



$$w_1 \quad [0.03 \quad 0.05 \quad -0.1 \quad 0.8 \quad \dots]$$

$$w_2 \quad [0.2 \quad 0.8 \quad -0.3 \quad \dots]$$

$$w_3 \quad \left[ \begin{array}{c} \vec{w}_1 \cdot \vec{w}_2 \\ \frac{\vec{w}_1 \cdot \vec{w}_2}{\|\vec{w}_1\| \|\vec{w}_2\|} \end{array} \right] \quad \underline{d-\text{dim}}$$

$$w_1 \quad [-0.2 \quad -0.3 \quad \dots]$$

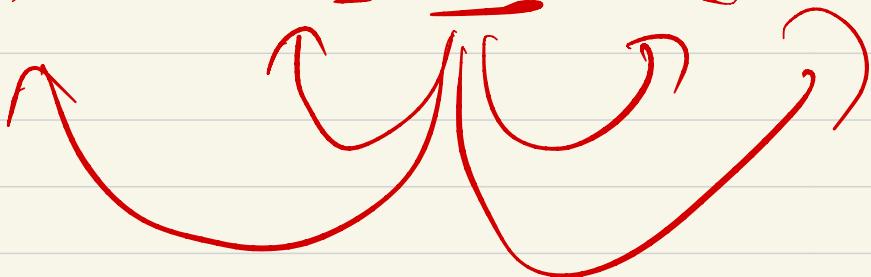
$$w_2 \quad [\dots]$$

$$w_2 \quad [\dots]$$

$$d = 100 / 300$$

w<sub>0</sub> w<sub>1</sub> w<sub>2</sub>

problems turning into bank's crises

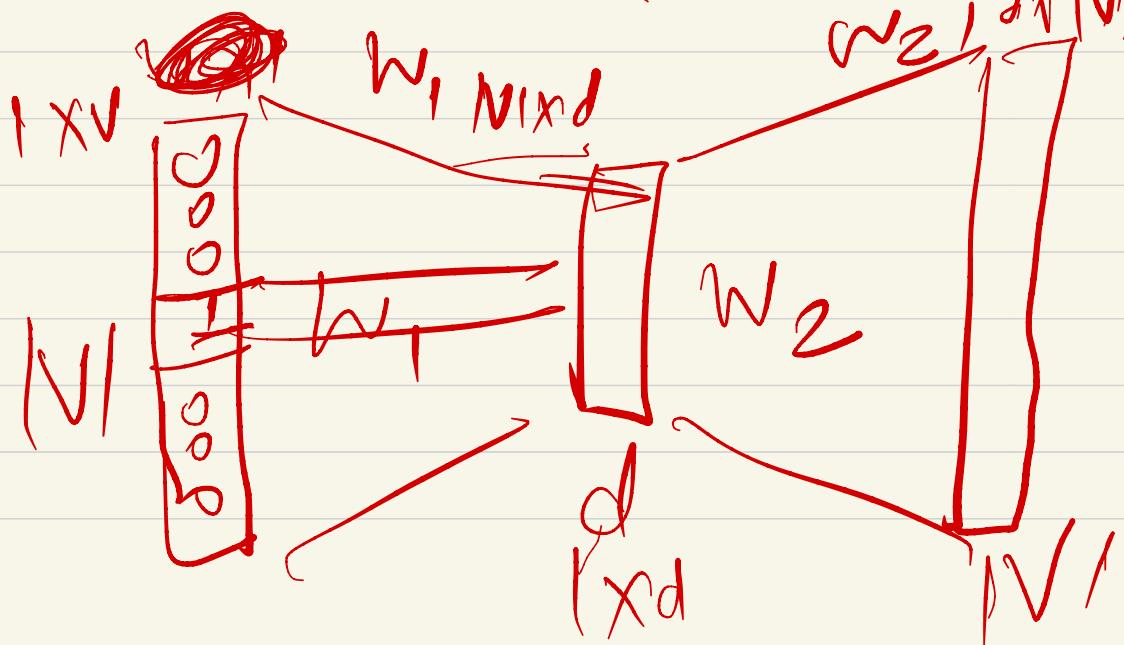
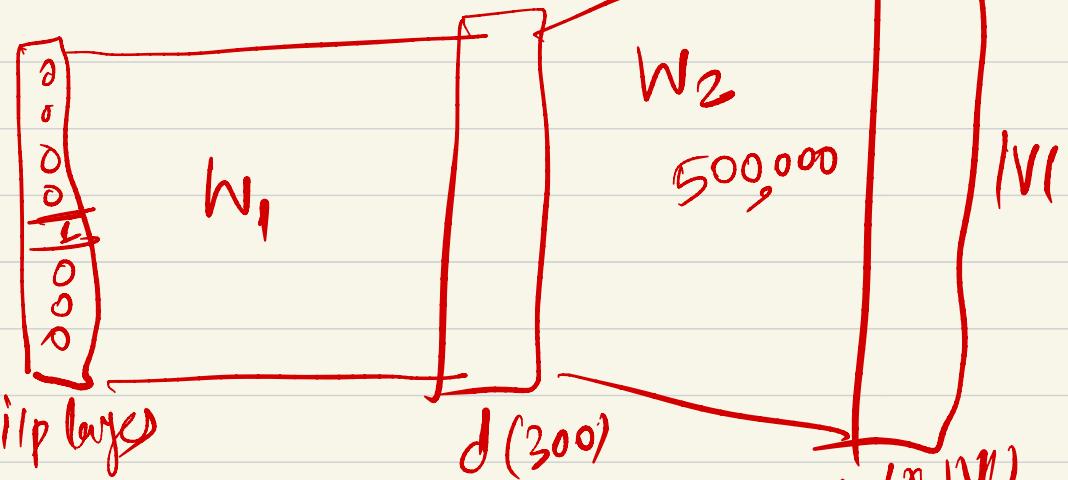


into

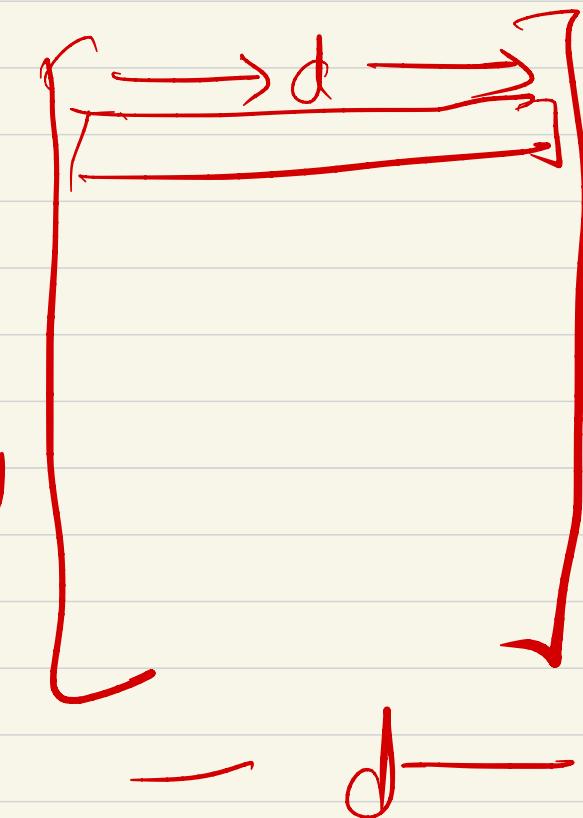
[ - ] [ - ] [ ]

[ ] [ - ] [ ]

