

Lecture NO.8 (TF-IDF)

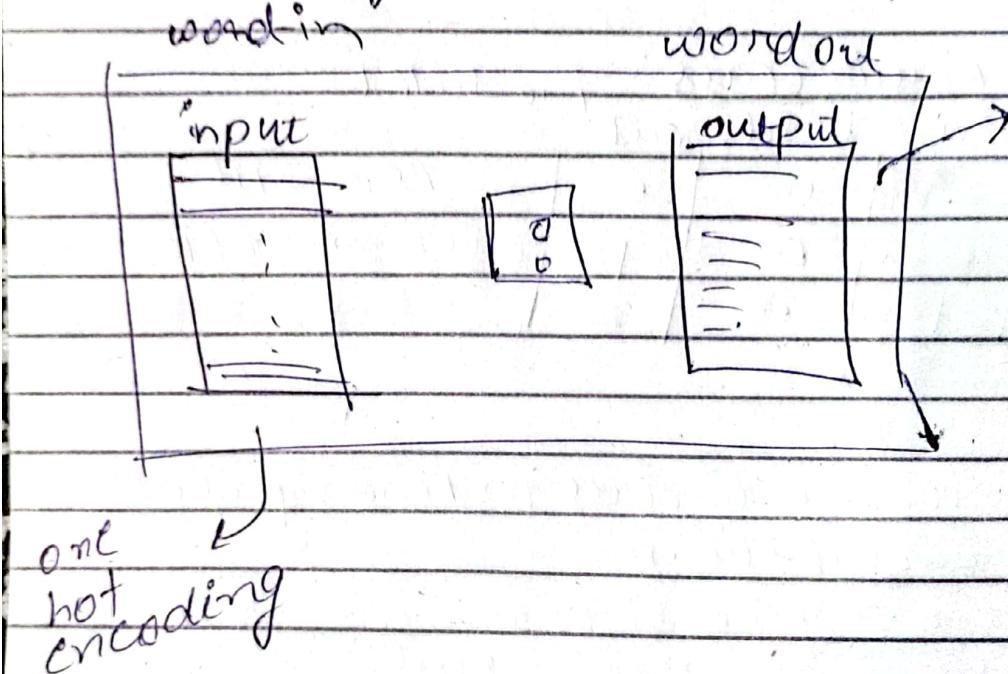
Lecture NO.9

→ word representation

↳ word embedding

↳ TF-IDF
Sparse vector

→ Embedding dimension are dimension controlled.



→ Types of words

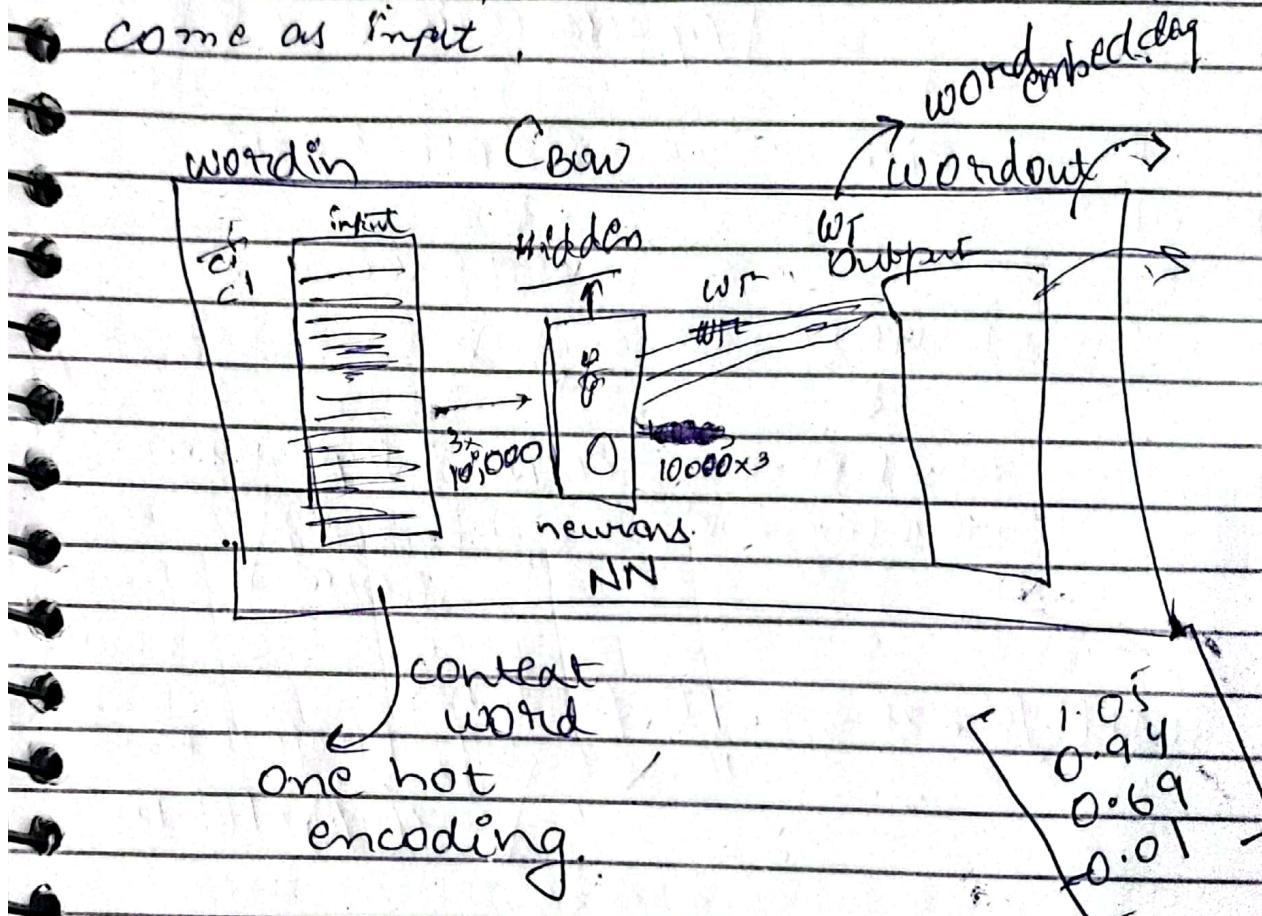
-) Context word
-) Target word.

