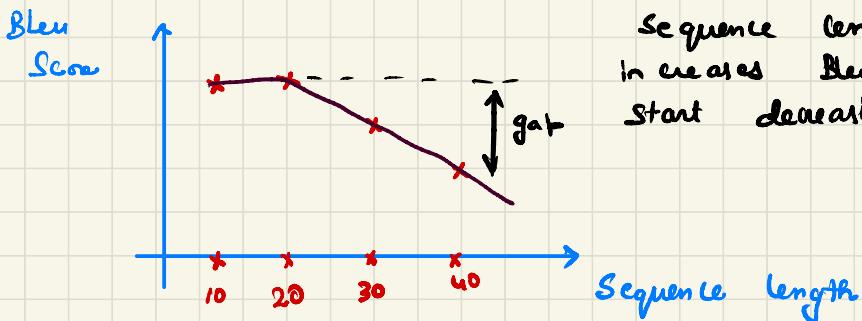
$$c > r : \text{BP} = 1$$

else:

$$\exp \left\{ 1 - \frac{r}{c} \right\}$$
$$= \exp \left\{ 1 - \frac{7}{6} \right\} = \frac{1}{6}^{-1} = \text{BP}$$

Bleu Score: $\text{BP} * \overline{\text{F}}$

$$\text{BP} * \frac{1}{2} \sum_{i=1}^2 \log p_i$$



Problem : Machine Translation.

x_1, \dots, x_m input sequence.

$\hat{y}_1, \dots, \hat{y}_n$ output sequence.

$n \neq m$

$\alpha_{1,1}$: how much attention should we give
for 1st i/p word to generate 1st o/p word.

$\alpha_{2,1}$:
for 1st i/p word to generate 2nd o/p word.

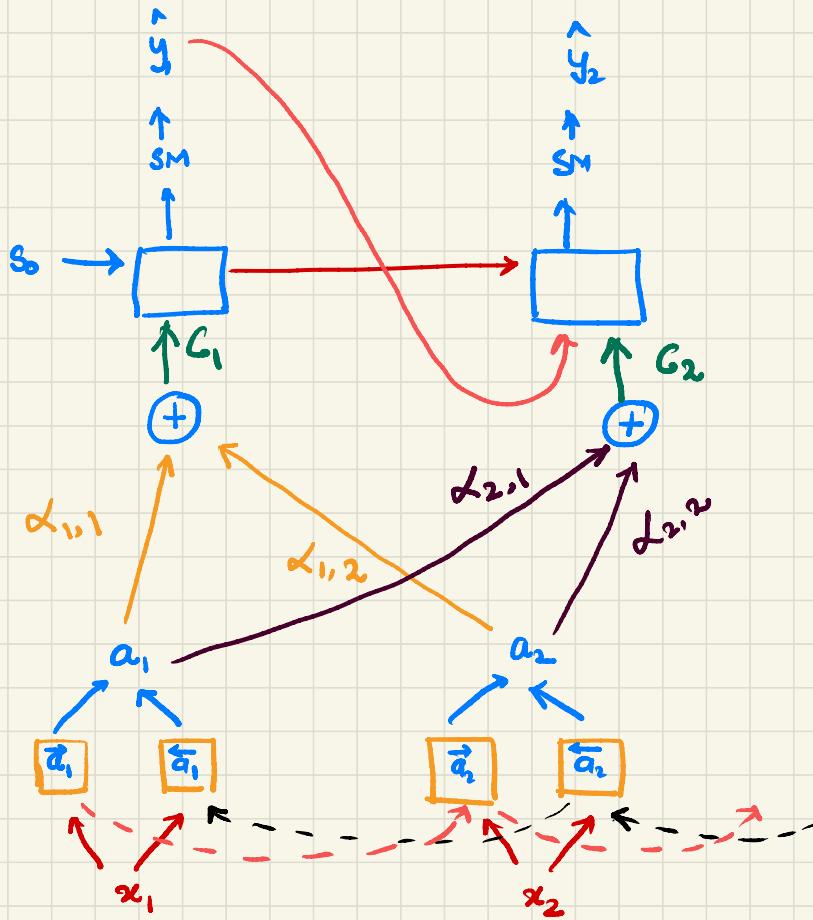
$\alpha_{i,j}$
to generate i/p word.