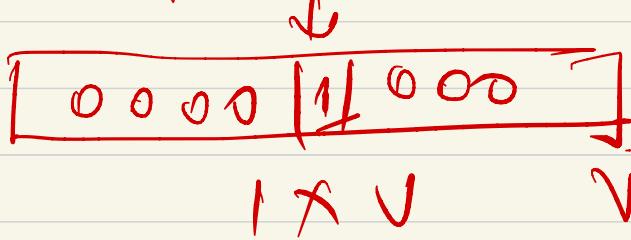
problems turning into bank's crises



$w_1$



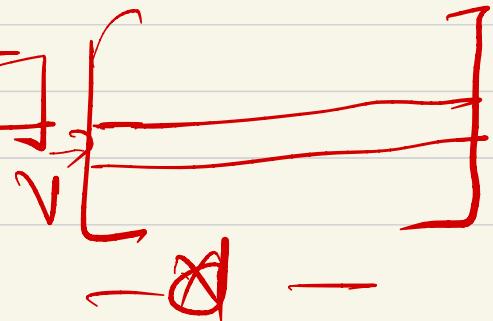
ip sector



$w_2$

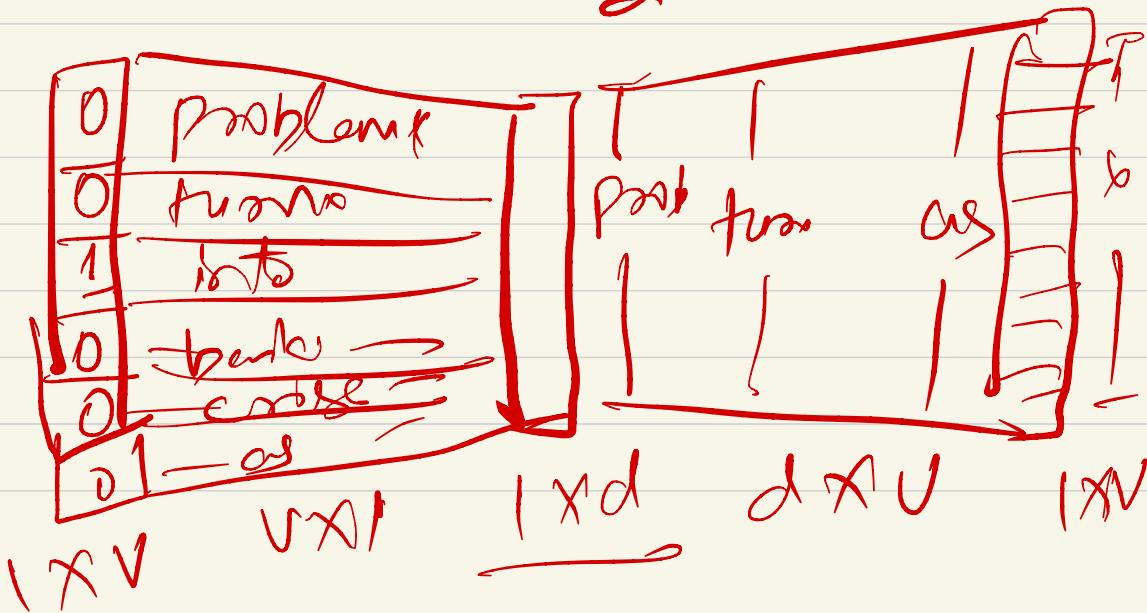


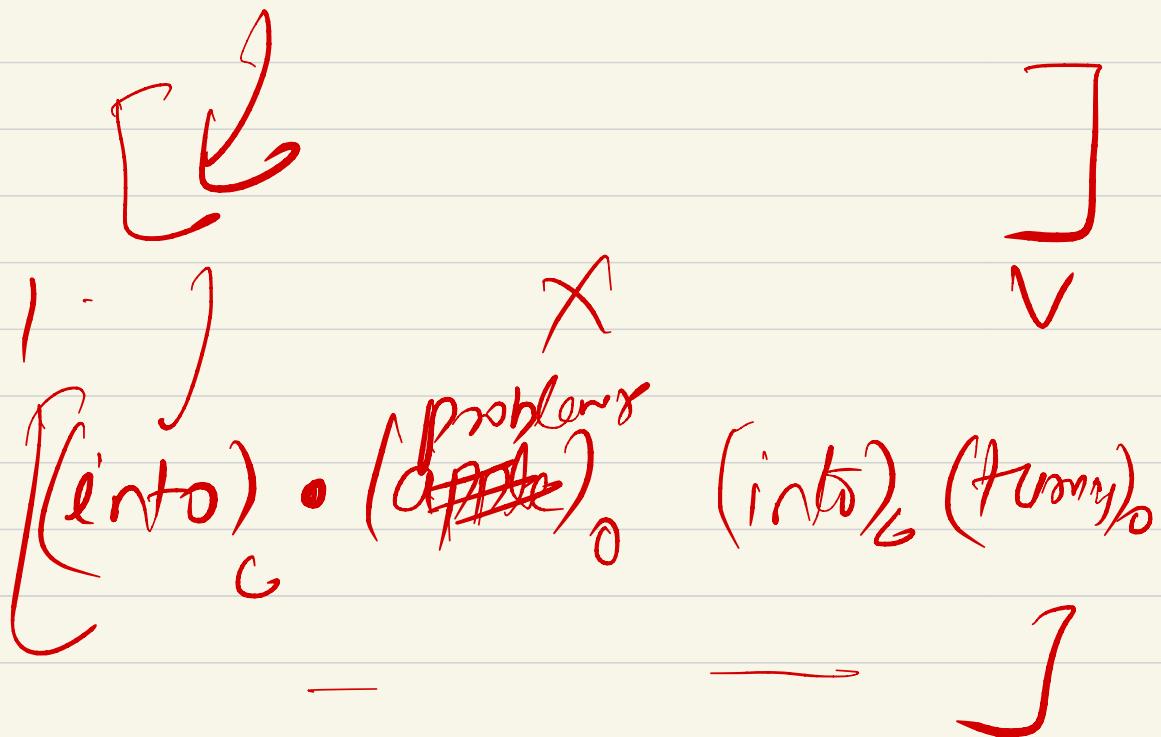
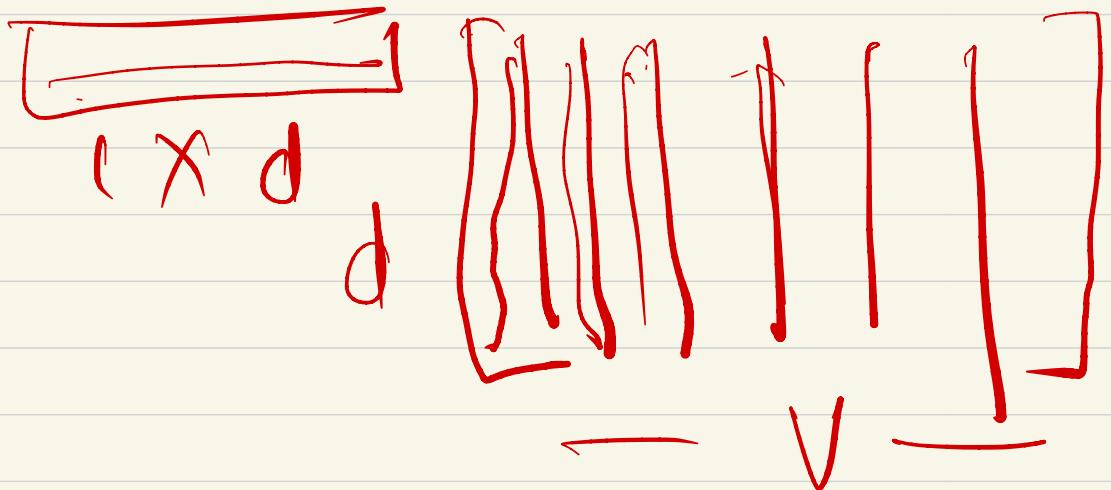
o/p vector



$$V = 6$$

problems [ ] [ ]  
storms [ ] [ ]  
into + [ ] [ ]  
bank's [ ] [ ]  
crises [ ] [ ]  
as [ ] -





$$0.5 \begin{bmatrix} 0.3 & 0.4 & 0.9 & -0.3 \end{bmatrix} \begin{bmatrix} 0.1 \\ \times \\ \checkmark \end{bmatrix}$$

problems turning into binary choices as



Pred.  $[0.5 \ 0.3 \ 0.4 \ 0.9 \ -0.3 \ 0.1]$

GT  $[1 \ 0 \ 0 \ 0 \ 0 \ 0]$

$$-\sum_i P_i \log Q_i$$

$$P = [0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0]$$

$$Q = [q_1 \ q_2 \ -q_3 \ -1]$$

$$-\log \underline{q_i}$$

$$\log \frac{1}{q_i}$$

# Softmax

$x_1$

- - -

$x_n$

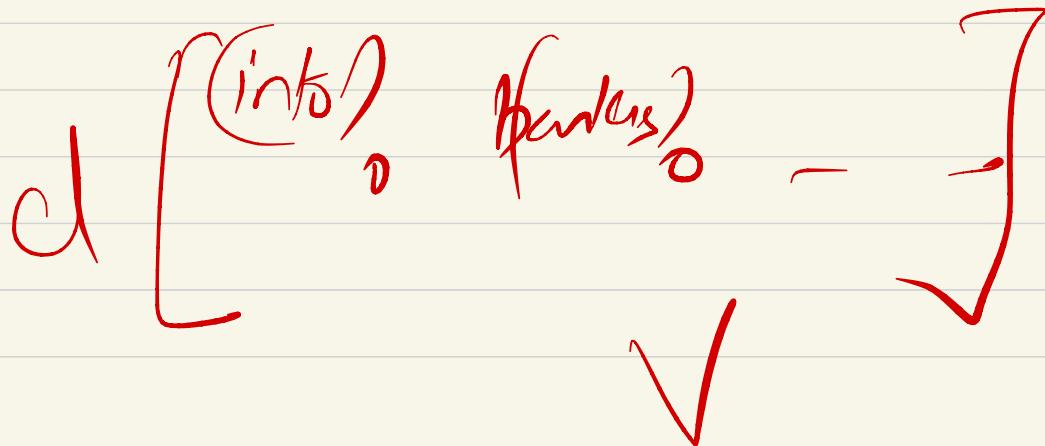
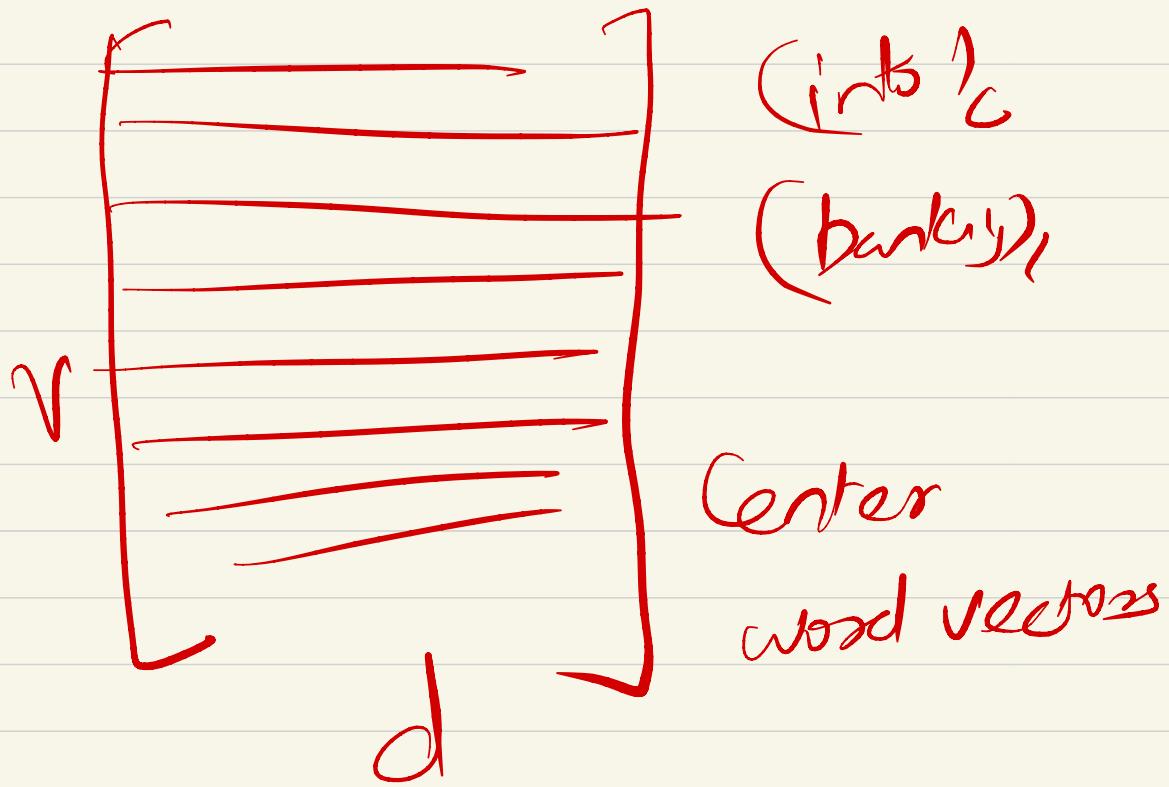