→ Types of word2Vec:

-) CBOW
-) Skip-gram

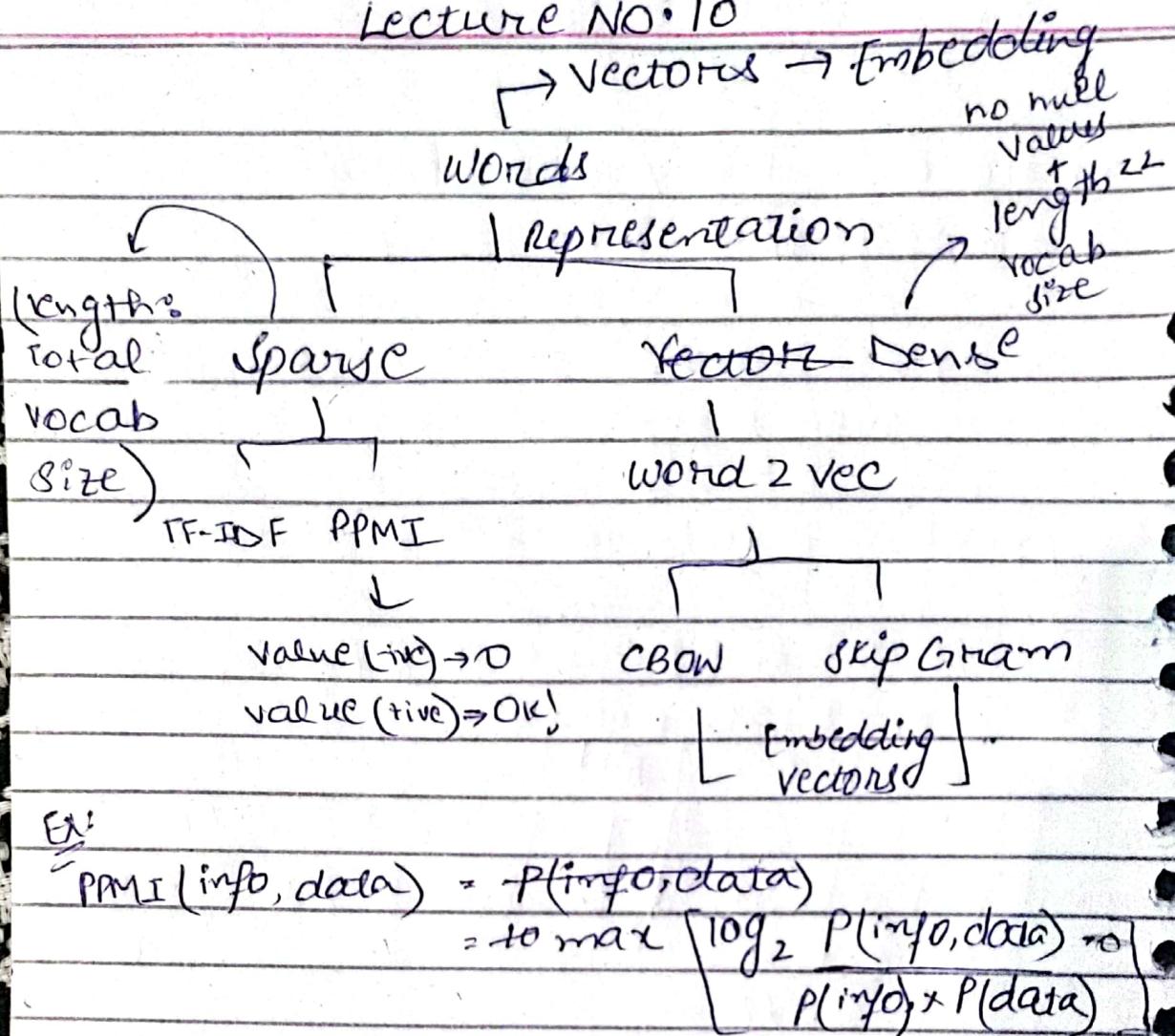
Self-Supervised!

- In case of CBOW context will come as input.



- training the model based on neighbouring words.
- learning → self supervision → no human involvement.
- No. of hidden units = Size of embedding
- $[0.7 \rightarrow 0.20.1]$ suit. for learning.

Lecture NO. 10



Ex:

$$\begin{aligned}
 \text{PRMI}(\text{info}, \text{data}) &= P(\text{info}, \text{data}) \\
 &= \log_2 \frac{P(\text{info}, \text{data})}{P(\text{info}) \times P(\text{data})}
 \end{aligned}$$

→ CBOW: (used when we have large corpus)

→ input: context words

→ output: target

→ single hidden layer is used.

→ weights (hidden → output)

→ skipGram (reverse of CBW):

→ input: target

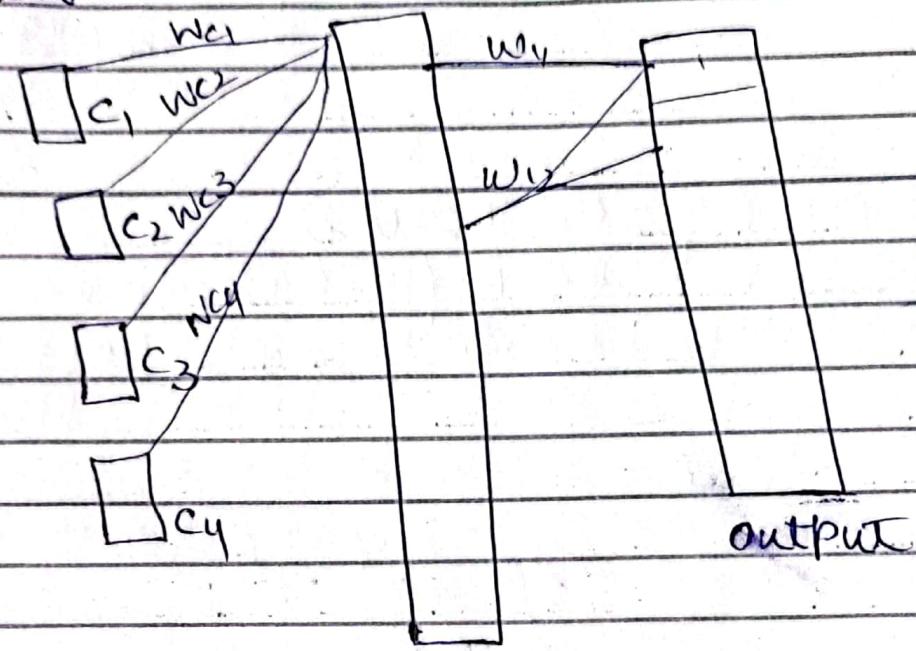
→ output: context word.

