

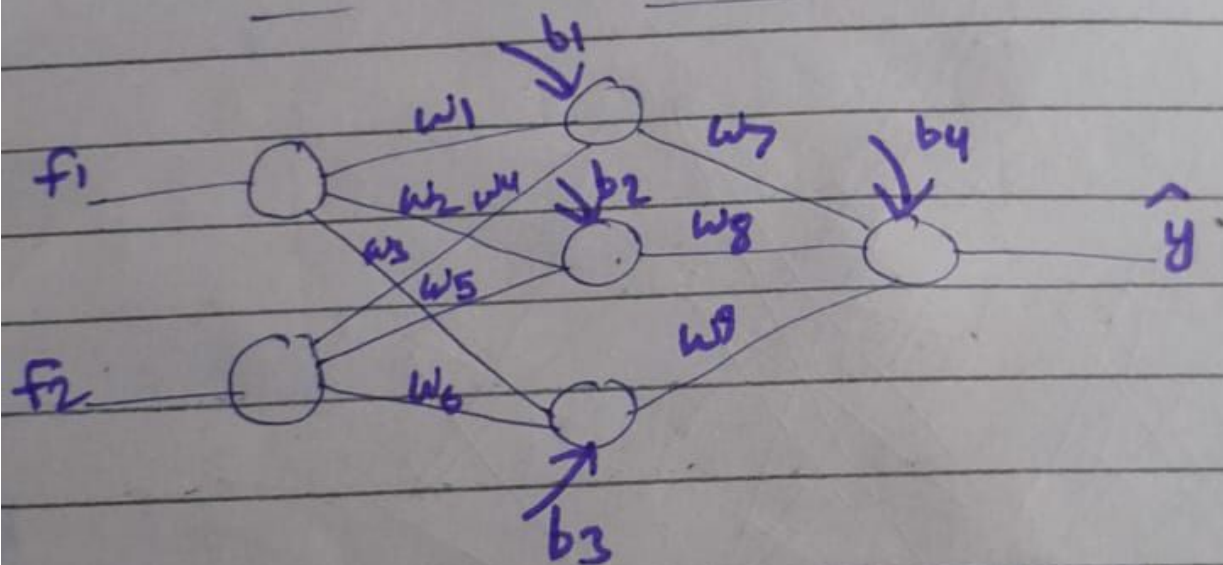
⇒ NLP in Deep Learning:

① ANN → Artificial Neural Network

Tabular data $\begin{cases} \rightarrow \text{Classification} \\ \rightarrow \text{Regression} \end{cases}$

e.g. House Price Prediction

f_1	f_2	y
House Size	No. of Rooms	Price



Here sequence of data doesn't matter we can interchange f_1 and f_2 . This is called Non-Sequential data.

② CNN (Convolutional Neural Network)
↳ Images → Image classification/
object detection

③ RNN → data (sequential data)

↓
Sequence/order is
important (changing
sequence change meaning)
Input

e.g (sequential data)

(i) Text Generation → This is an apple
↳ juice ^{output}

(ii) Chatbot conversation → Q and A
(Sequential Data → Sequence is also
important here)

ii) Language Translation

↳ [English] → [French]

(Here also sequence is important)

i) Autogeneration

Sales Data → Time Series (Datatime)