$$\frac{e^{x_1}}{\sum e^{x_i}}$$

$$\frac{e^{x_2}}{\sum e^{x_i}}$$

$$\frac{e^{x_n}}{\sum e^{x_i}}$$

$$\frac{e^{x_n}}{\sum e^{x_i}}$$



Passed through [relu] activation function  
 context window

2-d wr  
 in  $\begin{bmatrix} 1 & -1 \end{bmatrix}$   $\begin{bmatrix} 1 & 1 \end{bmatrix}$   $\begin{bmatrix} -2 & 1 \end{bmatrix}$   $\begin{bmatrix} 0 & 1 \end{bmatrix}$   $\begin{bmatrix} 1 & 0 \end{bmatrix}$

out

$$V = \frac{1}{2}$$

$$d = 2$$

$$\frac{e^{-3}}{2} \frac{e^{-1}}{\sqrt{2}} \frac{e^5}{\sqrt{2}} \frac{e^1}{\sqrt{2}} \frac{e^{-2}}{\sqrt{2}}$$

$$(1 \times 5) [-3 \quad -1 \quad 5 \quad 1 \quad -2]$$

$$-\log \frac{e^{-1}}{(e^{-3} + e^{-1} + e^5 + e^1 + e^{-2})}$$

O/p layer

Out vectors

$$-\log \left( \frac{e^{-1}}{e^{-3} + e^{-1} + e^5 + e^1 + e^{-2}} \right)$$

$$2 \times 2$$

hidden

$$5 \times 2$$

In vectors

through

$$1 \times 5$$

$$5$$

I/p layer

Cross-Entropy Error

$$P = [0, 1, 0, 0, 0] \rightarrow$$

Model prob  
dist.

$$q = [q_1, q_2, q_3, q_4, q_5]$$

$$-\sum_i p_i \log q_i = -\log(q_2)$$

## Softmax

$$\begin{matrix} 3 & 5 & -7 \\ \cancel{\frac{e^3}{e^3+e^5+e^{-7}}} & \cancel{\frac{e^5}{e^3+e^5+e^{-7}}} & \cancel{\frac{e^{-7}}{e^3+e^5+e^{-7}}} \end{matrix}$$

Softmax

$$\frac{e^3}{e^3+e^5+e^{-7}} \quad \frac{e^5}{e^3+e^5+e^{-7}} \quad \frac{e^{-7}}{e^3+e^5+e^{-7}}$$

(-hot)

$$\begin{array}{l} \text{hotel} = [0.1 \quad -0.1] \\ \text{motor} = [-0.1 \quad 0.1] \\ \text{chalk} = [-0.1 \quad -0.1] \end{array}$$

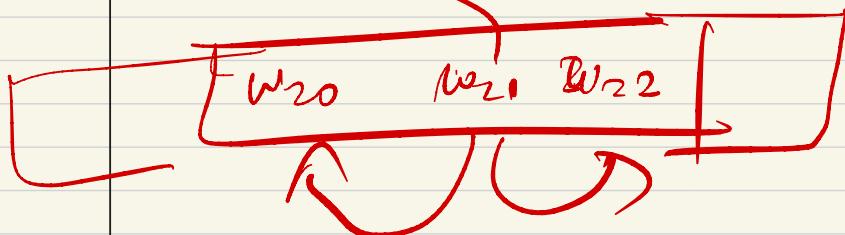
$$\begin{array}{l} \text{hotel} = [0.3 \quad -0.5] \\ \text{motor} = [0.4 \quad -0.3] \\ \text{chalk} = [-0.1 \quad 0.3] \end{array}$$

Red annotations:

- $0.27$  written above the first row.
- $-0.45$  written below the third row.

A large red brace groups all three equations together.

$\begin{bmatrix} \omega_1 & \omega_2 & \omega_3 \end{bmatrix}$



$d$        $\frac{300}{\pi}$        $\frac{100}{\pi}$

$w_{100}$

man: woman :: uncle: ??

man -	[0.3	0.7]
woman	[0.5	0.3]
uncle	[0.9	0.1]

man - woman  $\approx$  uncle -  $x$

$x \approx$  Uncle + woman - man

$$\approx [1.1 - 0.3]$$

V

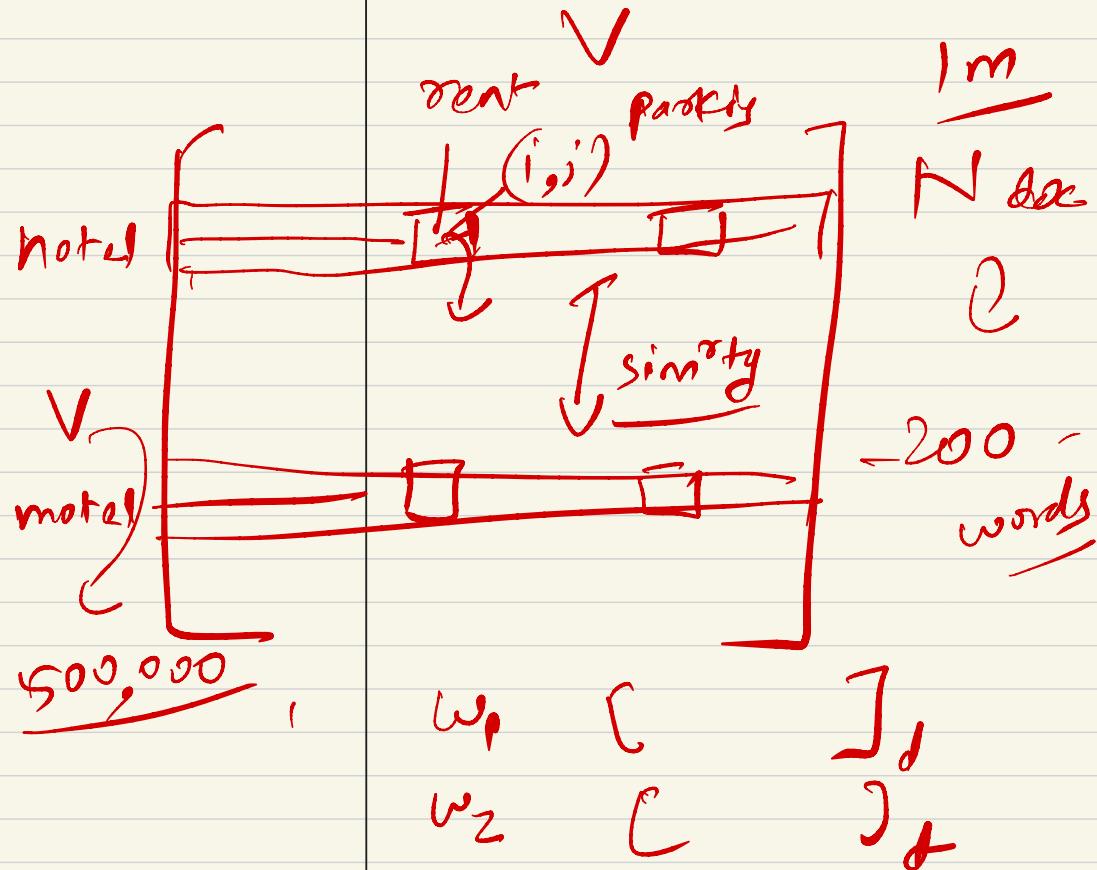
words

$$\underset{x}{\operatorname{argmax}} [1.1 - 0.3]x$$

$x = \text{aunt}$

# V words in vocab

co-occ.



$$\frac{1}{\sqrt{d}} \tilde{w}_i \cdot \tilde{w}_j \approx \text{cooc}(i, j)$$

$w_j$  [ ]  $J_d$

100



En Wiki

Hi Wiki

night [ ]

$\frac{2\pi}{\lambda}$  [ ]

noon [ ]

shuz [ ]

frame square → right  
Jumble Hi & En

mix 8 merge

Bat - flying mammal

wooden stick

Bat =

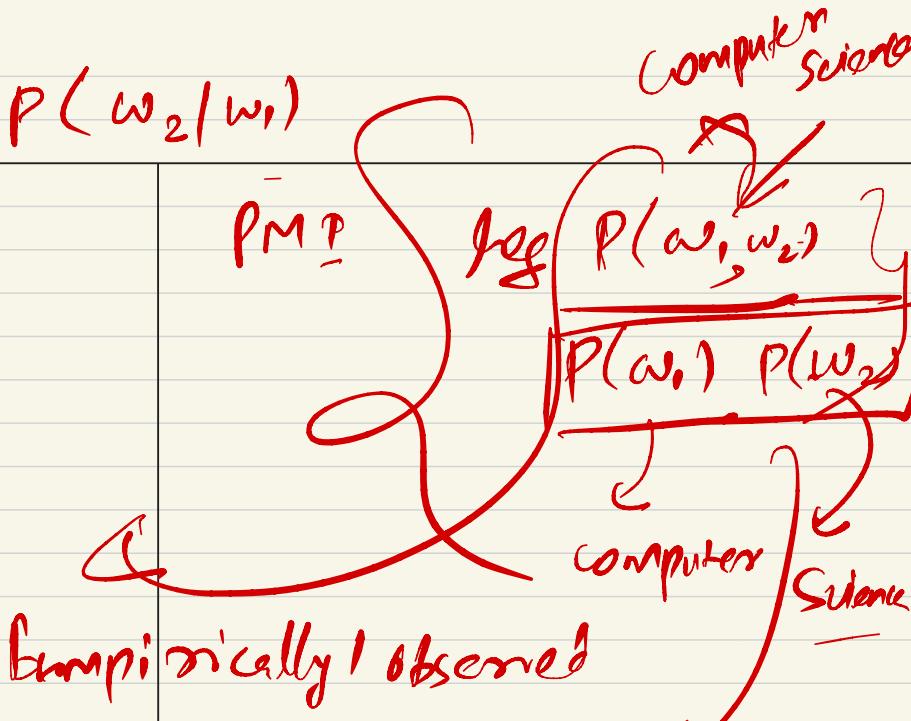


$\left[ \begin{array}{c} \text{bat} \\ \hline \end{array} \right]_1^T \text{ mar}$

$\left[ \begin{array}{c} \text{bat} \\ \hline \end{array} \right]_2^S$

$\left[ \begin{array}{c} \text{bat} \\ \hline \end{array} \right]_3^B$

$$P(w_2|w_1)$$



$$\frac{PM \geq 0}{2}$$

if they're independent

Computer science

Data

Class

mean word  
vector

$d_1$

Politics



$d_2$

Sports



$d_3$

Entertainment



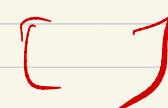
$d_4$

Politics



$d_5$

Fun!