$a_t = (\vec{a}_t, \overleftarrow{a}_t)$

input length : T_{2x}

$$\sum_{t=1}^{T_{2x}} \alpha_{1,t} = 1$$

$$G_1 = \sum_{t=1}^{T_{2x}} \alpha_{1,t} a_t \\ = \alpha_{1,1} * a_1 + \alpha_{1,2} * a_2 \dots$$



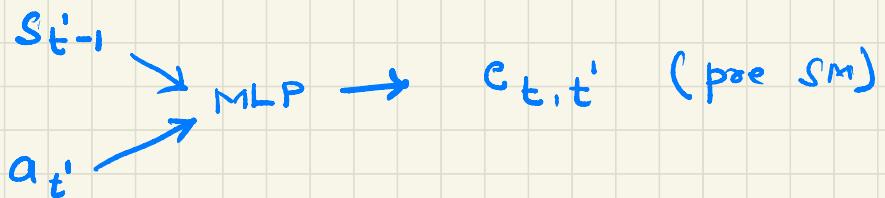
$\alpha_{1,2}$: amount of attention \hat{y}_1 should pay
for x_2

$\alpha_{1,2}$:

a_2

s_1

$\alpha_{t,t'}$: amount of attention \hat{y}_t should
pay to the input $x_{t'}$



$$\alpha_{t,t'} = \frac{\exp \{ c_{t,t'} \}}{\sum_{t'=1}^{T_x} \exp \{ c_{t,t'} \}}$$

Transformer flow architecture:

$x_1, x_2, x_3, x_4, x_5 \rightarrow$ input in word embedding.

$A_1, A_2, A_3, A_4, A_5 \rightarrow$ Attention.



		Query (Q)	key (k)	Value (v)
x_1	A_1	q_1	k_1	v_1
x_2	A_2	q_2	k_2	v_2
x_3	A_3	q_3	k_3	v_3
x_4	A_4	q_4	k_4	v_4
x_5	A_5	q_5	k_5	v_5

q_i is like a question

$q_3 \cdot k_1$: how good is 1st word for 3rd word's question.

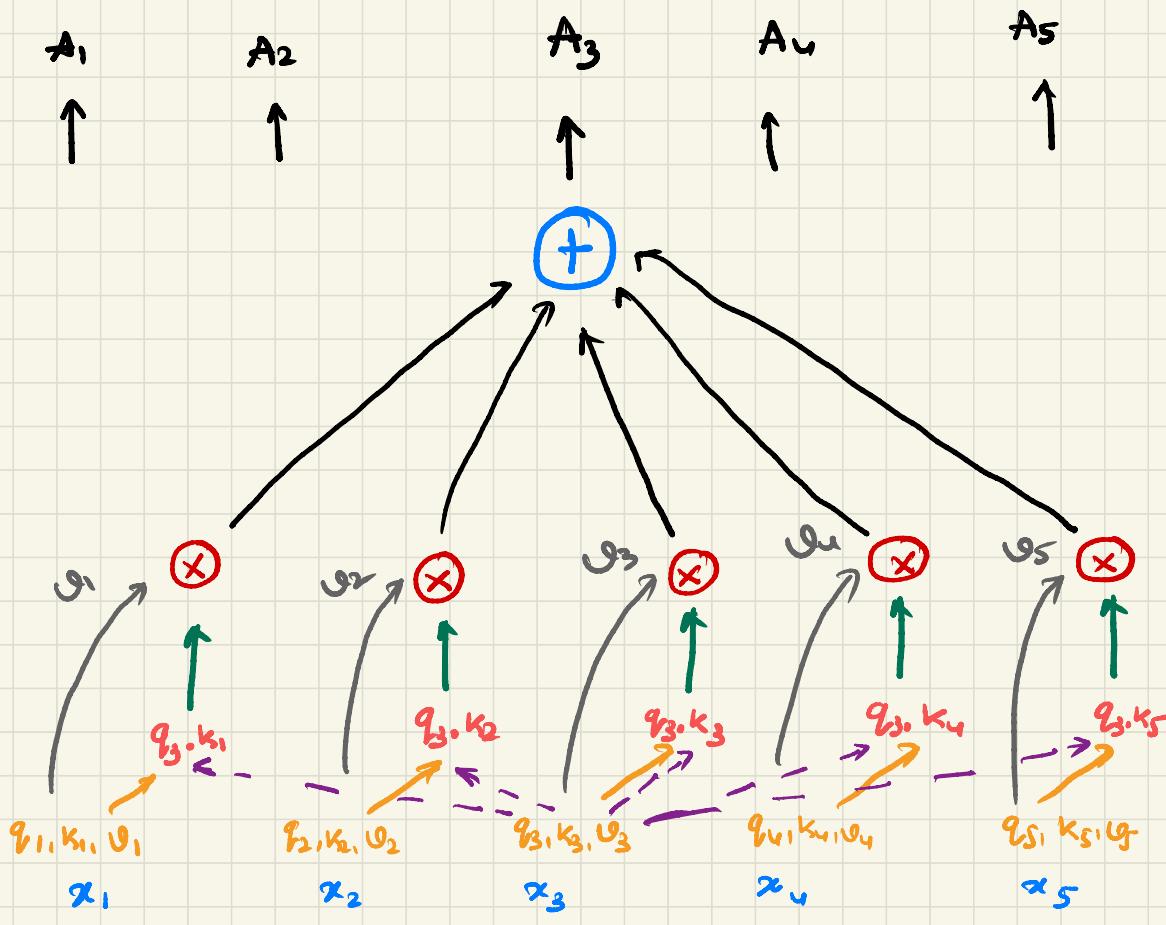
$q_i \cdot k_j$: how good is the jth word for ith question.