Sales Forecasting → Sequence matters

Can we solve with ANN?
↳ Sequential Data
e.g. Sentiment Analysis

Text	Output
1) The food is good	1
2) The food is bad	0
3) The food is not good	0

Step 01: Text → Vectors

We have 4 unique words after ignoring stop words ⇒ Vocabulary size = 4

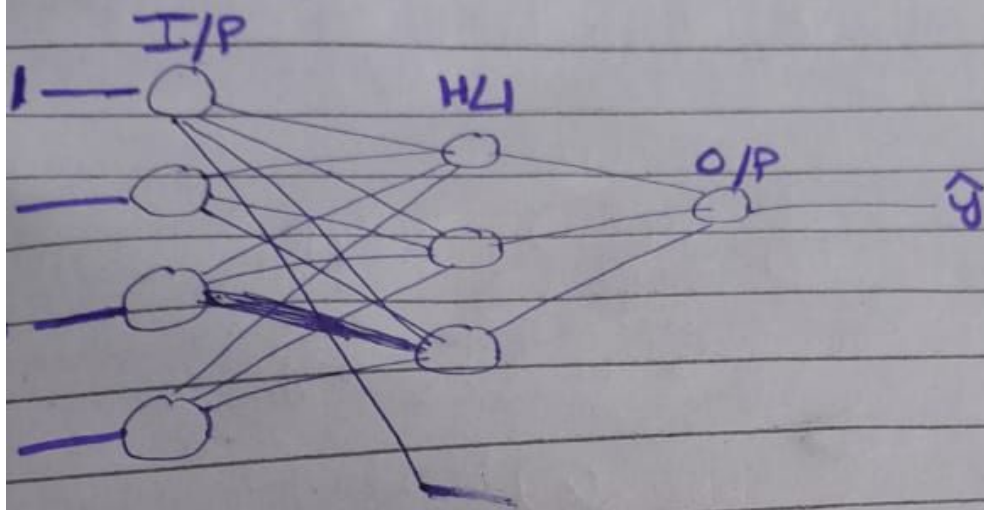
Bag of Words (BOW)

	food	good	bad	not
1	1	1	0	0
2	1	0	1	0
3	1	1	0	1

Text Data → Sequence Information is important

Here in ANN, we are losing

Sequence Information \rightarrow Meaning/context is losing.



Text Data:

Sequence of Information

"The food is good"

Here we are giving whole sentence at once in ANN so sequence lost.

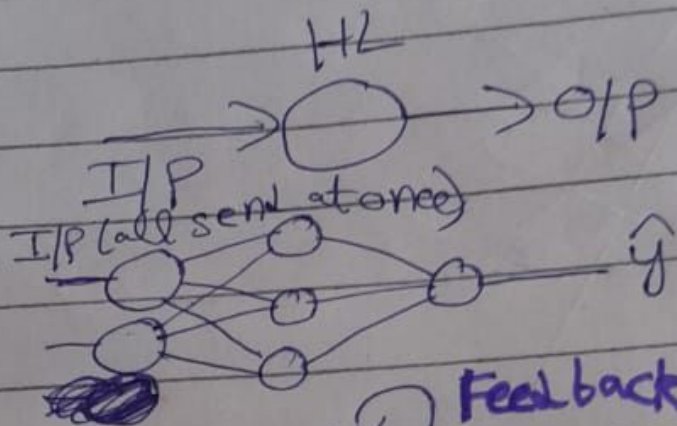
The best practice is that we should give word by word one word at a time to maintain sequence and context.

That's why to solve this kind of problem we use

Simple RNN (Recurrent Neural Network)
↳ Data is sequential

ANN:

e.g



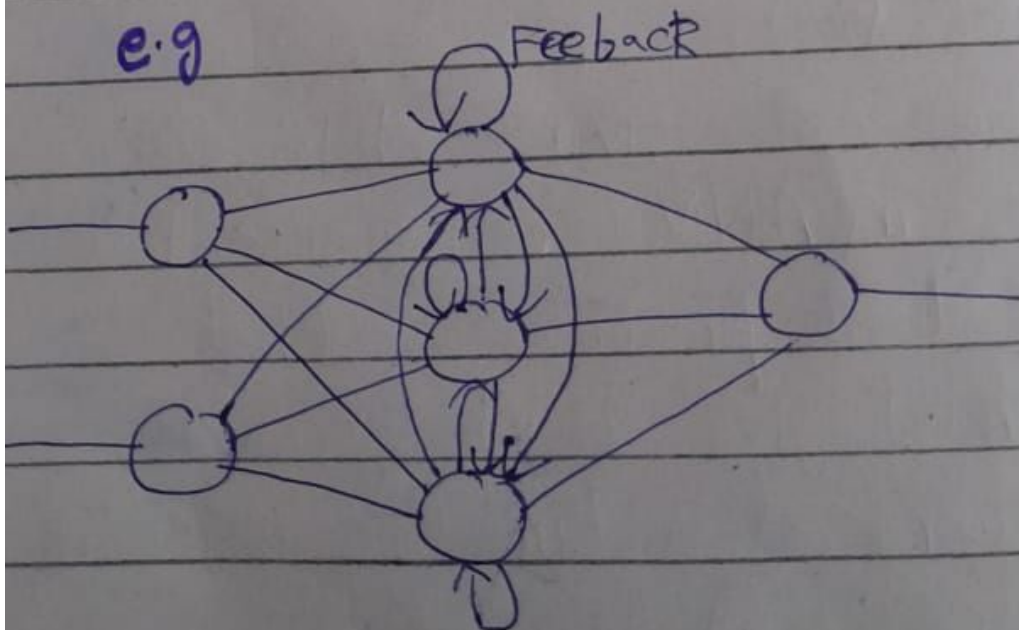
RNN:

I/P

Feedback

O/P

e.g



- Actually in RNN, a neuron shares its output or feedback to itself

and all other neurons in the same hidden layer.

eg

x_{11} x_{12} x_{13} x_{14}

The food is good

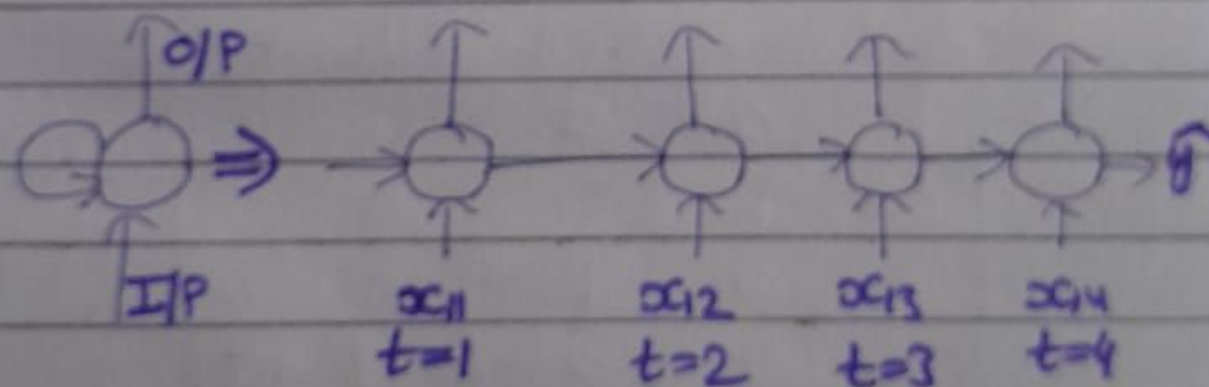
- In RNN, we pass each word at a particular time stamp (one word at a time)

At $t=1$ x_{11}

$t=2$ x_{12}

$t=3$ x_{13}

$t=4$ x_{14}



⇒ Working of Simple RNN with Forward Propagation:

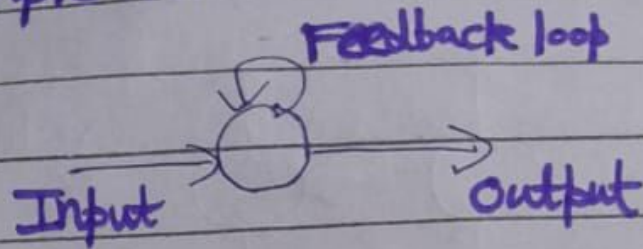
Dataset

Text	Output
^{x_{t1}} The ^{x_{t2}} food ^{x_{t3}} is <u>good</u>	1
The food is <u>bad</u>	0
The food is <u>not</u> good	0