→ used when we have corpus of small size

→ weights (input → hidden)

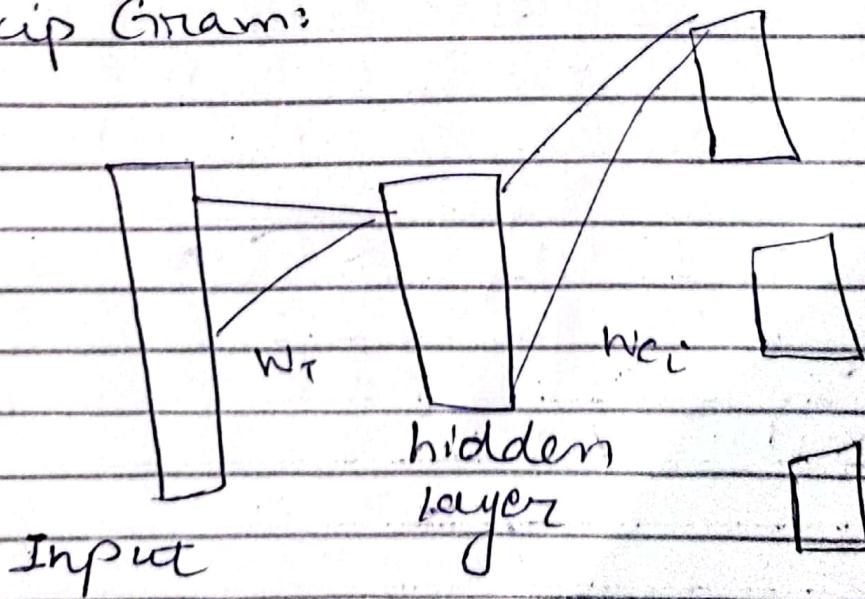
$W_t \times W_o$ one hot encoding

CBOW:



hidden
layer

Skip Gram:



Input

hidden
layer

THAT'S IT

→ Method used for training: Negative samples

Word Bag Representation

↳ Semantic Representation

↳ Dense

Lecture NO. 11

! Language Models!

Problems:

- Large Corpus → occurrence of a certain word is less
(Sparse)

In this case, Backoff is used.

information by
encoding /

- LM → looks into positional information
NLM doesn't do this.
↓
Neural Language Model.

Lecture NO.12

NN (Neural Network)

Ex:

	One-hot
$V = \{$	1 0 0 0 0 0 0
The,	0 1 0 0 0 0 0
Students,	0 0 1 0 0 0 0
opened,	0 0 0 1 0 0 0
their,	0 0 0 0 1 0 0
Laptop,	0 0 0 0 0 1 0
book,	0 0 0 0 0 0 1
exam.}	0 0 0 0 0 0 1

{The student opened their}

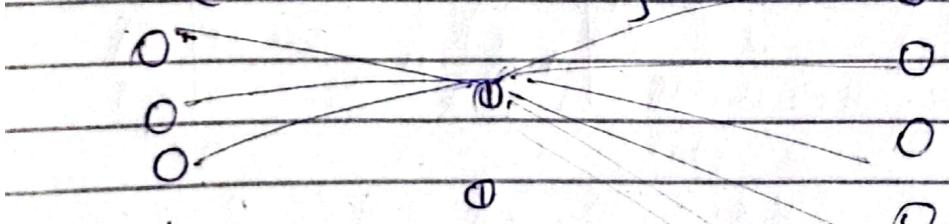
e₁ e₂ e₃ e₄

Embedding of 1st row: W_H 7x4

3 0 4 4	7 7 0 0
---------------	---------------

-8 0 -8 0	2 5 2 5
-----------------	---------------

$$E = [e_1 \ e_2 \ e_3 \ e_4]$$



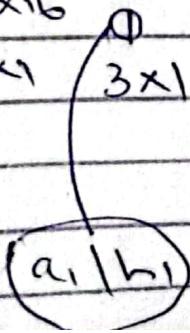
$$\therefore w_1 = 3 \times 16$$

$$\therefore b = 3 \times 1$$

:

0

16x1



$$w_2 = 7 \times 3$$

$$B = 7 \times 1$$

0

7x1

$$a_1 = (w_1 \cdot I) + B$$

$3 \times 16 \cdot 16 \times 1 + 3 \times 1$
 3×1

$$h_1 = \sigma(a_1)$$

16x1 input

3	
0	
4	
4	
7	
7	
0	
0	
-8	
0	
-8	

$w_1 = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \\ 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 & 5 \end{bmatrix}$
 8×16

