

NLP In Deep Learning

Sunday, 29 December 2024 9:10 PM

NLP in Deep Learning:-

Text Data \rightarrow Vectors \rightarrow Numerical Rep.

- ① OHE - one hot encoding
- ② BOW
- ③ TF-IDF
- ④ word2vec, Avg word2vec
e.g. sentiment Analysis, Text classification

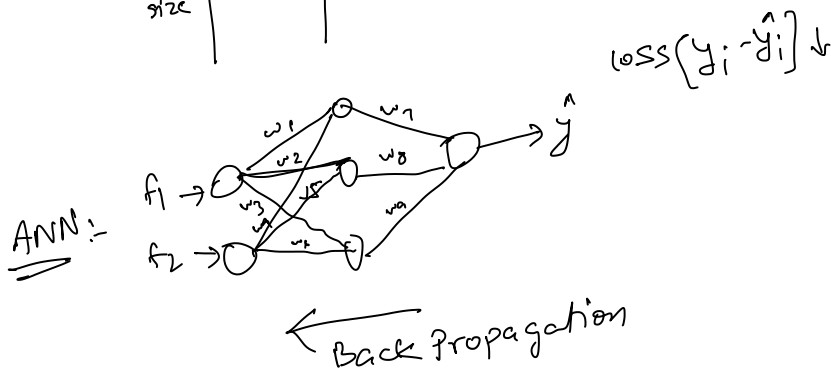
ANN \rightarrow Artificial Neural Network

↓

Classification Regression

Eg. House Price Prediction

f_1 size	f_2 rooms	y Price
---------------	----------------	--------------



② CNN \rightarrow convolutional Neural Network

\hookrightarrow Image classification.

Data - Image, Video frames.

③ Data - Sequential Data

① text generation

eg. I/P This is a Apple O/P juice
sequence

eg. Chatbot conversion Q/A
 I/P Question \rightarrow o/p. Answer
 O/P Answer

you...

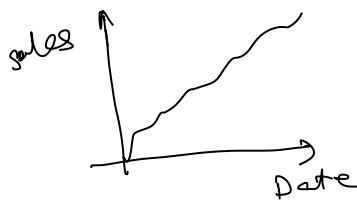
Sequential Data \rightarrow The food is good.

Eg Language Translation:



Eg. Auto suggestion
LinkedIn, email.

Eg. Sales Data.



Sales forecasting
on future date

Can we use ANN to solve this problem?
 \rightarrow which has sequential data.

NLP: Generative AI \rightarrow LLM, Multi Model.

① Simple RNN \Rightarrow LSTM/GRU RNN

\downarrow
Bidirectional RNN

\downarrow
Encoder Decoder

\downarrow
 \leftarrow Transformers \leftarrow Self Attention

Can we use ANN \rightarrow sequential data.

Dataset:- Sentiment Analysis

Text	O/p
the food is good	1
the food is bad	0
The food is not good	0

① Text preprocessing \rightarrow Text \rightarrow Vectors

Vocab: $\rightarrow 4$

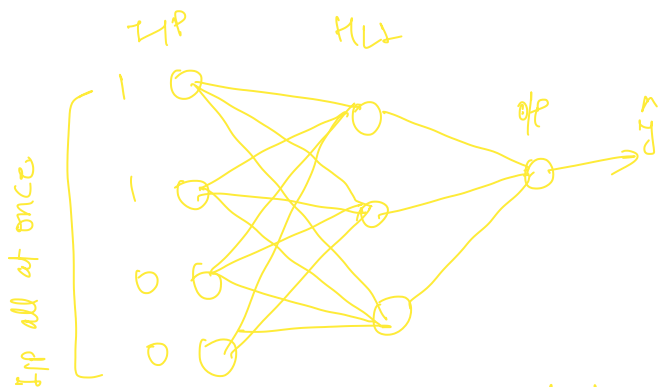
BOW: convert every sentence we convert to vectors.

	food	good	bad	not
S1	1	1	0	0
S2	1	0	1	0
S3	1	1	0	1

} I/P.