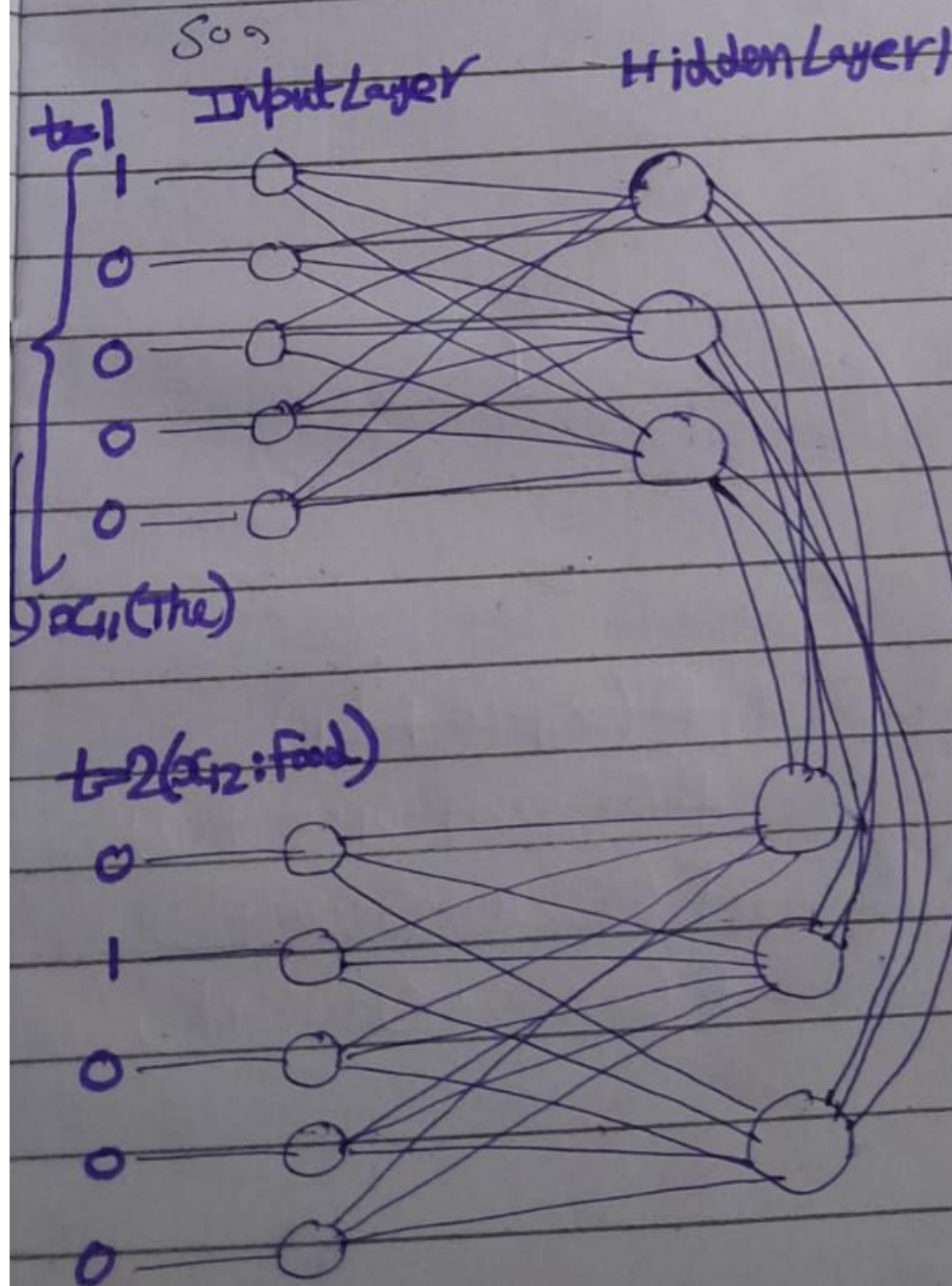
Words → Vectors (using one Hot encoding)