$B_1 = \begin{bmatrix} 7 \\ 0 \\ 0 \end{bmatrix}$
 3×1

$a_1 = w_1 \cdot I + B_1$

23	+ 7	= 30
-54	0	-54
-135	0	-135

$h_1 = \sigma \left(\begin{bmatrix} 30 \\ -54 \\ -135 \end{bmatrix} \right) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

a_2

$$\hat{y} = w_2 \cdot h_1 + B_2$$

$7 \times 3 \quad 3 \times 1$

$\overbrace{\quad\quad\quad}^{7 \times 1}$

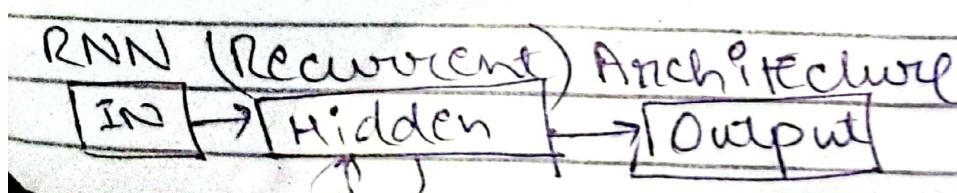
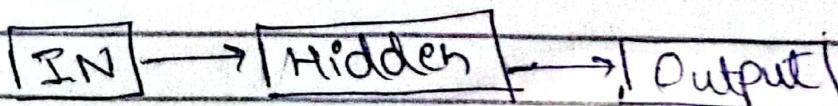
$$\hat{y} = \text{softmax}(a_2)$$

$$w_2 = \begin{bmatrix} 1 & 1 & 3 \\ 2 & 2 & 1 \\ 2 & 3 & 1 \\ 1 & 4 & 2 \\ 2 & 4 & 4 \\ 7 & 5 & 4 \\ 2 & 5 & 5 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

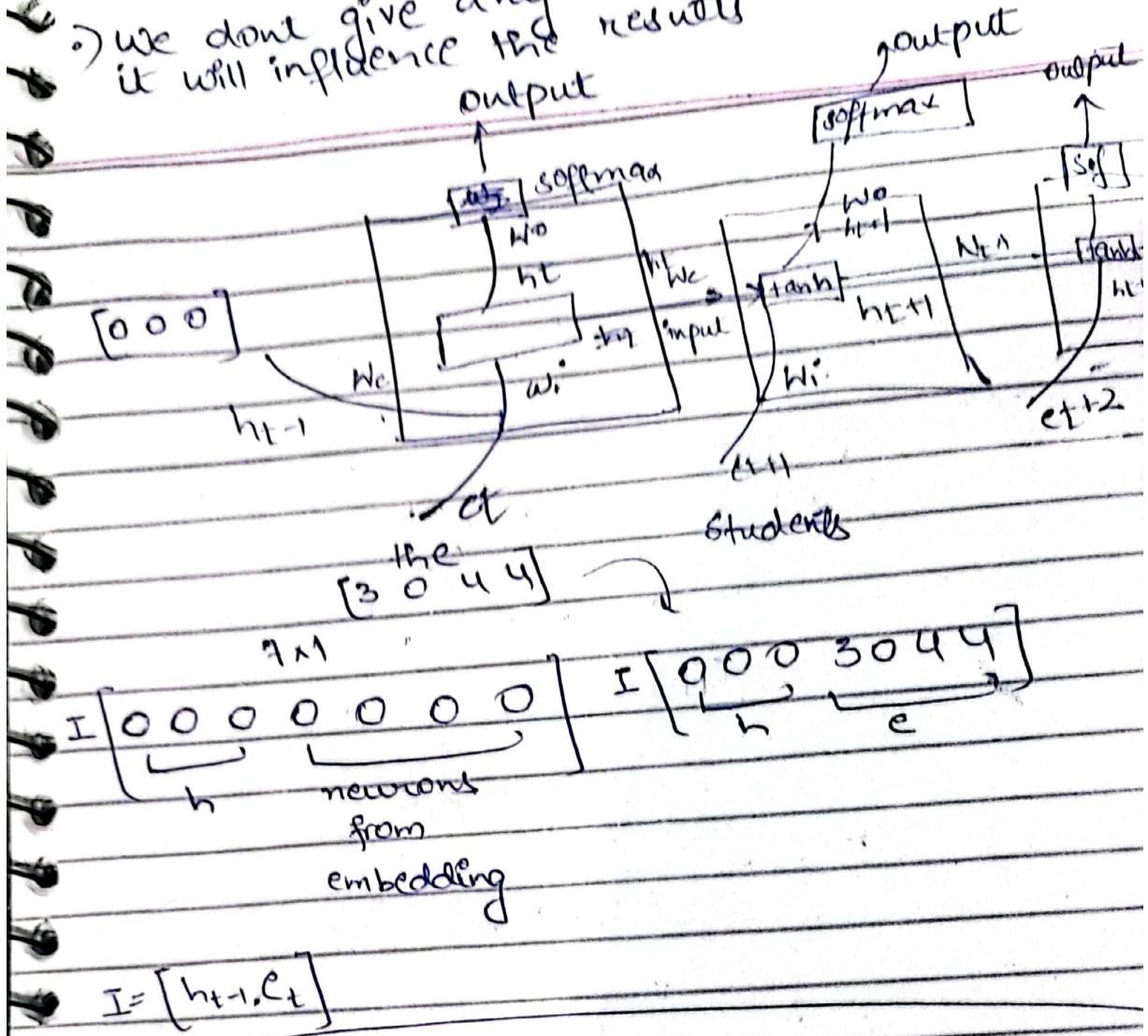
$$a_2 = \begin{bmatrix} 1 \\ 2 \\ 2 \\ 1 \\ 2 \\ 7 \\ 2 \end{bmatrix} \Rightarrow \vec{y} = \text{softmax} \begin{bmatrix} 1 \\ 2 \\ 2 \\ 1 \\ 2 \\ 7 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} 0.0024 \\ 0.0065 \\ 0.0065 \\ 0.0024 \\ 0.0065 \\ 0.9691 \\ 0.0065 \end{bmatrix} \Rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \Rightarrow \text{Books}$$

Simple NN Architecture



→ we don't give any random value bcz it will influence the results



$$I = [h_{t-1}, c_t]$$

→ RNN types:

- one-to-one
- many-to-one
- one-to-many
- many-to-many

Formalizing

$$\begin{aligned} a_t &= W_c \times h_{t-1} \cdot W_i \times e_t + b_i \\ &= W_i \cdot [h_{t-1}, e_t] + b_i \end{aligned}$$

as RNN gives a sequence

as computation is stored in hidden units, and is

$$h_t = \tanh(a_t) \rightarrow \text{'The'}$$

$$a_{t+1} = W_i \cdot [h_t, e_{t+1}] + b_i$$

$$h_{t+1} = \tanh(a_{t+1}) \rightarrow \text{'Student'}$$

The student

used further afterwards

$$a_{t+2} = W_i \cdot [h_{t+1}, e_{t+2}] + b_i$$

$$h_{t+2} = \tanh(a_{t+2})$$

↓
'The students opened'

$$a_{t+3} = W_i \cdot [h_{t+2}, e_{t+3}] + b_i$$

$$h_{t+3} = \tanh(a_{t+3})$$

↓

'The students opened their'
softmax

$$e = \text{output} = [W_o \cdot h_{t+3} + b_o]$$

vector

Disadvantage:

- None the words, None the context
is disturbed

Lecture NO. 13

RNN one-hot

$V = \{ \text{the},$
 students
 opened,
 their,
 laptops,
 books,
 exams \}

" " " " "

" " " " "

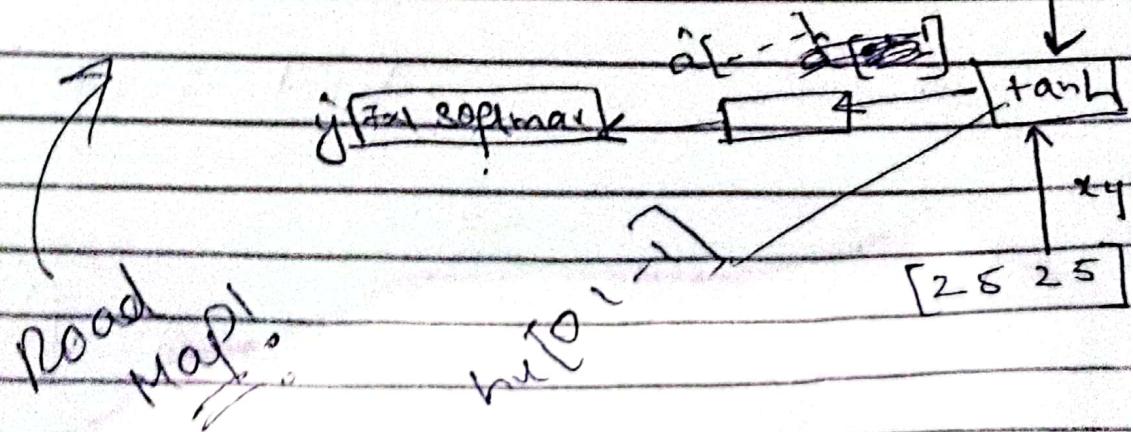
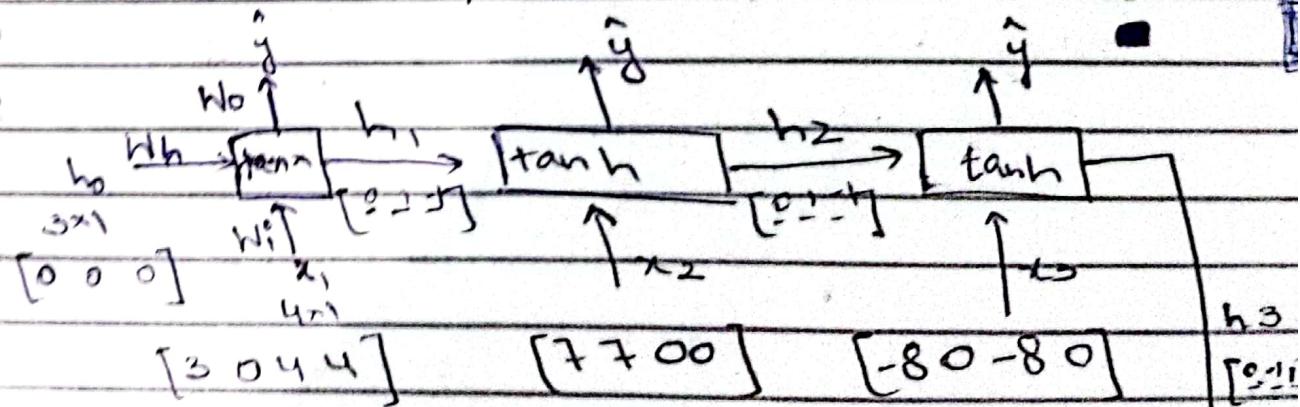
" " " " "

Embedding of 'the' $\begin{bmatrix} 3 & 0 & 4 & 4 \end{bmatrix} \dots$

" " " students $\begin{bmatrix} 7 & 7 & -8 & 0 \end{bmatrix} \dots$

" " " their $\begin{bmatrix} 2 & 5 & 2 & 5 \end{bmatrix} \dots$

" " " opened $\begin{bmatrix} -8 & 0 & -8 & 0 \end{bmatrix} \dots$



$$\begin{matrix} x_1 \\ x_2 \\ \hline -1 \\ 1 \end{matrix}$$

Ex:

$$I_1 = [h_0, x_1] = [0 \ 0 \ 0 \ 3 \ 0 \ 4 \ 4]^T \quad 7 \times 1$$

T = 0

$$h_1 = [0 \ 0 \ 0] \quad \text{Fixed.}$$

T = 1

$$h_1 = \tanh(W \times I_1 + B) \quad 3 \times 7 \quad 7 \times 1 \quad 3 \times 1$$

$$= \tanh \left(\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ -3 & -3 & -3 & 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & -2 & -2 & -2 & -2 \end{bmatrix} \times \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \right)$$

$$= \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \tanh \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix} \quad 3 \times 1$$

$$= \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}$$

$$I_2 = [h_1, x_2] = [0 \ 1 \ -1 \ 7 \ 7 \ 0 \ 0]^T$$

T = 2

$$h_2 = \tanh(W \times I_2 + B)$$

$$= \tanh \left(\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ -3 & -3 & -3 & 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & -2 & -2 & -2 & -2 \end{bmatrix} \times \begin{bmatrix} 0 \\ 1 \\ -1 \\ 7 \\ 7 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \right)$$

