

$$\frac{\partial L}{\partial x} = 0$$

$$y = (\underline{\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_4}) + b$$

$$Y = \overline{w_0 x + f(0)}$$

$$J = \frac{(Y - \hat{Y})^2}{2}$$

$$\downarrow P\omega_n = \omega_n - \mathcal{N} \left( \frac{\partial L}{\partial \omega_n} \right)$$

$$\omega_n = \omega_n - n \text{ (true)}$$

$$W_n = w_n - \eta(-ve)$$

$$= w_n + n (+ve)$$

$$\tilde{w}_y = w_y - n \frac{\partial L}{\partial w_y}$$

$$\frac{\partial L}{\partial \omega_y} = \frac{\partial f}{\partial y} - \frac{\partial \hat{y}}{\partial \omega_y}$$

$\frac{\partial L}{\partial w_n}$

$$\omega_4 = \underline{\underline{\omega}}$$

$$\frac{dl}{dw} < 0$$

$$\frac{\partial L}{\partial \hat{y}} = -2(y - \hat{y})$$

$$= -1 \begin{pmatrix} y \\ y \end{pmatrix}$$

$$\cancel{A} = \textcircled{A}$$

$$\hat{y} = w_4 f(0)$$

$$\frac{dy}{dx} = f(a)$$

$$w_1 = \dots$$

$$L = \frac{(y - \hat{y})^2}{2}$$

$$\omega_1 = \omega_1 - \eta \frac{\partial L}{\partial \omega_1}$$

$$\frac{\partial y}{\partial \omega_1} = \frac{\partial y}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \omega_1}$$