	The	food	good	bad	not
$The(x_{t1}) ⇒$	[1	0	0	0	0]
$ood(x_{t2}) ⇒$	[0	1	0	0	0]
$od(x_{t3}) ⇒$	[0	0	1	0	0]
$ad(x_{t4}) ⇒$	[0	0	0	1	0]
$t(x_{t5}) ⇒$	[0	0	0	0	1]

• Simple RNN Architecture

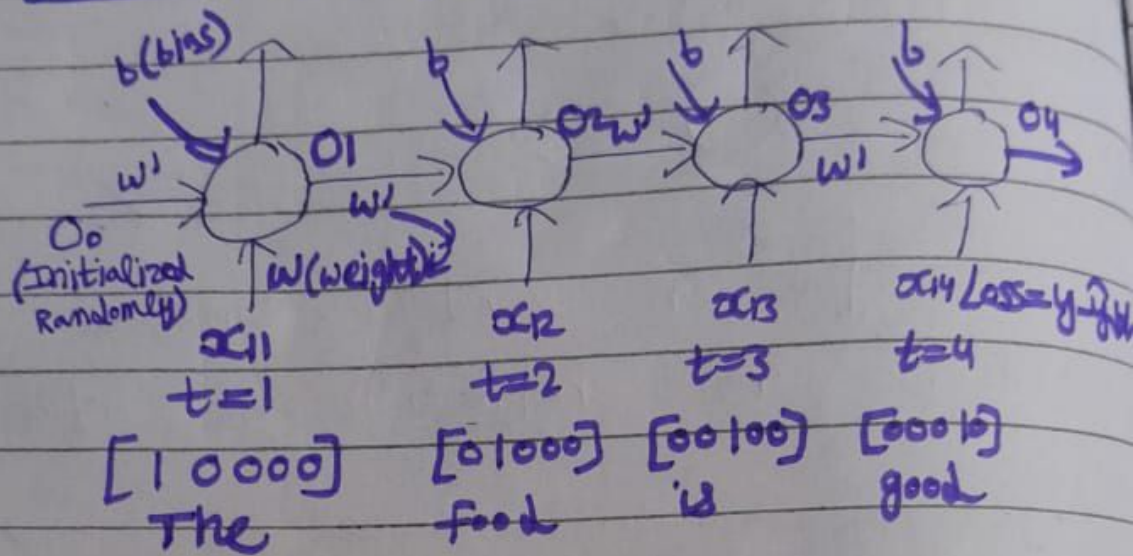


As we pass one word at a time
and one word is representing by 5 inputs



RNN
(You can see
at $t=2$ when
food is passed
the Hidden Layer
also has context
of previous
word the)

Forward Propagation with Time:



Dataset

Text

The food is good

$x_{11} \ x_{12} \ x_{13} \ x_{14}$

Output

1

Forward Propagation:

$$O_1 = f(x_{11} \cdot w + b_1)$$

$$O_2 = f(x_{12} \cdot w + O_1 \cdot w' + b_2)$$

$$O_3 = f(x_{13} \cdot w + O_2 \cdot w' + b_3)$$

$$O_4 = f(x_{14} \cdot w + O_3 \cdot w' + b_4)$$

$f \Rightarrow$ Activation function

RNN Back Propagation with time:

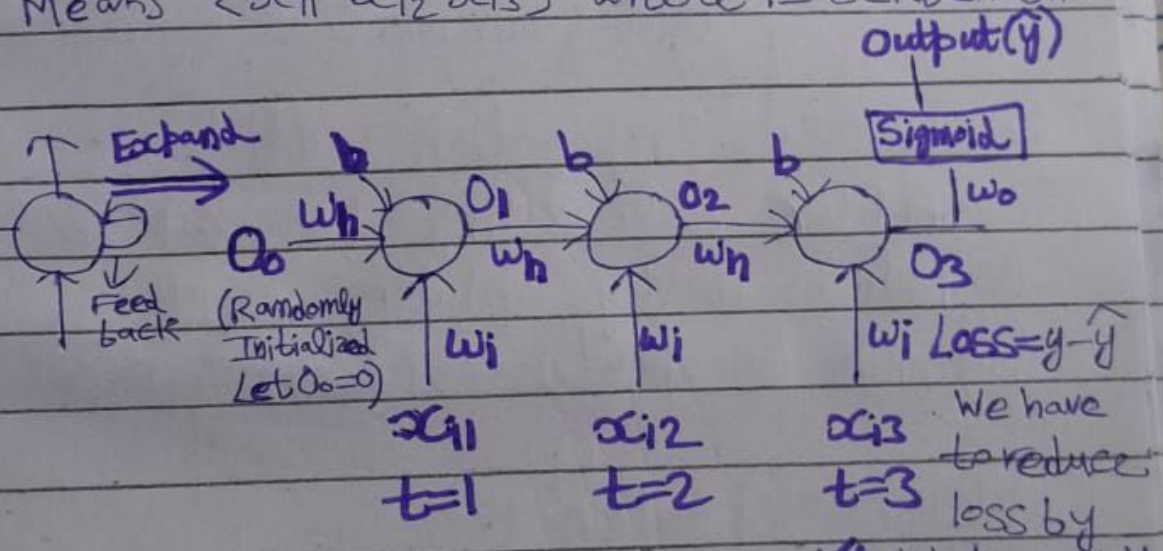
Let we have sentences

$S_1 \langle x_{11} \ x_{12} \ x_{13} \rangle$

$S_2 \langle x_{21} \ x_{22} \ x_{23} \rangle$

$S_3 \langle x_{31} \ x_{32} \ x_{33} \rangle$

Means $\langle x_{i1} \ x_{i2} \ x_{i3} \rangle$ where $i = \text{Sentence no.}$



Forward Propagation:

→ Activation function

$$O_1 = f(x_{i1} * w_i + O_0 * w_h + b)$$

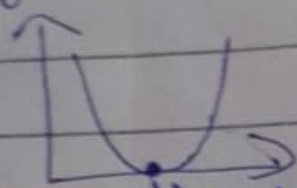
$$O_2 = f(x_{i2} * w_i + O_1 * w_h + b)$$

$$O_3 = f(x_{i3} * w_i + O_2 * w_h + b)$$

$$\hat{y} = \sigma(O_3 * w_o)$$

Back Propagation
updating weights
[w_i, w_h, w_o]

We will repeat forward and back propagation until convergence happens means until we reach global minima in gradient descent



Global minima (Loss is minimum)

Backward Propagation with time:

After calculating loss we have to reduce it by updating weights
Weights to update are: $[w_i, w_h, w_o]$

Weight Update Formula:

$$w_{\text{new}} = w_{\text{old}} - \eta \left[\frac{\partial L}{\partial w_{\text{old}}} \right]$$

Learning rate η Derivative / Slope of Gradient Descent

① Updating w_o :

$$w_{o_{\text{new}}} = w_{o_{\text{old}}} - \eta \frac{\partial L}{\partial w_{o_{\text{old}}}}$$

Chain rule of derivative

$$\boxed{\frac{\partial L}{\partial w_{old}} = \frac{\partial L}{\partial y} * \frac{\partial y}{\partial w_{old}}}$$

② Update w_h [Hidden Layer weights]

Time Stamps

$$w_{h_{new}} = w_{hold} - \eta \frac{\partial L}{\partial w_{hold}}$$

We have 3 time stamps $t=1, 2, 3$

$t=3$

$$\frac{\partial L}{\partial w_{hold}} = \left[\frac{\partial L}{\partial y} * \frac{\partial y}{\partial o_3} * \frac{\partial o_3}{\partial w_h} \right]$$

$t=2 \Rightarrow$

$$+ \left[\frac{\partial L}{\partial y} * \frac{\partial y}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial w_h} \right]$$