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\tanh \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix}$$

$T=3$

$$I_3 = [h_2, x_3] = \begin{bmatrix} 0 & 1 & -1 & -8 & 0 & -8 & 0 \end{bmatrix}$$

$$h_3 = \tanh \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ -3 & -3 & -3 & 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & -2 & -2 & -2 & -2 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ -1 \\ -8 \\ 0 \\ -8 \\ 0 \end{pmatrix}$$

$$\tanh \begin{pmatrix} 0 \\ -16 \\ 32 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 7 \end{pmatrix} = \tanh \begin{pmatrix} 0 \\ -16 \\ 39 \end{pmatrix}$$

$$= \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix}_{3 \times 1}$$

~~28
24
21~~

$T=4$

$$I_4 = [h_3, x_4] = \begin{bmatrix} 0 & -1 & 1 & 2 & 5 & 2 & 5 \end{bmatrix}$$

$$h_4 = \tanh \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ -3 & -6 & -3 & 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & -2 & -2 & -2 & -2 \end{pmatrix} = \begin{pmatrix} 0 \\ -1 \\ 1 \\ 2 \\ -5 \\ 2 \end{pmatrix}_{8 \times 1}$$

$$\tanh \begin{pmatrix} 0 \\ 1 \\ -1 \\ -18 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 7 \end{pmatrix} = \tanh \begin{pmatrix} 0 \\ 1 \\ -1 \\ -21 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix}$$



: Bias dependent on neurons.

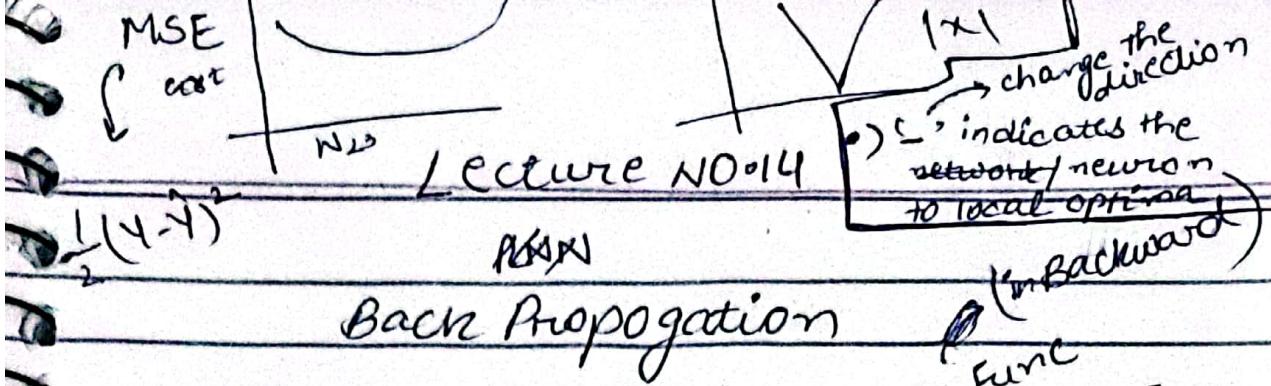
$$\vec{y} = \vec{W}\vec{x} + \vec{b} \rightarrow \text{softmax}(\vec{W}\vec{x} + \vec{b})$$

$$= \text{softmax} \left(\begin{array}{ccc|c|c} 1 & -3 & -3 & 2 \\ 1 & -2 & 3 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 2 & 2 & 0 \\ 1 & 3 & -5 & 0 \\ 1 & 2 & 2 & 0 \\ 1 & -2 & -2 & 0 \end{array} \right)$$

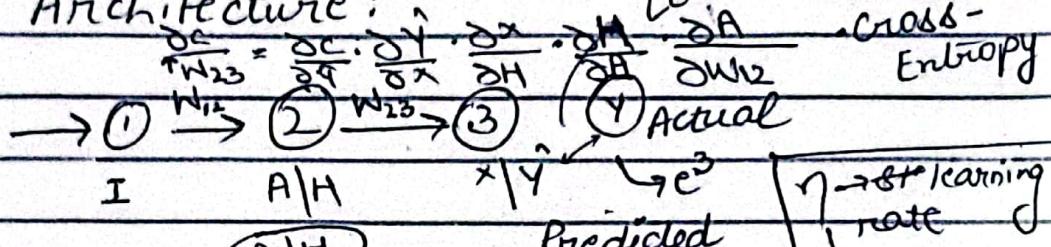
$$= \text{softmax} \left(\begin{array}{ccc|c|c} 1 & 7 & 2 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \\ 8 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{array} \right)$$

$$= \text{softmax} \left(\begin{array}{ccc|c|c} 3 & 7 & 2 & 9.1 \times 10^{-4} & \text{Argmax} \\ -5 & 2 & 0 & 5.0 \times 10^{-100} & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 10 & 0 & 0.99 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{array} \right)$$

$[0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0]$ laptops



NN Architecture:



Activation unit: Sigmoid

Forward:

$$A = I \cdot W_{12}$$

$$H = F(A) = \sigma(A)$$

$$X = H * W_{23}$$

$$\hat{Y} = F(X) = \sigma(X)$$

Backward:

$$W_{23}^{\text{new}} = W_{23}^{\text{old}} - \eta \frac{\partial C}{\partial W_{23}}$$

$$\frac{\partial C}{\partial W_{23}} = \frac{\partial C}{\partial Y} \cdot \frac{\partial Y}{\partial X} \cdot \frac{\partial X}{\partial H} \cdot \frac{\partial H}{\partial A} \cdot \frac{\partial A}{\partial W_{12}}$$

For

$$\frac{\partial L}{\partial W_{23}}$$

using chain Rule

$$\frac{\partial a}{\partial c} = \frac{\partial a}{\partial b} \cdot \frac{\partial b}{\partial c}$$

