

3<sup>rd</sup> March 2025:

Sequences:

MLP      Fixed length data  
CNN      3

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

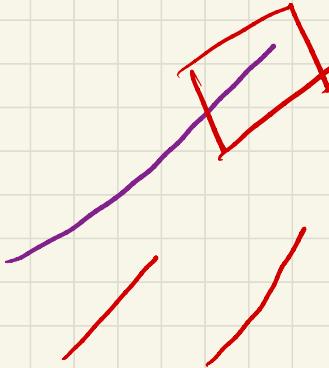
→ Images of Varying resolution.

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	?

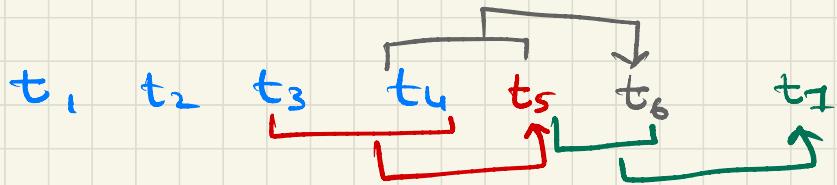


how much to go back?

taking Only the past 2 days.

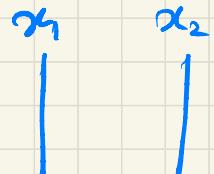
$$t_5 = b + w_1 t_4 + w_2 t_3$$

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

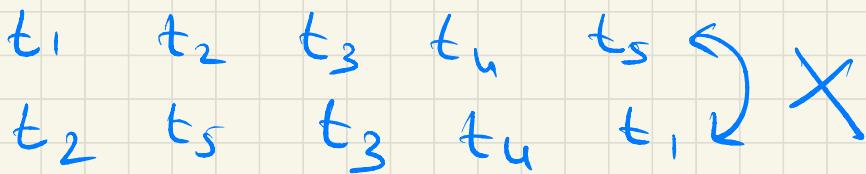


Auto      Regressive models.

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"

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

$p=1 :$

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

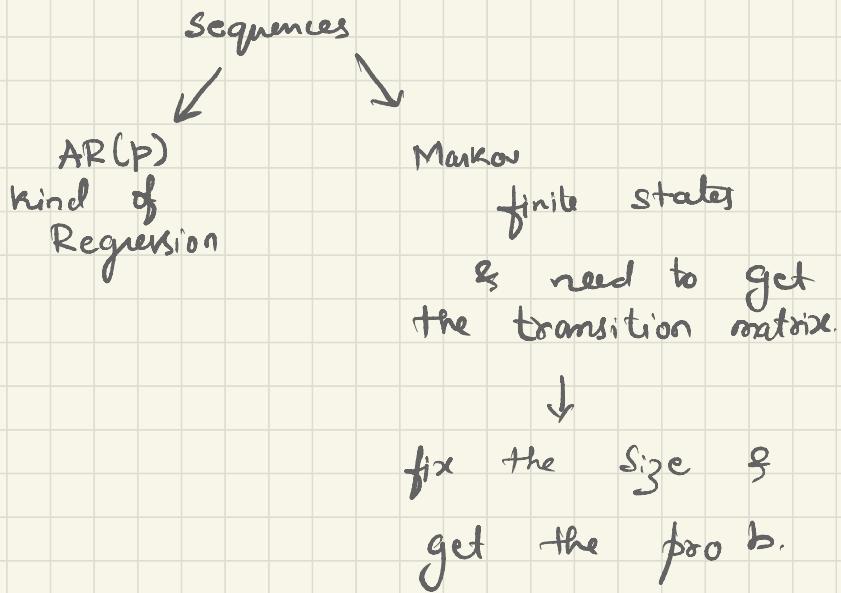
$$x_2 = b + w_1 x_1$$

$$x_3 = b + w_1 x_2$$

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.