→ the aim of this operation is to get more info needed to help to get most useful info

$$A(q, k, v) = \sum_i \frac{\exp\{q \cdot k_i\}}{\sum_j \exp\{q \cdot k_j\}} \cdot v_i$$

Softmax

$$q_1 \Rightarrow \frac{\exp\{q_{11}, k_{11}\}}{\sum_j \exp\{q_{1j}, k_{1j}\}} \cdot v_1 + \frac{\exp\{q_{11}, k_{12}\}}{\sum_j \exp\{q_{1j}, k_{1j}\}} \cdot v_2 + \dots$$



$$q_i \in \mathbb{R}^{d_k}$$

$$k_i \in \mathbb{R}^{d_k}$$

$$v_i \in \mathbb{R}^{d_v}$$

$$x \in \mathbb{R}^{n \times d_{\text{model}}}$$

n : num of example
 d_{model} : dim of embedding

$$Q = X W_Q$$

$$(4 \times 70) (70 \times 50)$$

$$V = X W_V$$

$$(4 \times 70) (70 \times 60)$$

$$K = X W_K$$

$$(4 \times 70) (70 \times 50)$$

$$d_K = 50$$

$$d_V = 60$$

$$n = 4$$

$$d_{\text{model}} = 70$$

$$x \in \mathbb{R}^{4 \times 70}$$

$$\begin{bmatrix} -x_1- \\ -x_2- \\ -x_3- \\ -x_4- \end{bmatrix}_{4 \times 70} = X$$

$$\begin{bmatrix} -q_1- \\ -q_2- \\ -q_3- \\ -q_4- \end{bmatrix}_{4 \times 50} = Q$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ K_1 & K_2 & K_3 & K_4 \\ 1 & 1 & 1 & 1 \end{bmatrix}^T = K$$

$$\begin{bmatrix} -q_1- \\ -q_2- \\ -q_3- \\ -q_4- \end{bmatrix} \quad \begin{bmatrix} 1 & 1 & 1 & 1 \\ k_1 & k_2 & k_3 & k_4 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} q_1 k_1 & q_1 k_2 & q_1 k_3 & q_1 k_4 \\ q_2 k_1 & q_2 k_2 & q_2 k_3 & q_2 k_4 \\ q_3 k_1 & q_3 k_2 & q_3 k_3 & q_3 k_4 \\ q_4 k_1 & q_4 k_2 & q_4 k_3 & q_4 k_4 \end{bmatrix} = Q \cdot K^T$$

Attention (Q, K, V) = $\text{Softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$.

Scaled dot Product attention.