$\partial C / \partial Y$ = cost = cross-entropy

$$\rightarrow \frac{\partial C}{\partial Y} = -Y \log \hat{Y} - (1-Y) \log (1-\hat{Y})$$

$$\rightarrow \frac{\partial C}{\partial Y} = -Y \log (\hat{Y}) - (1-Y) \log (1-\hat{Y})$$

$$= -Y \times 1 \times (+1) - (1-Y) \times 1 \times (-1)$$

$$= -y + (1-y)$$

$$= \left[\frac{y}{(1-y)} \cdot \frac{(1-y)}{(1-y)} - y \right] \cdot \frac{y}{(1-y)}$$

$$\rightarrow \frac{\partial y}{\partial x} = y \times (1-y)$$

∂x ↗



Sigmoid

$$\begin{aligned}\frac{\partial c}{\partial y} \cdot \frac{\partial y}{\partial x} &= \left[\frac{(1-y)}{y} - y \right] \cdot [y \times (1-y)] \\ &= \frac{(1-y)}{(1-y)} \cdot y(1-y) - y \cdot y(1-y) \\ &= (1-y)y - y(1-y) \\ &= y - y \\ &= 0\end{aligned}$$

$$\rightarrow \frac{\partial x}{\partial w_{23}} = \frac{\partial (-H \cdot w_{23})}{\partial (w_{23})} = H$$

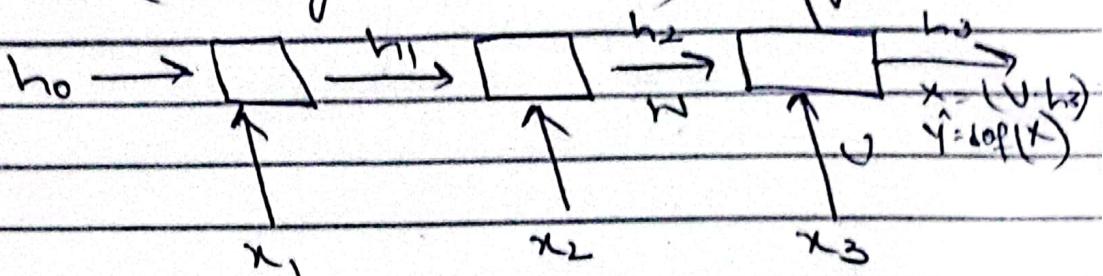
$$\rightarrow \frac{\partial L(y, \hat{y})}{\partial w_{23}} = (y - \hat{y}) \cdot H$$

$$\rightarrow \frac{\partial c}{\partial w_{23}} = \underbrace{\frac{\partial c}{\partial y}}_{} \cdot \underbrace{\frac{\partial y}{\partial x}}_{} \cdot \underbrace{\frac{\partial x}{\partial H}}_{} \cdot \underbrace{\frac{\partial H}{\partial A}}_{} \cdot \underbrace{\frac{\partial A}{\partial w_{23}}}_{} \quad \text{L L L L L}$$

$$= (\hat{y} - y) \cdot W_{23} \cdot H(1-H) \cdot I$$

(updated)

RNN (Backpropagation through Time)



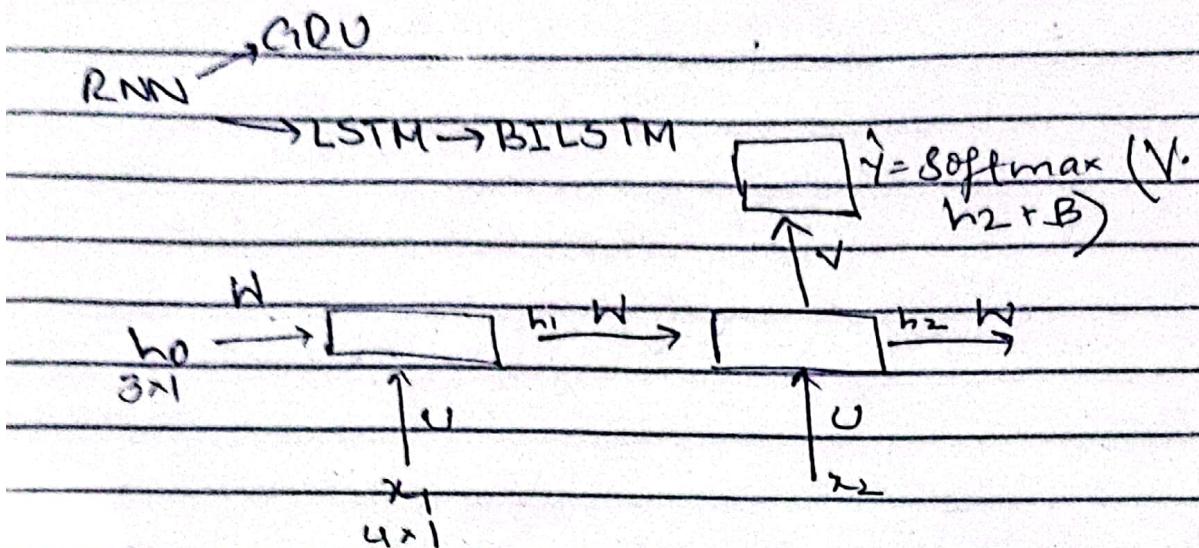
$$\frac{\partial C}{\partial V} = \frac{\partial C}{\partial y_3} \cdot \frac{\partial y_3}{\partial x} \cdot \frac{\partial x}{\partial V}$$

- Difference b/w NN & RNN is that RNN has the concept of shared weights

$$\begin{aligned} \frac{\partial C}{\partial W} &= \frac{\partial C}{\partial y_3} \cdot \frac{\partial y_3}{\partial x} \cdot \frac{\partial x}{\partial h_2} \cdot \frac{\partial h_2}{\partial W} + \\ &\quad " \quad " \quad " \quad " \cdot \frac{\partial h_2}{\partial h_1} \cdot \\ &\quad \frac{\partial h_1}{\partial W} \end{aligned}$$