1

1

-

$d_{100}$

Politics C1  
Sports C2  
Fun! C3

$\delta_1$   $\delta_2$

$S$  unique words in  $d$   
 $P$

$\sum_{i=1}^n$

$C^B$   
mean word vector  
mean word vector

$(d+3) \times h$

3 2 -4

3-dim

$h \times 3$

activation

act.

$h$ -dim

$d \times n$

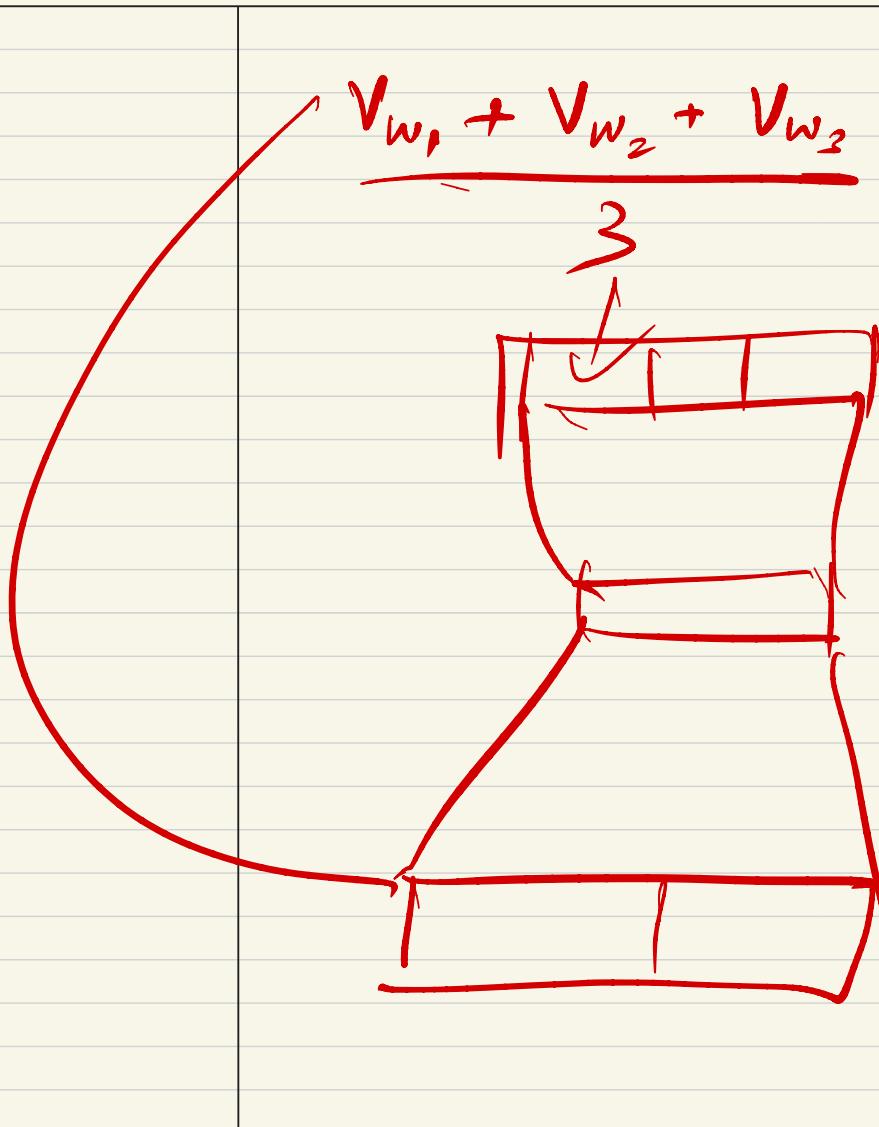
$\delta_1 \rightarrow P$

$d$ -dim

mean word vector

Up layers

$w_1$      $w_2$      $w_3$



V

$$w_1 = [0 \underline{1} 0000]$$

$$w_2 = [0 0 0 \underline{1} 00 -]$$

$$w_3 = [-1 0 0 0 0 0 -]$$

3

n

$$[1 1 0 1 0 0 0 0]$$

multi-hot

V

$$[1 1 0 \underline{1} 0 0 0]$$

multi-hot

# 10000 words

apple

1

doc.



3

apple eat

eat

500

drive

zebra

10000

[1-0-1-0-1-0-] 10000

500



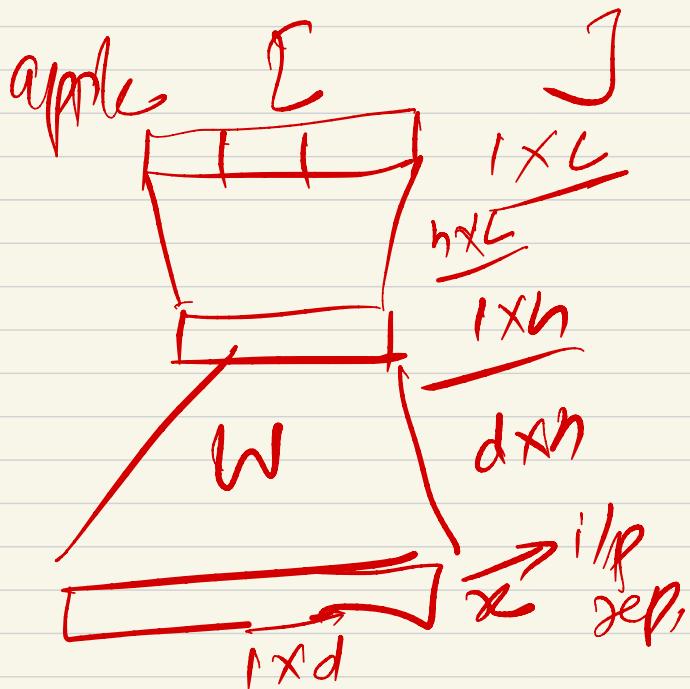
Word Vectors  
= Word Embeddings

3

Word2Vec

300d/100d

Clone



Multi-hot vector

$P(\text{the students opened their books})$

$$= \underbrace{P(\text{line})}_{P(\text{opened} \mid \text{the students})} \underbrace{P(\text{students} \mid \text{line})}_{P(\text{books} \mid \text{the student opened their books})}$$

$P(\text{opened} \mid \text{the students})$

$P(\text{books} \mid \text{the student opened their books})$

$P(\text{books} \mid \text{the student opened their books})$

$\# \text{ the}$

$\sum_{\omega} \# \omega$

$\omega$

$C(\text{the student})$

$C(\text{the})$

$C(\text{the students opened})$

$C(\text{the students})$

$x_1 \quad x_2 \quad \dots \quad x_{t-1} \quad \underline{x_t}$

1st order

$\approx$

Marks

$P(x_t | x_1, \dots, x_{t-1})$

$P(x_t | x_{t-1})$

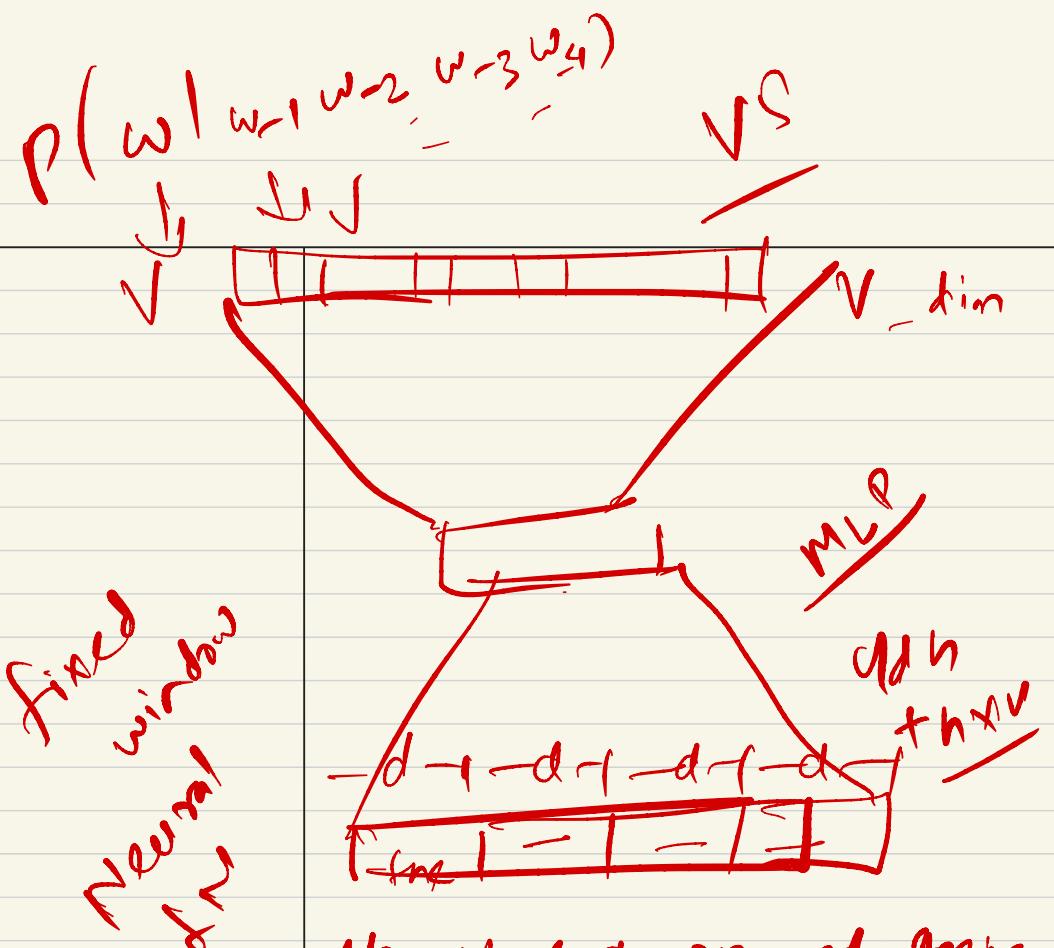
$P(\text{the students opened their books})$