6<sup>th</sup> Mar 2025:

- ✓ fixed length i/p → fixed length o/p  
→ variable length i/p → fixed length o/p

### Named Entity Representation: NER

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

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_6$

$x_t$ :  $t^{\text{th}}$  feature

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

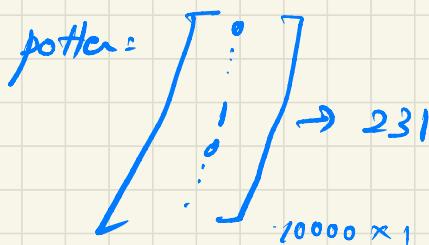
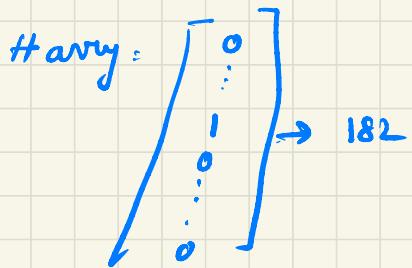
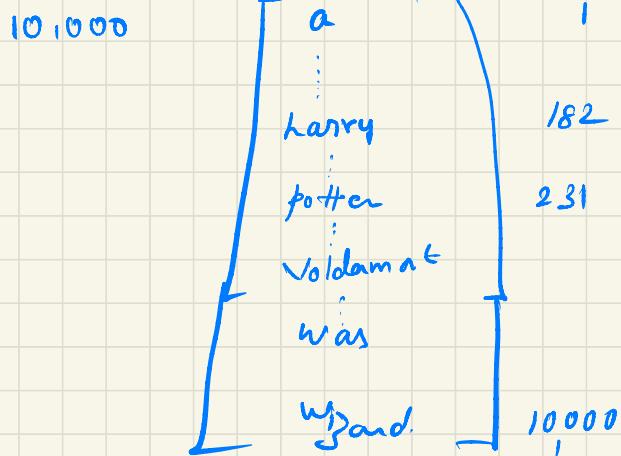
In this specific case of NER.

$y_1$ ,  $y_2$  ...,  $y_6$

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

Construct Vocabulary:

num of unique tokens (Size of vocabulary) is

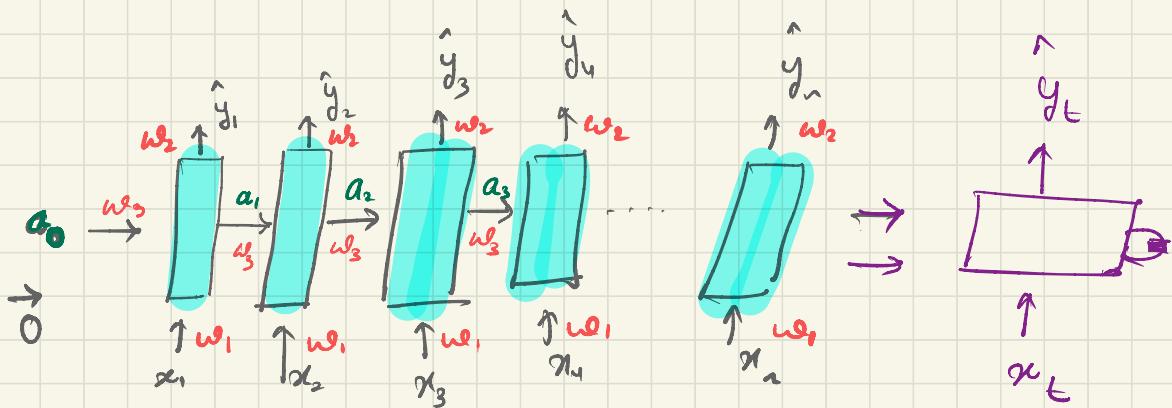


model to map from  $X \rightarrow Y$ .

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.

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

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

10<sup>th</sup> 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.