- MS(a) of RNN \rightarrow 1) vanishing gradient problem
2) exploding gradient

Lecture NO. 15



$$\begin{aligned} W' &= [W \cdot U] & h_1 &= \tanh(W' \cdot I + B) \\ \cancel{3 \times 7 \times 3} & I = \begin{bmatrix} h_0, x_1 \end{bmatrix} & 3 \times 1 & 3 \times 1 \\ 3 \times 1 & 3 \times 1 \quad 4 \times 1 & & \\ h_1 &= \tanh \left(\begin{bmatrix} W_{3 \times 3} & U_{3 \times 1} & B_{1 \times 1} \end{bmatrix} \right) & 3 \times 1 & 3 \times 1 \end{aligned}$$

$$\frac{\partial C}{\partial V} = \frac{\partial C}{\partial Y} * \frac{\partial Y}{\partial V}$$

$$\frac{\partial C}{\partial W} = \frac{\partial C}{\partial Y} * \frac{\partial Y}{\partial h_2} * \frac{\partial h_2}{\partial W} + \frac{\partial C}{\partial Y} * \frac{\partial Y}{\partial h_2}$$

$$* \frac{\partial h_2}{\partial h_1} * \frac{\partial h_1}{\partial W} +$$

Too much multiplication may cause the resultant value to be small causing following problems

-) Vanishing Gradient Descent)
-) Exploding Gradient

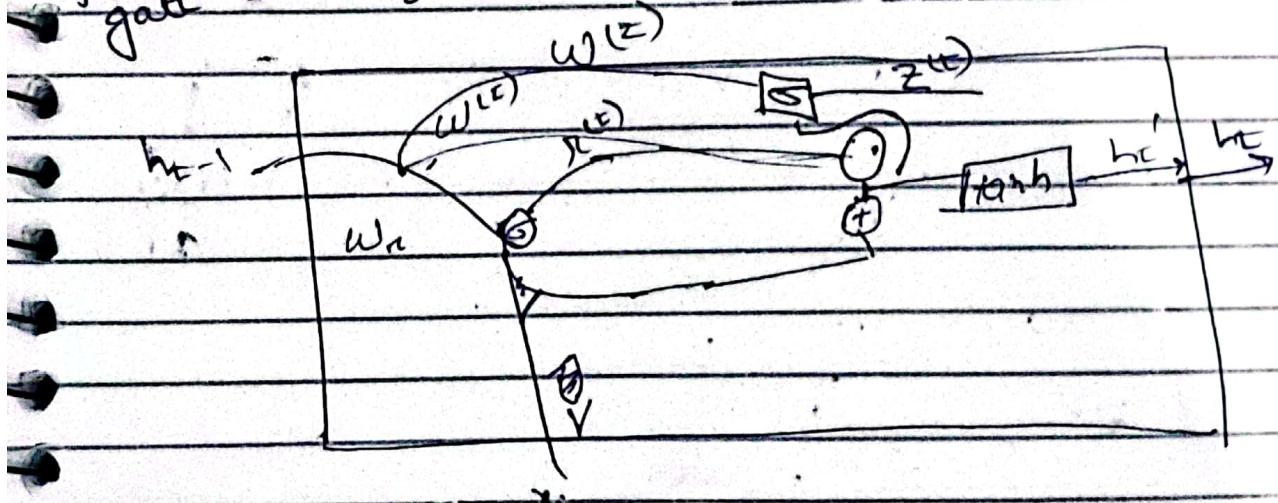
GRU:

Gates:

-) update → to update prior context
-) Reset → prior context $\xrightarrow{\text{forget}}$ how much

\downarrow similar to forget gate

$$i_t = \sigma \left[w^{(i)} h_{t-1} + u^{(i)} z_t \right]$$



$$h_t' = \tanh \left(w^{(r)} h_{t-1} + w^{(u)} z_t \right).$$

Ex: $h_{t-1} = [7 \ 1 \ 7 \ 1 \ 7] \quad \text{same size}$

$$z_t = [1 \ 1 \ 0 \ 1 \ 0]$$

\downarrow
forget: $[7 \ 1 \ 0 \ 1 \ 0]$

Preserving 7 forgetting rest

•) weight matrix of reset, update & final output will be different.

•) update

$$z^{(t)} = \sigma (W^{(z)} h_{t-1} + U^{(z)} x_t)$$

Sigmoid

$$h_t = z^t + h_{t-1}^* (1 - z^t) * h_{t-1}$$

Ex:

$$\begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix} \xrightarrow{z^t} \begin{bmatrix} 1 & 0 & 1 & 1 \end{bmatrix}$$

$$r_t^{(t)} = \sigma (W^{(r)} h_{t-1} + U^{(r)} x_t) \Rightarrow \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$$

$$z^{(t)} = \sigma (W^{(z)} h_{t-1} + U^{(z)} x_t) \Rightarrow \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$$

$$h_t' = \tanh (U^{(h)} x_t + W^{(h)} h_{t-1}) \Rightarrow \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$$

$$h_t = z^{(t)} \odot h_{t-1} + (1 - z^{(t)}) \odot h_t' \Rightarrow \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$$

Equation
GRU

Dimensions
according to
example.

FOR

$$n_0 = \boxed{7 \mid 7 \mid 7}$$

$$x_1 = \boxed{3 \mid 0 \mid 4 \mid 4}$$

$$n^{(1)} = \boxed{1 \mid 0 \mid 0 \mid 0}$$

$$z^{(t)} = \boxed{x_0 \mid x_0 \mid x_0} \quad 1 - z^{(t)} = \boxed{\text{[redacted]}} \quad \boxed{1 \mid 1 \mid 1}$$

GRU came after LSTM
↳ improved version of LSTM

Lecture NO. 16

LSTM



solved the problems of RNN

Gates:

- 1) Forget → removing all the previous
- 2) Input
- 3) Output

NATURAL LANGUAGE PROCESSING