$= P(\text{the}) P(\text{student})$   
 $\quad \quad \quad - P(\text{books})$

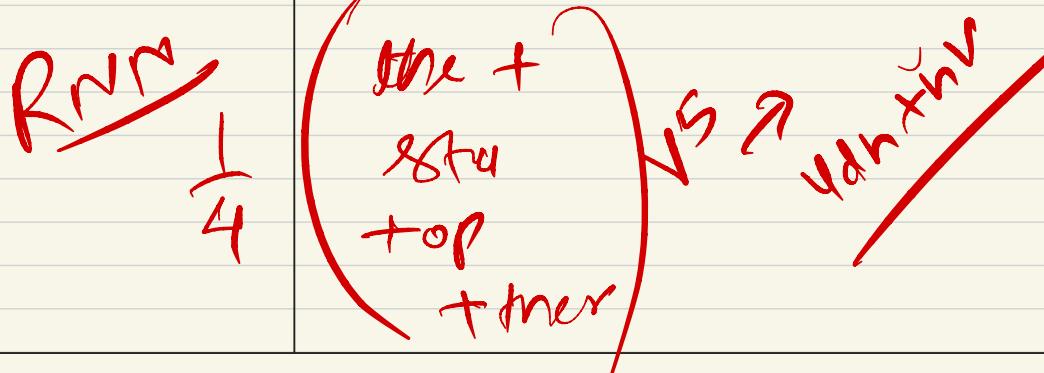
Unigram LM

bigram LM

$P(\text{the})$   
 $P(\text{students}) P(\text{the})$



the student opened their



$s_1$

$s_2$

$d_1$  ] mean  
 $d_2$  embedder

$$\left( \frac{w_1 + w_2 - \dots - w_n}{n} \right) \left( \frac{w'_1 + \dots + w'_m}{m} \right)$$

Score(I made her duck) ✓  
(Score(I'm expert or duck)) ✓

Language  
Model  
English



$$P(I \text{ made her duck})$$

Unigram  
LM

$$\frac{P(I), P(made) \cdot P(her)}{P(duck)}$$

Bigram

$$P(I|S) P(made|I) \dots$$

Trigram

$$P(w_1|w_2)$$

n-gram  
NLP

$$V^2 \quad V \quad \checkmark$$

$$P(\underbrace{I \text{ made}}_{\text{I}}, \underbrace{\text{her}}_{\text{her}}, \underbrace{\text{duck}}_{\text{duck}})$$

$$= P(I) \cdot P(\text{made} | I)$$

$$\begin{aligned} & P(\underbrace{\text{her}}_{\text{her}}, \underbrace{\text{I made}}_{I \text{ made}}) = \\ & P(\text{duck} | \text{made her}) \\ & \downarrow \\ & \frac{c(I \text{ made her})}{c(I \text{ made})} \end{aligned}$$

$$P(I'm \text{ right or duck})$$

I made her duck.

$$\arg \max_{w \in V} P(w | \text{made her}) \rightarrow w?$$

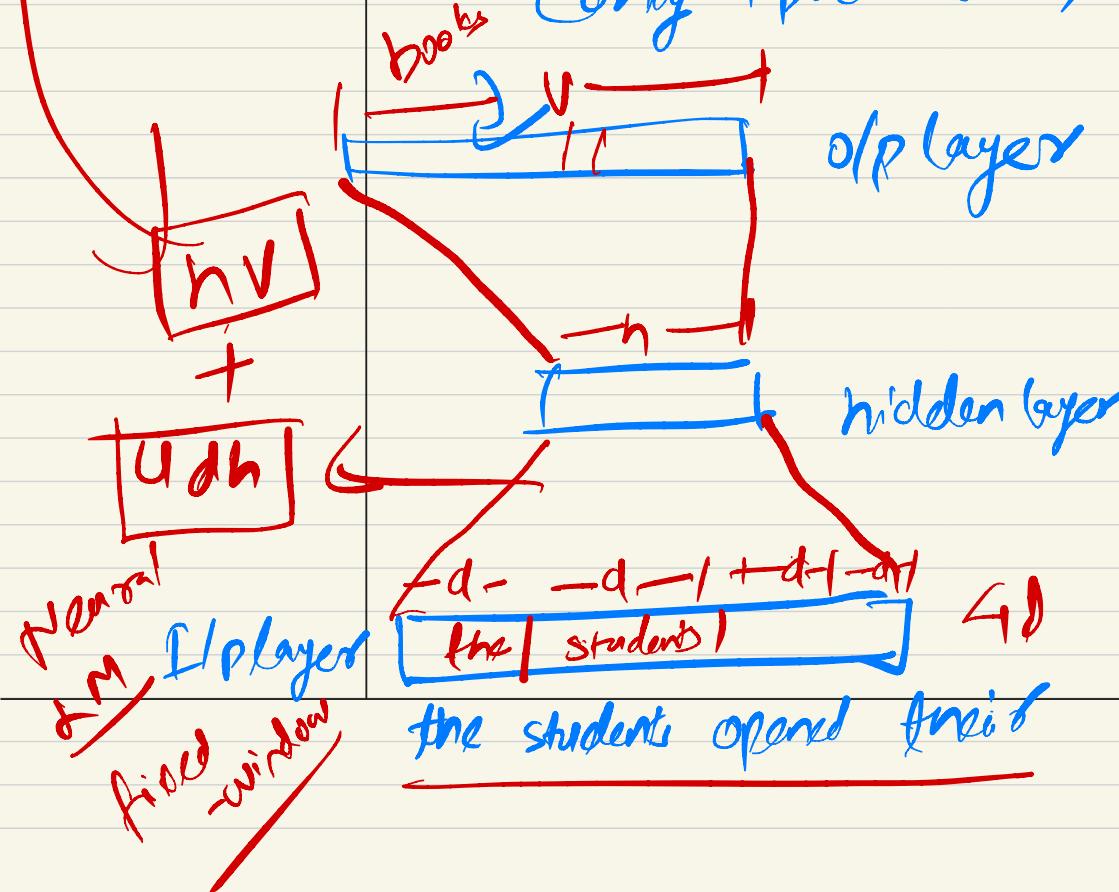
Negativprob

# Problems

1. storage  $O(V^m)$  n-gram  
5-gram  $O(V^5)$

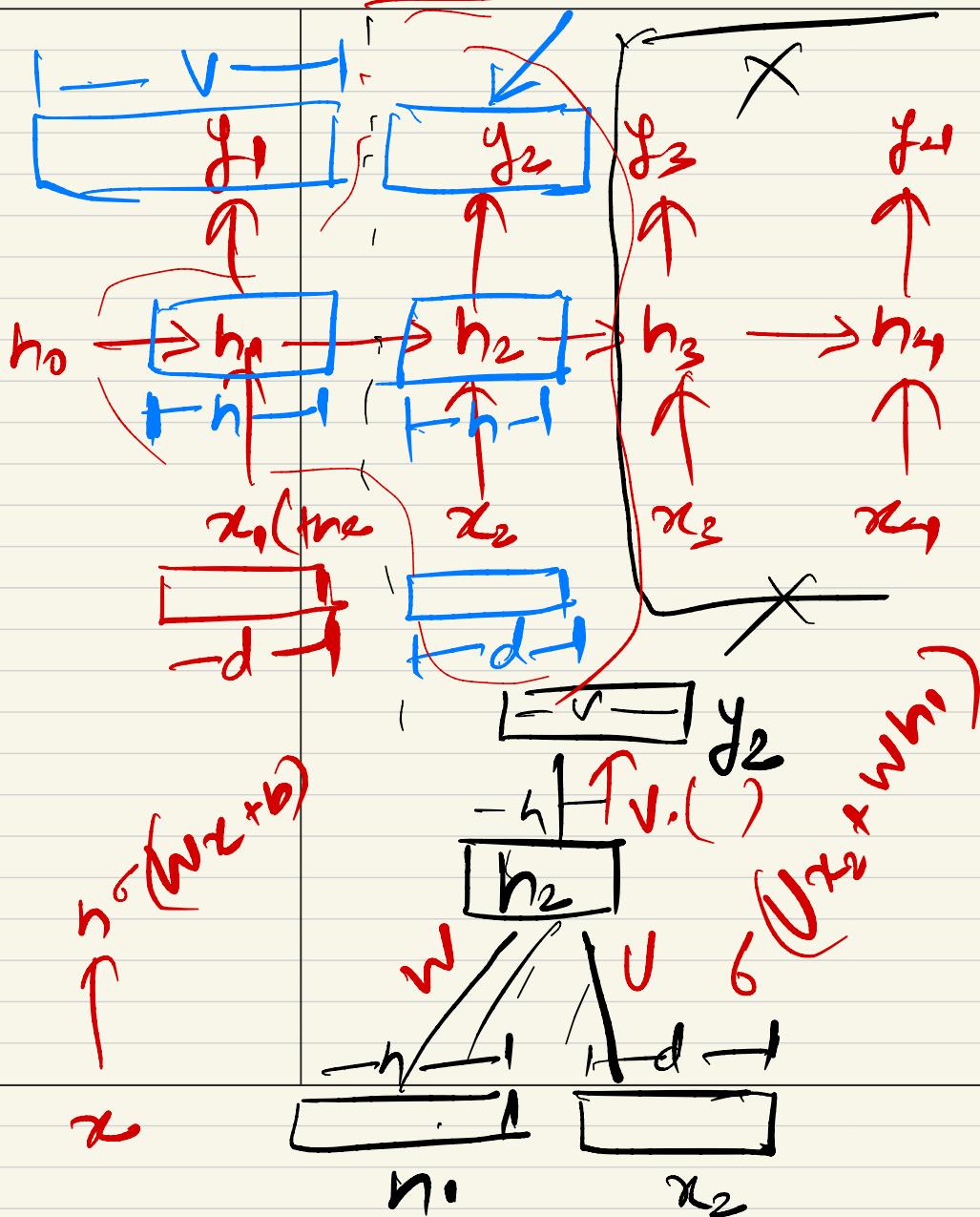
2. Limited context

(only 4 prev. words)