Text Data - Sequence Information is important

meaning of sentence is lost.



In case of ANN the data input is all at once, so we are losing the sequence information.

In case of Text Data:- we need to give one word at a time then backpropagation should happen & then another word.

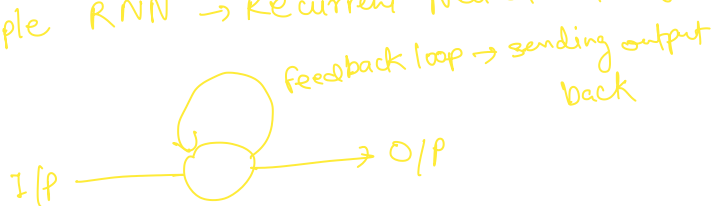
This food is good
t=1 t=2 t=3 t=4

Google Translation:-

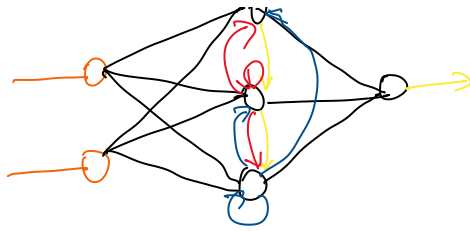
word to word translation takes place

There is another neural network

Simple RNN \rightarrow Recurrent Neural Network



n



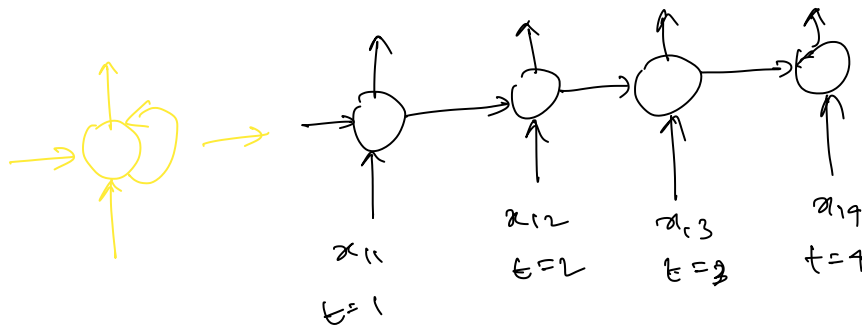
x_{11} x_{12} x_{13} x_{14}
the food is good

$t=1$ $x_{11} \rightarrow$ o/p

$t=2$ x_{12}

$t=3$ x_{13}

$t=4$ x_{14}



Working of Simple RNN with forward propagation:-

Dataset

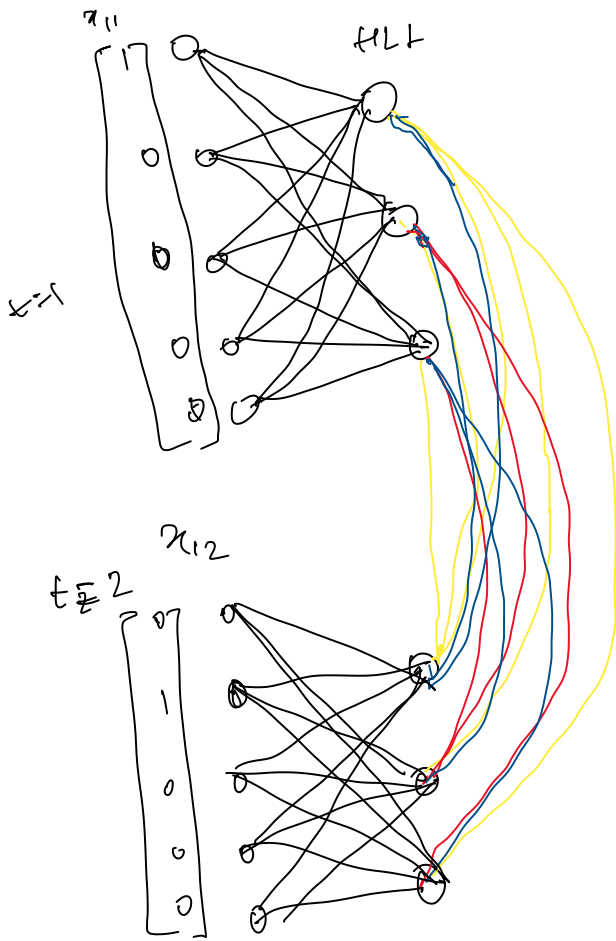
	text	o/p
s_1	the <u>food</u> is <u>good</u>	1
s_2	the food is <u>bad</u>	0
s_3	The food is <u>not</u> good	0