$$\frac{\partial \hat{y}}{\partial \omega_1} = \omega_1 f'(0)$$

$$\hat{y} = \omega_1 f(0)$$

$$o = \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3$$

$$\frac{\partial o}{\partial \omega_1} = x_1$$

$$f(x) = \frac{1}{1 + e^{-x}}$$

	Gender	Ticket Price	Age	Survived
A	M	£500	40	0
B				0
C				

$x_1 = \dots$   
 $x_2 = 0$   
 $x_3 = 0.67$

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

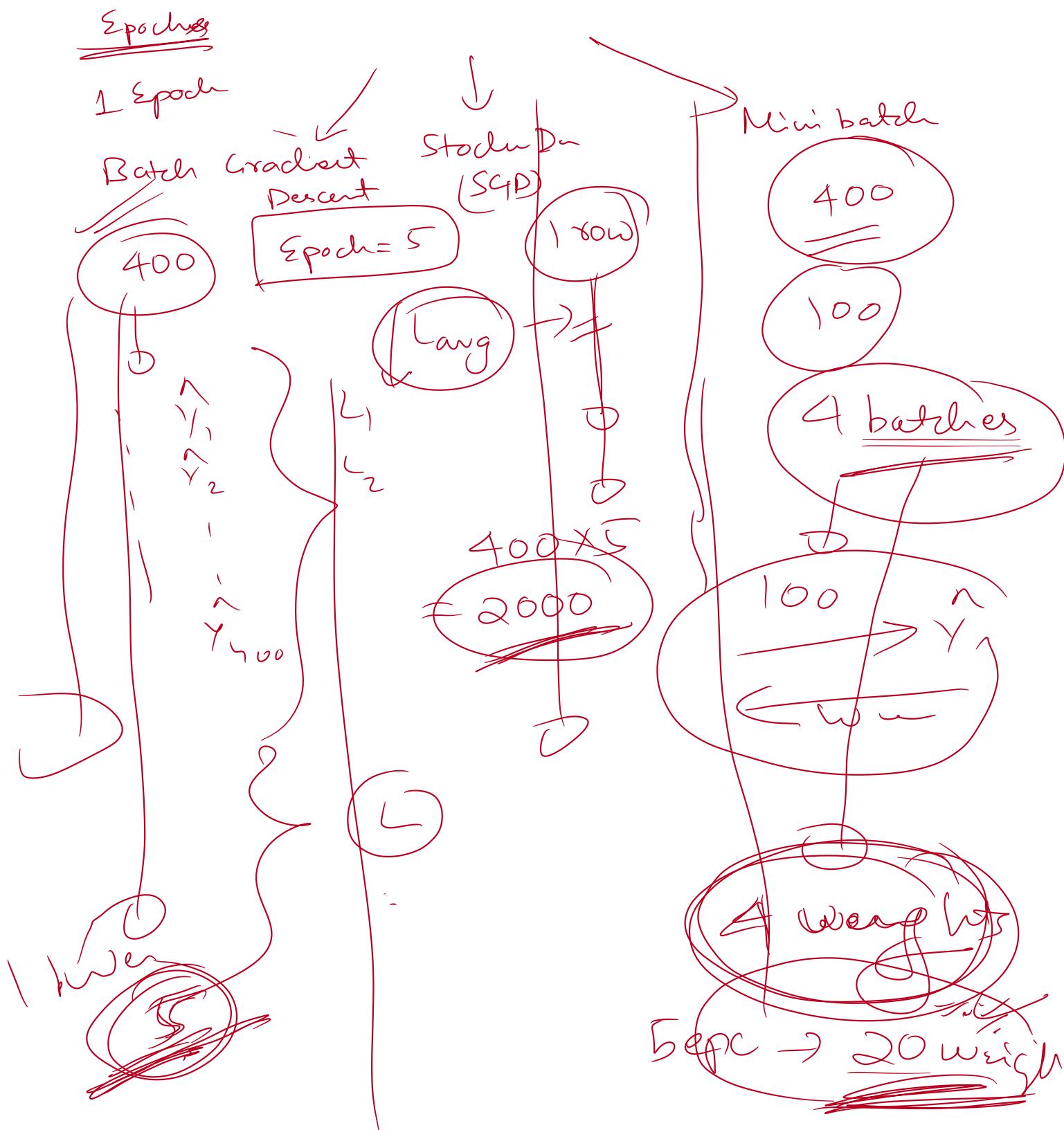
$$(1) L = -1 \log 0.67 - (1-0.67) \log(1-0.67)$$