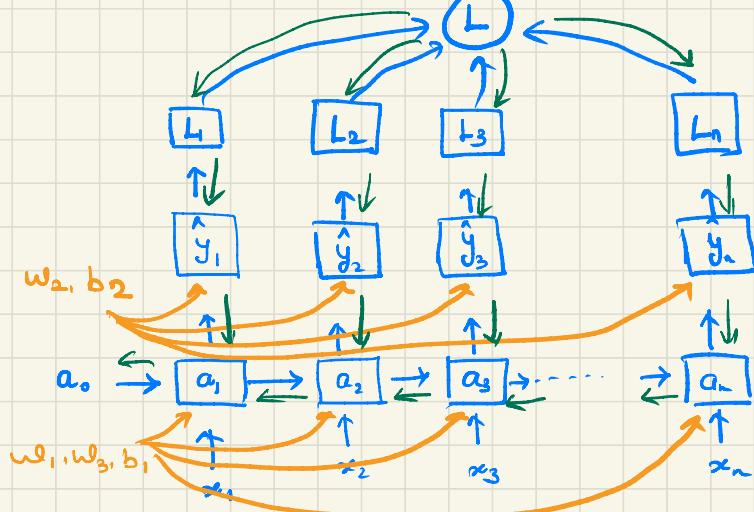
0      0      0 0      0 0 1

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 :

$$L_1(\hat{y}_i, y_i) = -y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$$

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

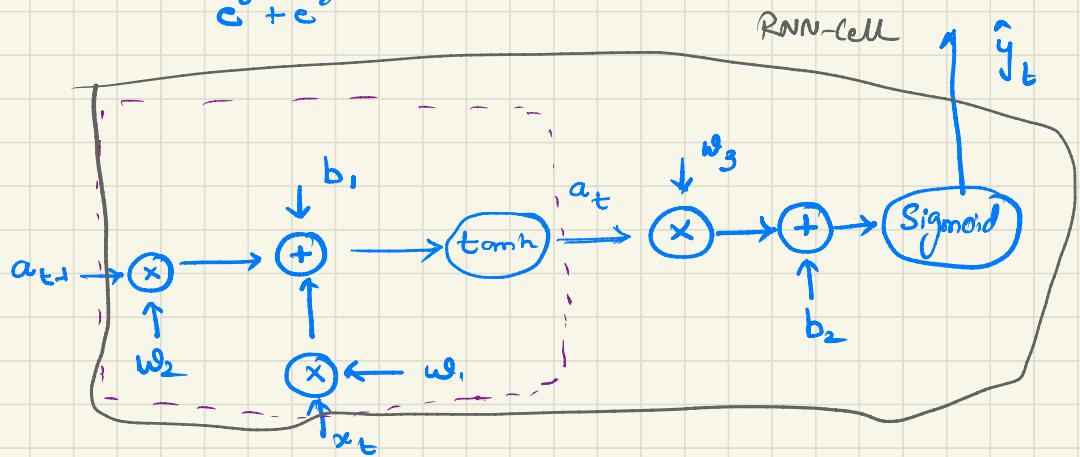
$$\begin{bmatrix} \hat{y}_i \\ 1-\hat{y}_i \end{bmatrix} \begin{bmatrix} \hat{y}_i \\ 1-\hat{y}_i \end{bmatrix}$$

11<sup>th</sup> March 2025:

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

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

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



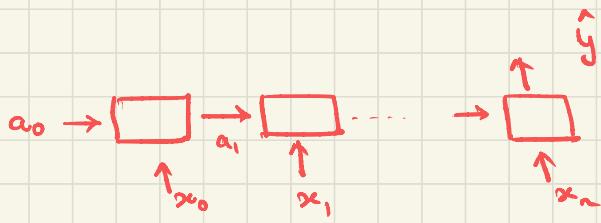
NER model we are looking is a Many to Many model.  
& size of inputs & output were same.

→ Movie Ratings

input: Review

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

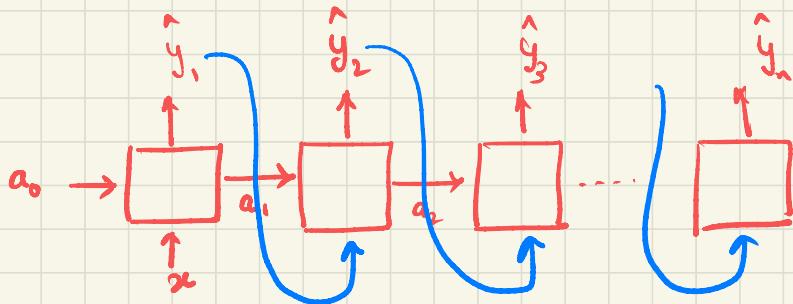
Many to One Mapping.