OKE

$\begin{bmatrix} [1 & 0 & 0 & 0 & 0], \\ [0 & 1 & 0 & 0 & 0], \\ [0 & 0 & 1 & 0 & 0], \\ [0 & 0 & 0 & 1 & 0], \\ [0 & 0 & 0 & 0 & 1] \end{bmatrix}$

o/p

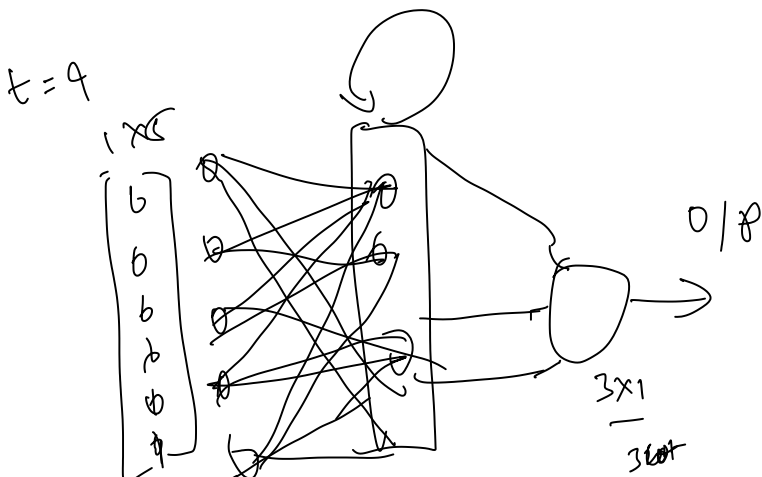
HP \rightarrow  \rightarrow IT



$t=3$

0	0
0	0
0	0
...	...

$t=4$



$5 \times 3 = 15 \text{ wt.} + 1 \text{ b.}$
 $3 \times 3 = 9 \text{ wt.} + 1 \text{ b.}$
 $\Rightarrow 30 \text{ parameters are there.}$

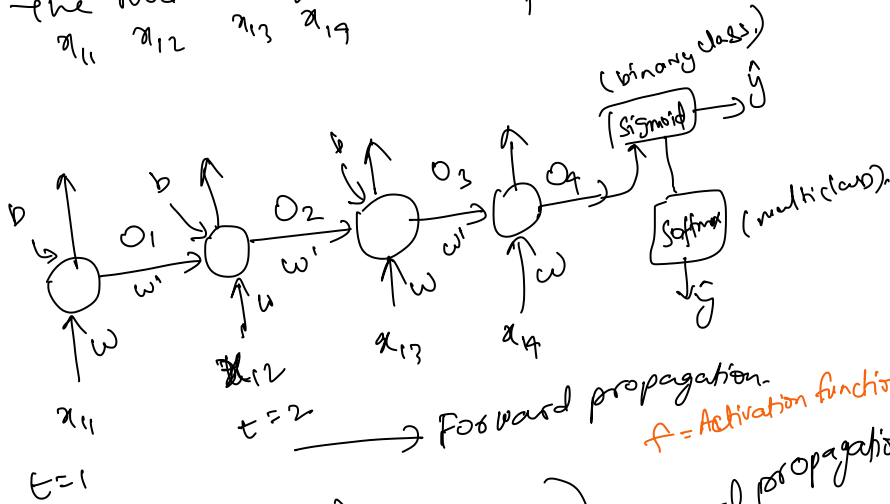
Dataset :-

the food is good.