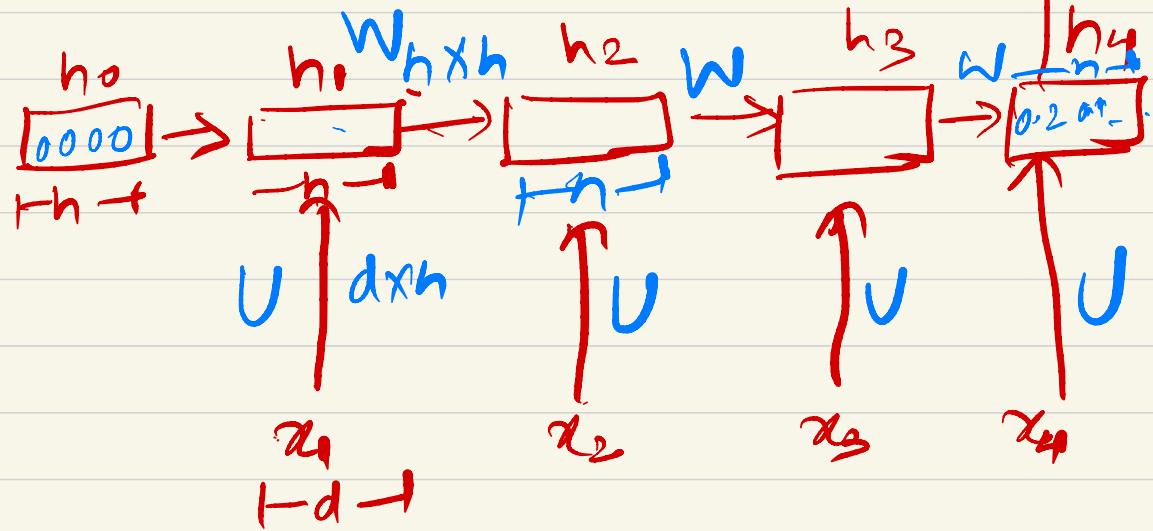
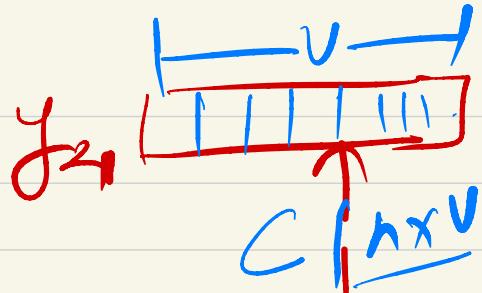


the students opened their —

## RNN



RNN



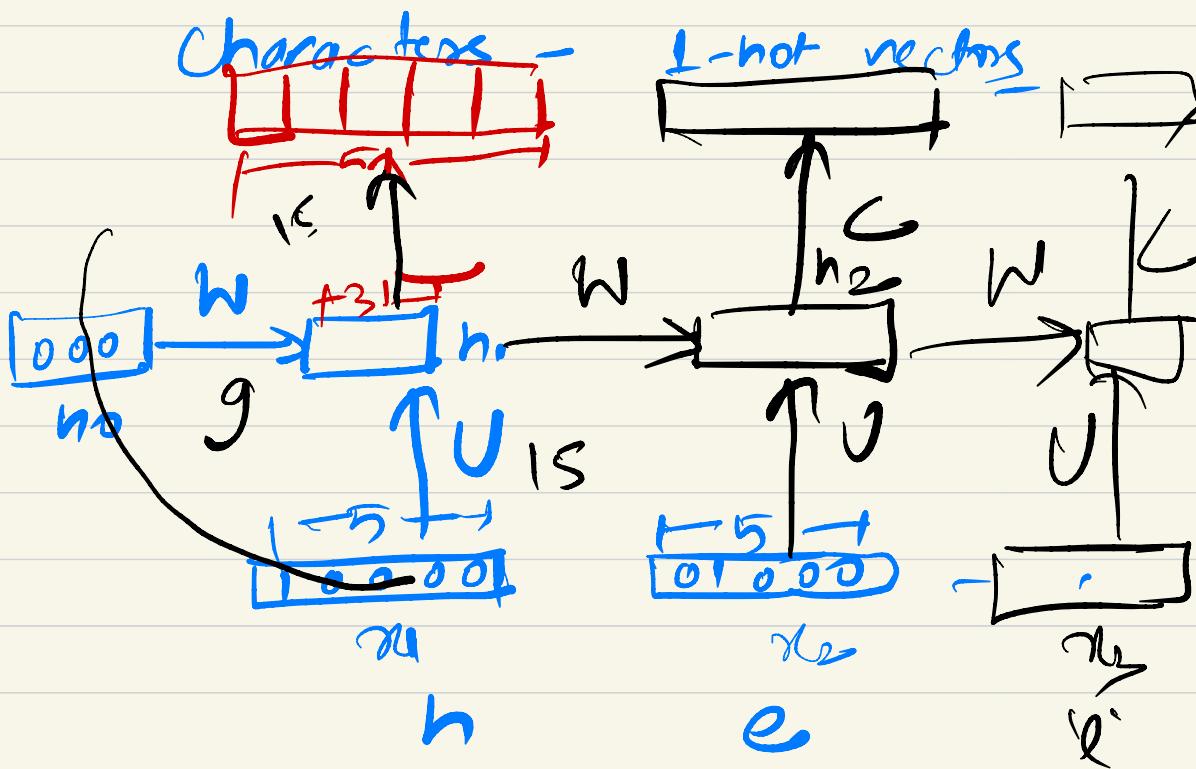
$U, W, C$

$d \times h + h \times h + h \times v$

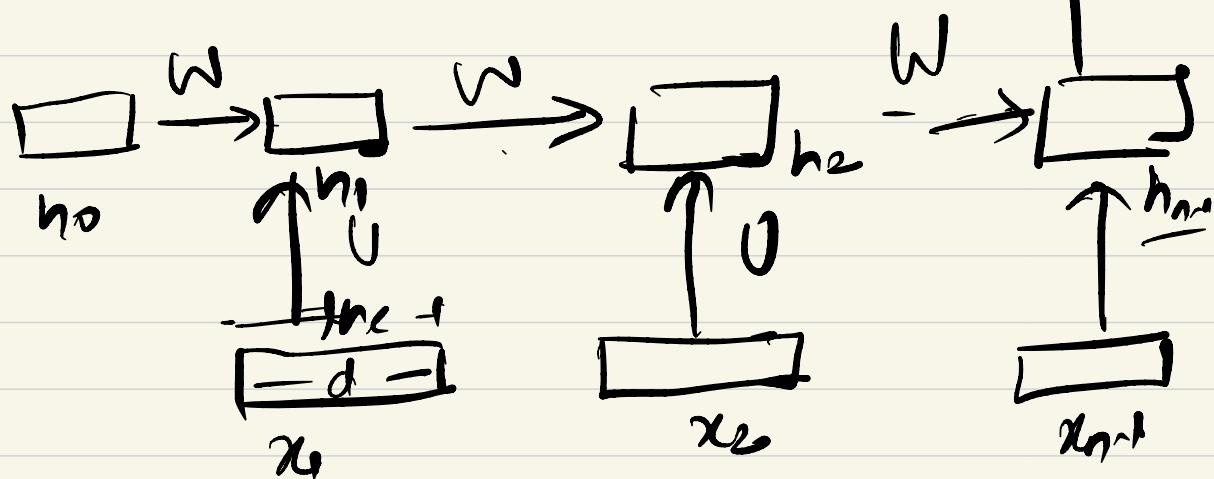
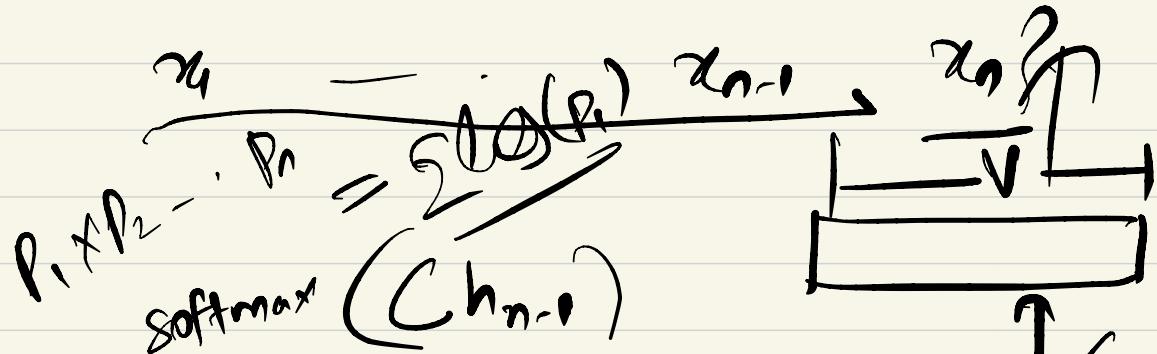
RNN for next character prediction