$t=1 \Rightarrow$

$$+ \left[\frac{\partial L}{\partial y} * \frac{\partial y}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_h} \right]$$

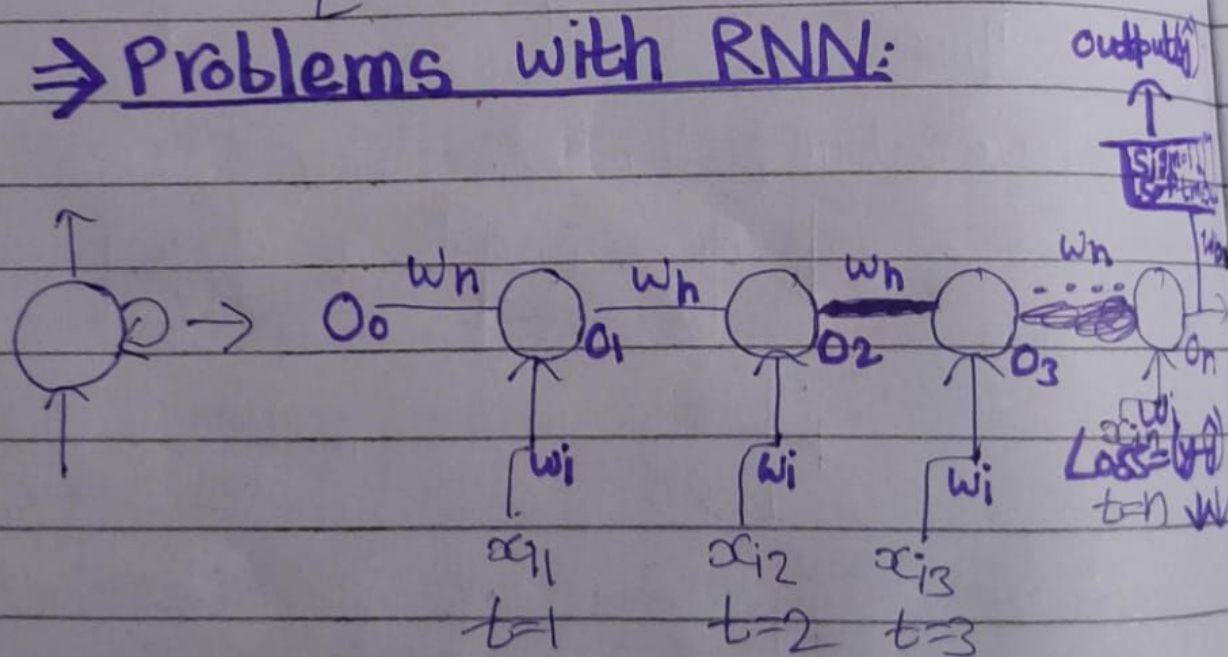
Updating weights w_i \rightarrow Time stamp

$$w_{i_{new}} = w_{i_{old}} - \eta \frac{\partial L}{\partial w_{i_{old}}}$$

We have 3 timestamps $t=1,2,3$

$$\begin{aligned} \frac{\partial L}{\partial w_{iold}} &= \left[\frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial w_{iold}} \right] + \\ t=2 &\Leftarrow \left[\frac{\partial L}{\partial \hat{y}} + \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial w_{iold}} \right] + \\ t=3 &\Leftarrow \left[\frac{\partial L}{\partial \hat{y}} + \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial o_2} + \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{iold}} \right] \end{aligned} \Rightarrow t=1$$

⇒ Problems with RNN:



ANN → Vanishing Gradient Problem

Let we have dataset:

Text

Output

The food is good

The food is bad

1

0

} Small sentences

Let we have to do Text Generation
I like to play —

- * Now the word which is to be generated/predicted can be dependant on any word of sentence.