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

O/P

1



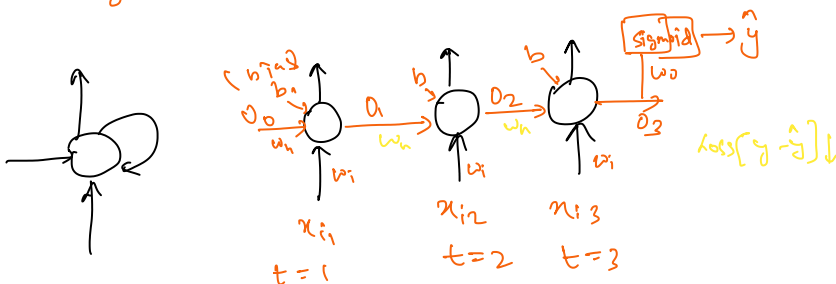
$t=1, O_1 = f(x_{11} \cdot w + b_1)$
 $t=2, O_2 = f(x_{12} \cdot w + O_1 \cdot w' + b)$
 $t=3, O_3 = f(x_{13} \cdot w + O_2 \cdot w' + b)$
 $t=4, O_4 = f(x_{14} \cdot w + O_3 \cdot w' + b)$

forward propagation.

$f = \text{Activation function.}$

RNN Back Propagation with time \rightarrow

How the back propagation happens and weights are updated in backpropagation.



Forward Propagation:-

$O_1 = f(x_{11} w_i + O_0 w_h + b)$
 $O_2 = f(x_{12} w_i + O_1 w_h + b)$
 $O_3 = f(x_{13} w_i + O_2 w_h + b)$
 \vdots
 $(n \times w_o)$

To reduce the loss we update the weights

$[w_i, w_h, w_o]$ & then we need to

$$\hat{y} = \sigma(L^{-1})$$

do backward propagation to reach global minima.

Backward Propagation with time:-

update $[w_i, w_h, w_o]$

① weight updation formula:-

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

→ slope of gradient descent.

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

based on chain Rule:

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{\text{old}}}$$

② update w_h [Hidden layer weights]:

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

$$t=1, 2, 3 \quad \frac{\partial L}{\partial w_{h_{\text{old}}}} = \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_h} \right] + \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_h} \right]$$

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

③ updating weights w_i :

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

$$\frac{\partial L}{\partial w_{i_{\text{old}}}} = \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_{i_{\text{old}}}} \right] + \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_{i_{\text{old}}}} \right]$$

$$+ \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_{i_{\text{old}}}} \right]$$

ANN - Vanishing gradient problem.

We can face the same problem in RNN too.

① The long term dependency cannot be captured by ANN with accuracy.

$$t=1 \quad \frac{\partial L}{\partial W_{hold}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_2} \cdot \frac{\partial O_2}{\partial O_1} \cdot \frac{\partial O_1}{\partial W_{hold}}$$

suppose if the sentence length is 50 words

$$at t=1, \quad \frac{\partial L}{\partial W_{hold}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_{50}} \cdot \frac{\partial O_{50}}{\partial O_{49}} \dots \frac{\partial O_2}{\partial O_1} \cdot \frac{\partial O_1}{\partial W_{hold}}$$

$$\begin{aligned} \frac{\partial O_2}{\partial O_1} &= \frac{\partial \sigma(x_{i3} \cdot W_i + O_2 \cdot W_h + b)}{\partial O_2} \\ &= \sigma'(1 \cdot W_h) \quad \text{Derivative of sigmoid} \\ &\quad \downarrow \\ &\quad [0 \rightarrow 0.25] \end{aligned}$$

At $t=1 \quad \frac{\partial L}{\partial W_{hold}} \approx 0$ [not participating much in updation of value]

⇒ The initial word is not playing significant role in the output.

⇒ The nearest words only have the significant impact.
the chain rule becomes bigger and approximates to 0.

This problem is known as Vanishing Gradient Problem.

→ In order to solve this we can use other activation function such as Relu.
Leaky Relu.

There is another solution:-

① LSTM RNN → Long short Term Memory RNN

② GRU RNN