5 chars = ['h', 'e', 'l', 'o', '#']  
 $\downarrow_{\text{RNN}}$

Recurrent state - 3-dim  
 $n \approx 3$



RNN LM



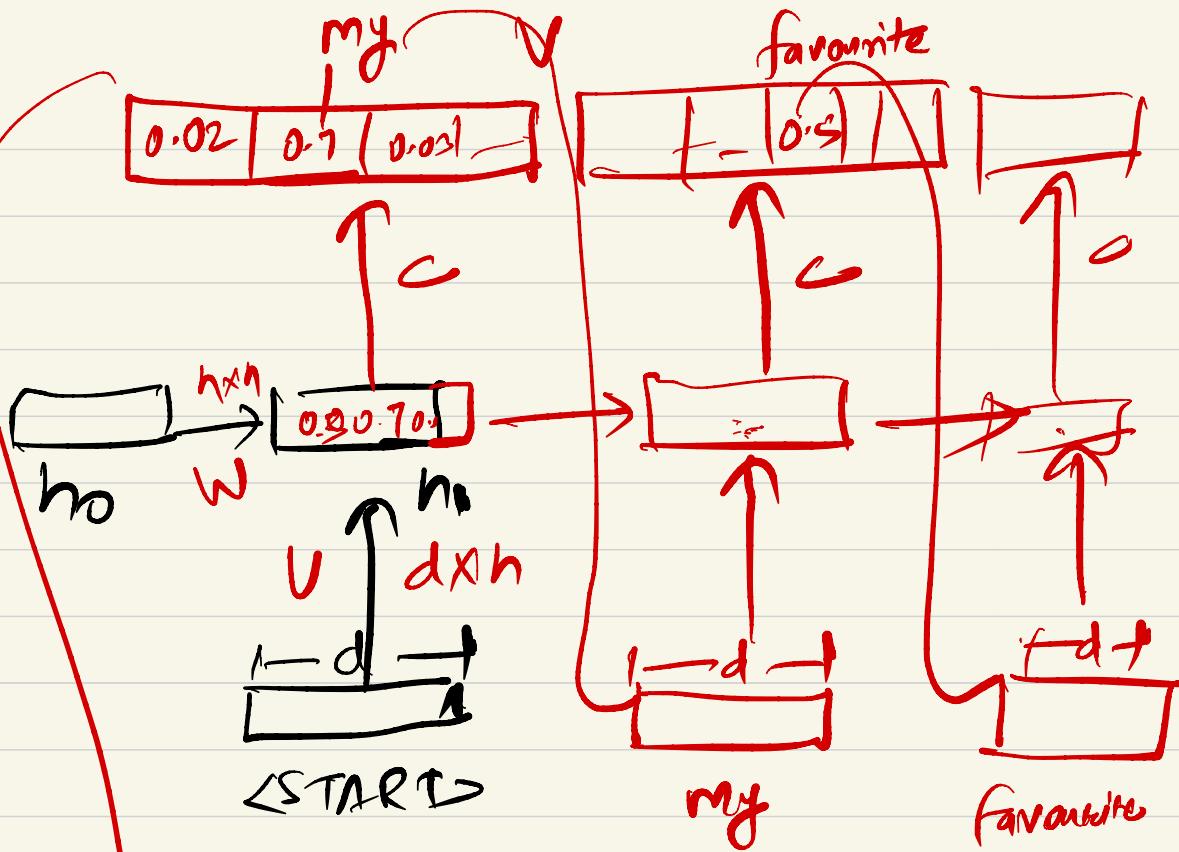
$$h_1 = \tanh(w_{h_0} + Ux_0)$$

$w_{h \times h}$      $U_d \times h$

BP TT

$$\frac{\partial L}{\partial w} \quad \frac{\partial L}{\partial b}$$

$L$      $w_{h \times h}$      $b_{h \times 1}$



Generation



of  
<STOP>

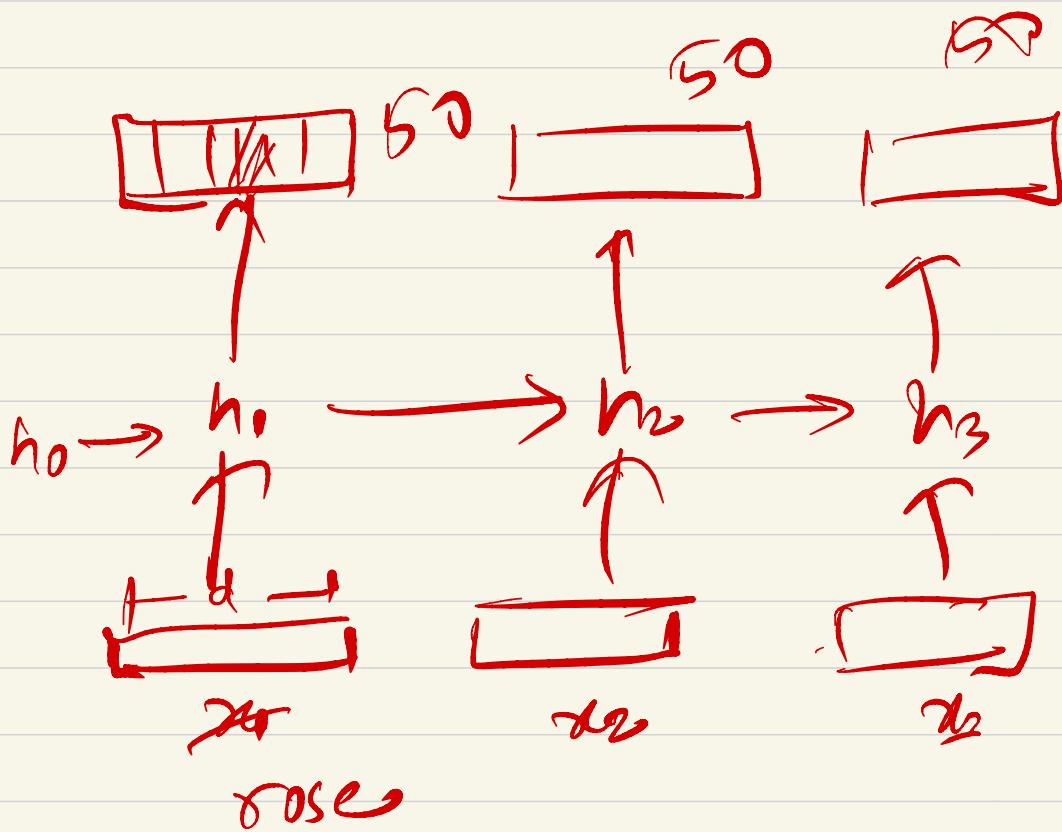


Sampling

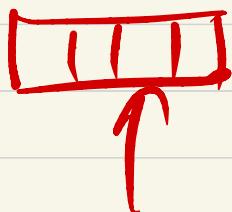
# Review for POS Tagging 50

The startled cat knocked - -

↓      ↓      ↓      ↓      - -

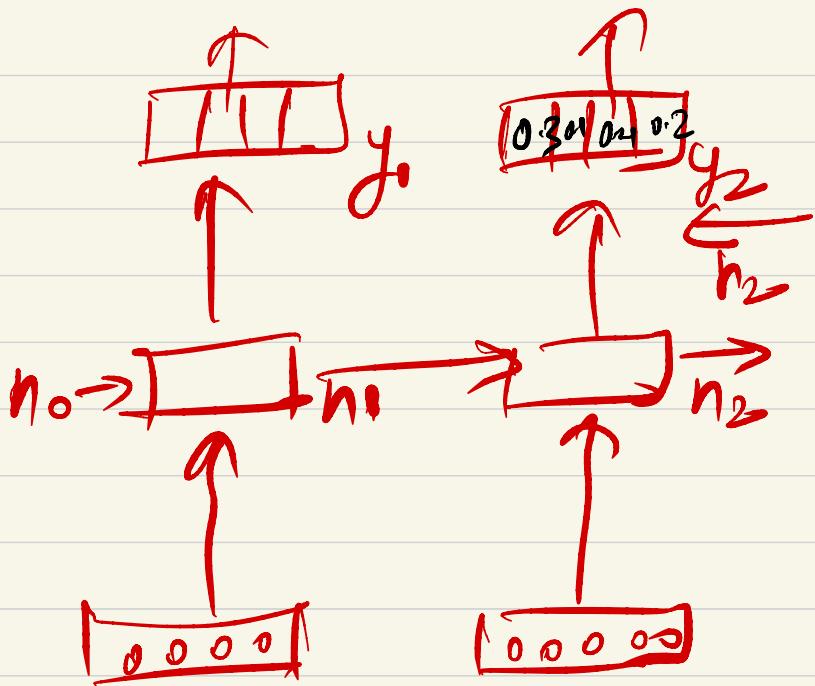


{ Name of a Person  
Location  
Organization  
Others + 1 }

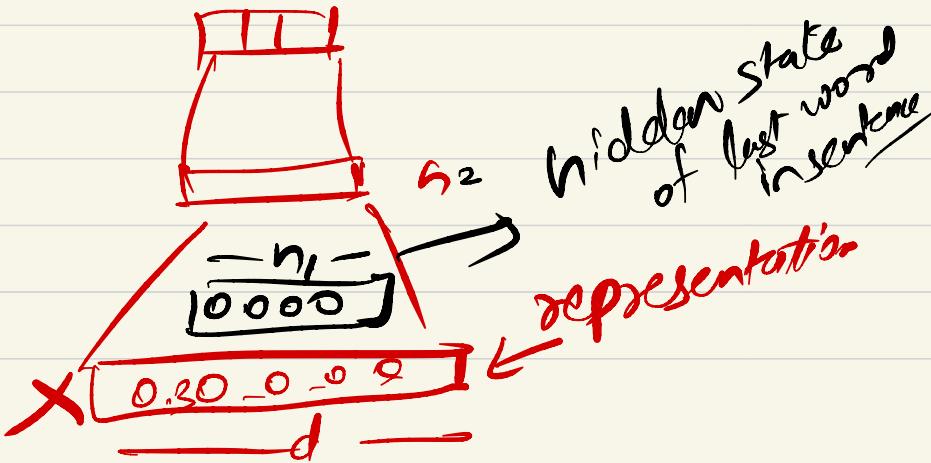


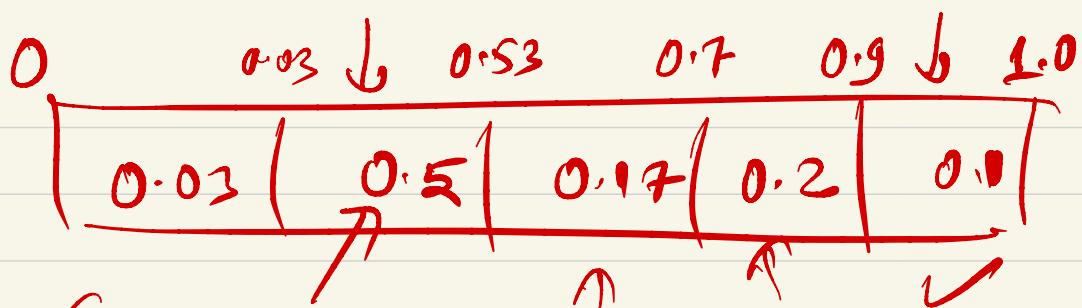
w<sub>1</sub>   w<sub>2</sub>   w<sub>3</sub>   ...

<sup>↑ Name</sup>   <sup>↑ Name</sup>  
Sachin   Tendulkar played  
a brilliant knock.