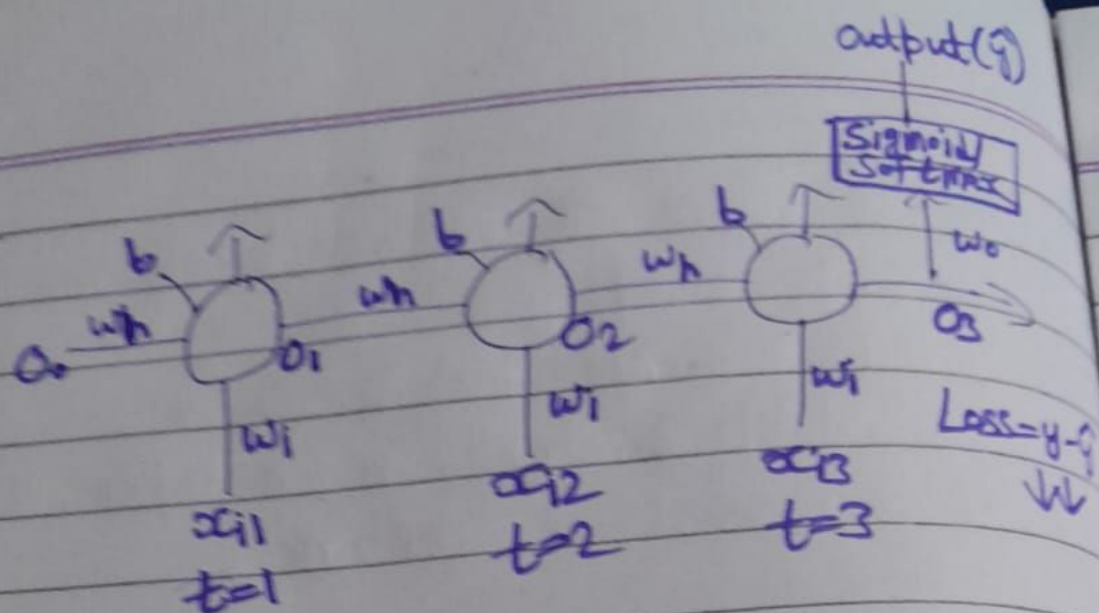
Now let we have a Long Sentence

"My name is Amir, and I like cricket, football and also like to make —"

- Now the dependency is huge as word to be predicted can be dependant on any word

Issue:

Long Term dependency can't be captured by simple RNN.



Backward propagation

Update weights: w_h, w_i, w_o

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial y} * \frac{\partial \hat{y}}{\partial w_{\text{old}}} \Rightarrow t=2$$

$$\frac{\partial L}{\partial w_{\text{old}}} = \left[\frac{\partial L}{\partial y} * \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial w_{\text{old}}} \right] \Rightarrow t=1$$

You can see that the length of word is just 3 and we got so lengthy chain rule

What if length of words = 50

At $t=2$

$$\frac{\partial L}{\partial w_{hot}} = \left[\frac{\partial L}{\partial y} * \frac{\partial \hat{y}}{\partial o_2} * \frac{\partial o_2}{\partial o_4} * \frac{\partial o_4}{\partial o_18} * \frac{\partial o_18}{\partial w_{hot}} \right] + \dots$$

$$\frac{\partial o_3}{\partial o_2} = \frac{\partial \sigma(x_{i3} * w_i + o_2 * w_h + b)}{\partial o_2} \ll o_3$$

We know

Derivative (sigmoid) = $[0 \text{ to } 0.25]$

When we keep on multiplying small values it will become ≈ 0

$t=1 \Rightarrow$ The word is not participating to update the weights value.

\Downarrow
Vanishing Gradient problem

\Rightarrow So the last words in long sentences are effective but the first words are not participating and may be the word to predict depend on first words.

⇒ So Simple RNN is unable to capture Long term dependency (when chain rule is very long)

To solve this problem we can use

- ① Relu, Leaky relu as activation functions
- ② LSTM RNN → Long Short Term Memory RNN
- ③ GRU RNN