→ given the base node, play some Raga.



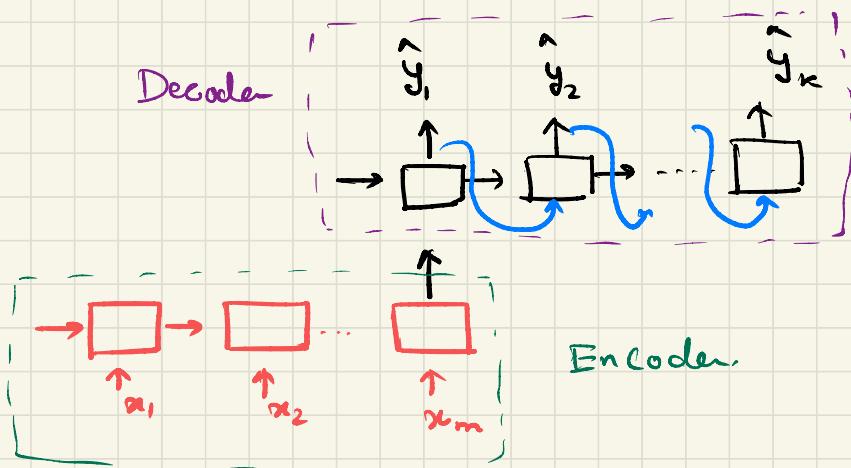
One to Many Problem.



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.

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

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

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.

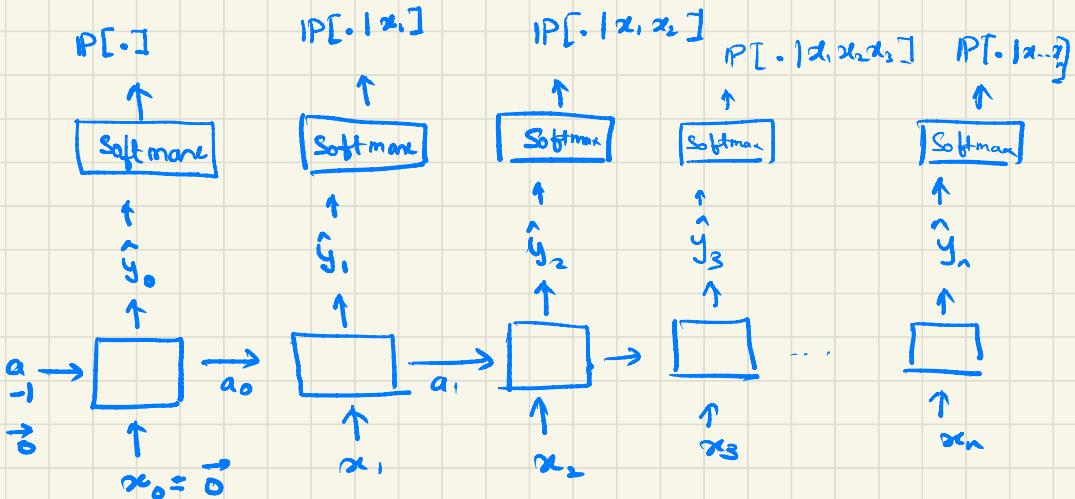
Output should choose one word out of all the words in the vocabulary.

$\therefore$  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 —

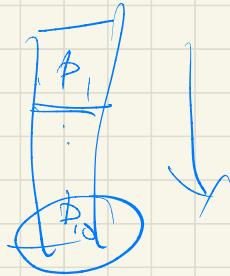
The orange —

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

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

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

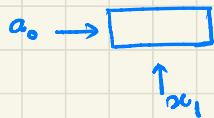
$$p_3 = \frac{e^{q_3}}{e^{a_1} + e^{q_2} + e^{q_3}}$$



12<sup>th</sup> 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  $5000 \times 5$

input @ every time step :

$5000 \times 5$

→ 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  $231 \times 5$

∴ 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)