Sachin Tendulkar played a  
Overall I enjoyed the movie a lot

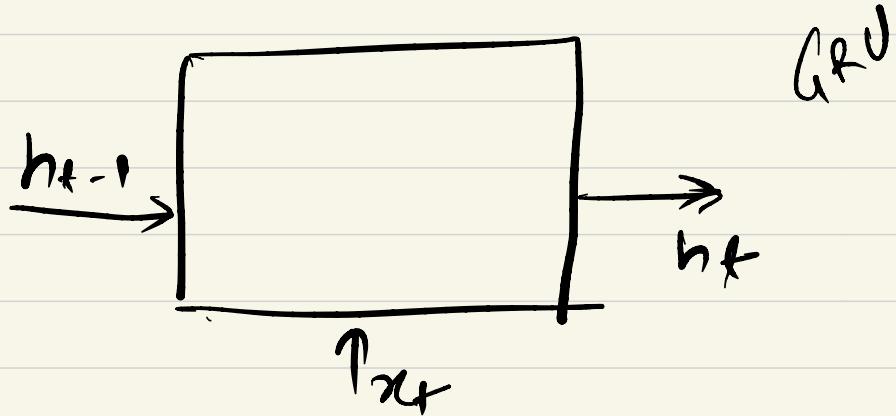




Generate a random number

between

$[0 - 1]$



Update gate  $u^{(t)} = \sigma(W_u h_{t-1} + U_u x_t + b_u)$

Reset gate  $r^{(t)} = \sigma(W_r h_{t-1} + U_r x_t + b_r)$

$$\tilde{h}_t = \tanh \left( W_n \underbrace{\left( r^{(t)} \cdot h^{(t-1)} \right)}_{\text{select useful parts of prev. hidden state}} + U_n x_t + b_n \right)$$

$$h_t = u_t \circ \tilde{h}_t + (1 - u_t) \circ h_{t-1}$$

$\rightarrow$  select useful parts of prev.  
hidden state

$$u_t = 0$$

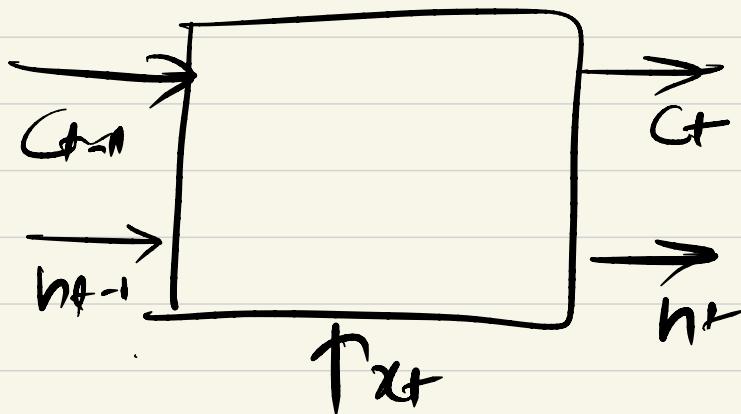
old context

What is kept from prev. hidden state & what is ignored.

$$\{r^{(t)} = 0 \Rightarrow \tilde{h}_t \rightarrow \text{only new content}$$

$$u_t = 1 \Rightarrow \text{only new}$$

# LSTM



forget gate      input gate

$$f_t = \sigma(w_f h_{t-1} + b_f x_t + b_f)$$

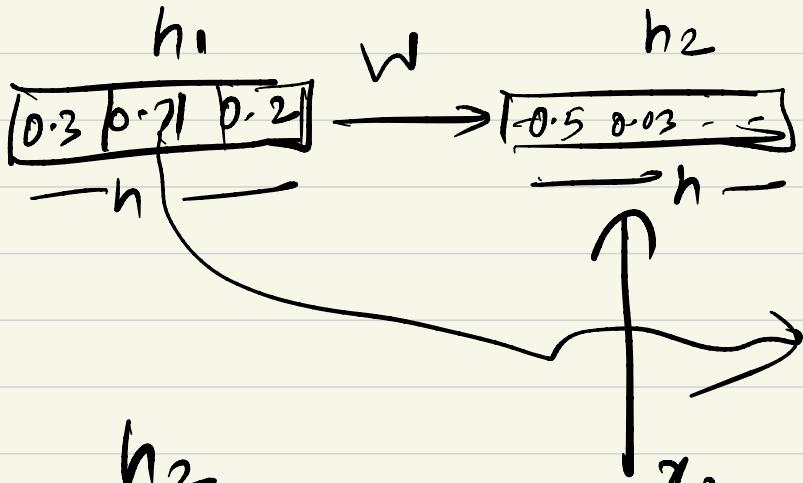
$$i_t = \sigma(w_i h_{t-1} + b_i x_t + b_i)$$

$$o_t = \sigma(w_o h_{t-1} + b_o x_t + b_o)$$

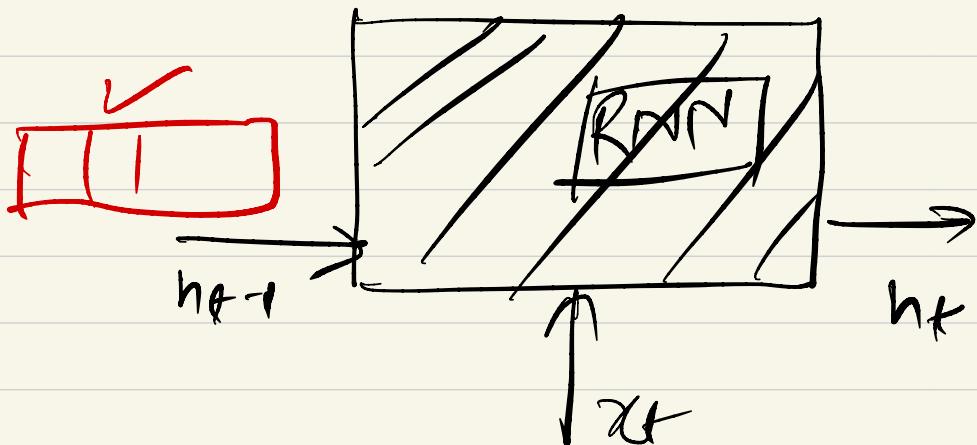
$$\tilde{C}_t = \tanh(w_c h^{(t-1)} + U_c x^{(t)} + b_c)$$

$$C_t = f^{(t)} \circ C^{(t-1)} + i^{(t)} \circ \tilde{C}^{(t)}$$

$$h^{(t)} = o^{(t)} \circ \tanh C^{(t)}$$

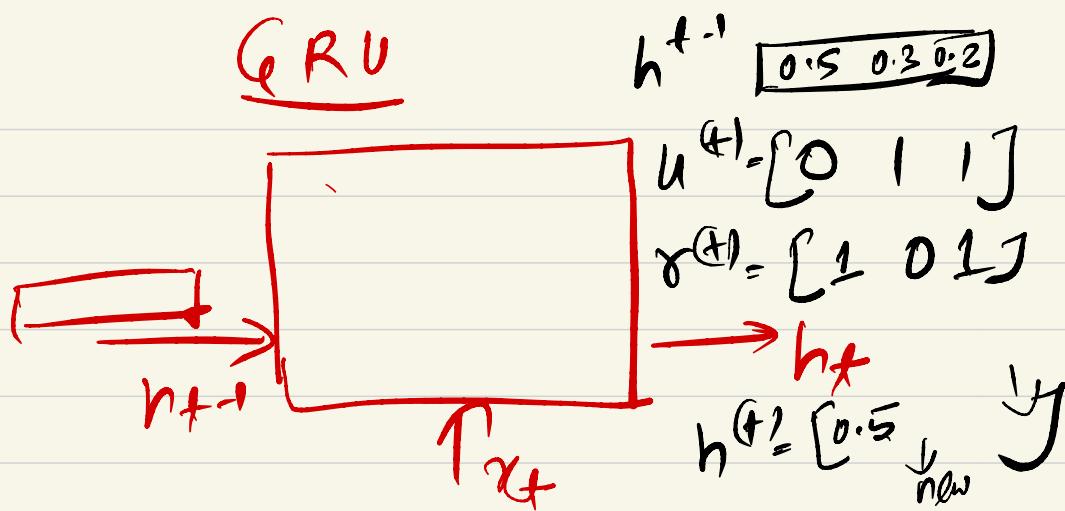


$$h_2 = \sigma(W h_1 + b x_2)$$



$$\text{Vanilla} - \sigma(U x_t + W h_{t-1})$$

# GRU

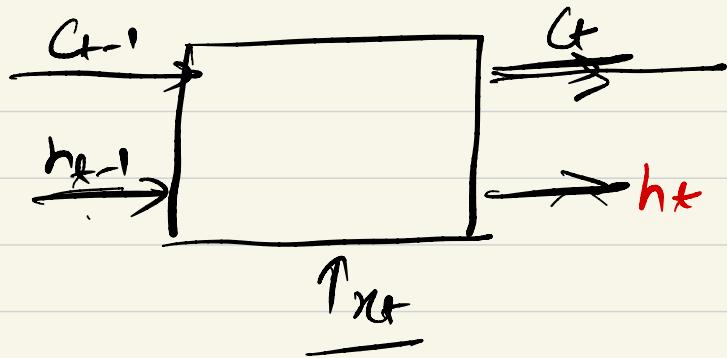


<u>Gates</u>	<u>Update</u>	<u>Reset</u>
$\begin{matrix} u^{(t)} \\ \downarrow \\ 0 \ 1 \end{matrix}$ <del><math>[0 \ 1 \ 0 \ 1 \ 0 \ 0 \ 1]</math></del>	$\sigma(W_u h_{t-1} + U_u x_t + b_u)$ $u^{(t)} = 0 \rightarrow \text{old info}$	$\sigma(U_r = 0 \quad U^{(t)} = 1)$ $\rightarrow \text{New info}$
	$\sigma(W_g h_{t-1} + U_g x_t + b_g)$	

$$\tilde{h}_t = \tanh(W_h(r^{(t)} \circ h_{t-1}) + U_h x_t + b_h)$$

$$r^{(t)} = [1 \ 0 \ 0 \ 1 \ 0 \ 1 \ 1]$$

$$h_t = (1 - u^{(t)}) \circ h_{t-1} + u^{(t)} \circ \tilde{h}_t$$



$$\begin{aligned}
 \text{Input Gate} & \quad i^{(t)} = \sigma(W_i h_{t-1} + U_i x_t + b_i) \\
 \text{forget Gate} & \quad f^{(t)} = \sigma(W_f h_{t-1} + U_f x_t + b_f) \\
 \text{Output Gate} & \quad o^{(t)} = \sigma(W_o h_{t-1} + U_o x_t + b_o)
 \end{aligned}$$

$$\tilde{C}^{(t)} = \tanh(W_C h_{t-1} + U_C x_t + b_C)$$

$$C^{(t)} = \underbrace{f^{(t)} \circ C_{t-1}}_{\text{Forget}} + \underbrace{i^{(t)} \circ \tilde{C}^{(t)}}_{\text{Update}}$$

$$h^{(t)} = o^{(t)} \circ \tanh(C^{(t)})$$

$$C^{(t)}$$

$$\boxed{\dots \downarrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow}$$

$$o^{(t)} \boxed{\begin{matrix} 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 \end{matrix}}$$

ELMO ✓

BERT ✓

LXNet} ↴

CS 224n

of NLP with

Deep Learning}

Retrieval Based

Chatbot

Face

Generative