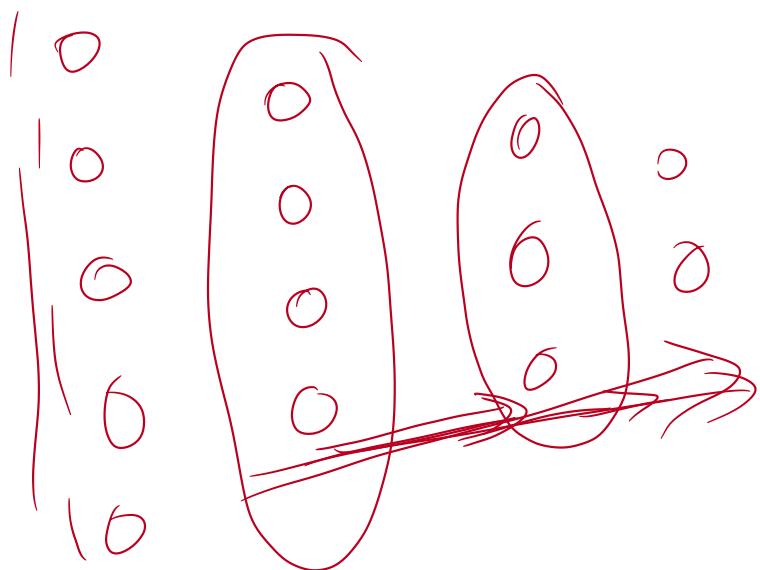
$$L = -\log 0.67$$

$$\phi = (\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3)$$

$$\hat{y} = \omega_0 f(0)$$

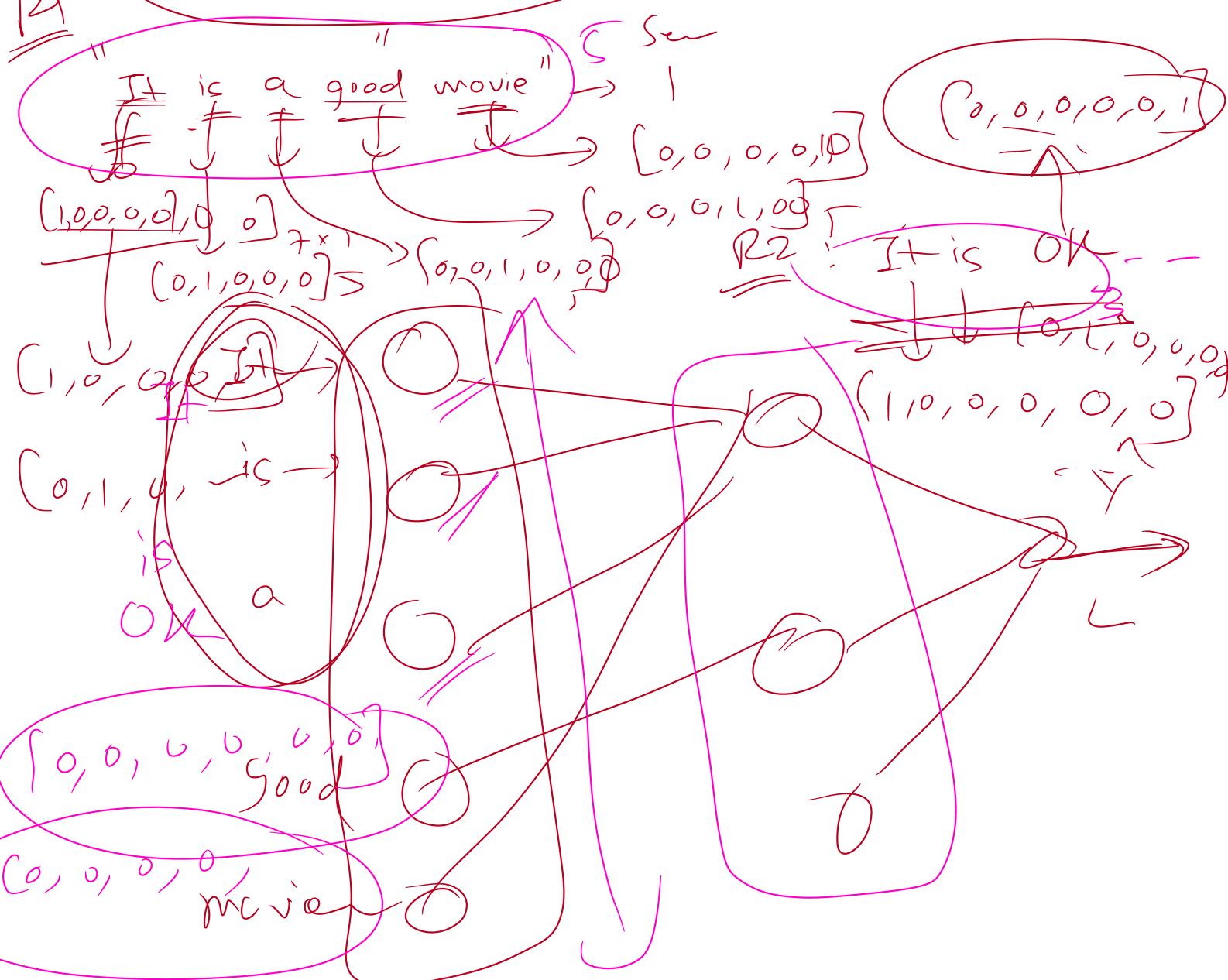
$$\hat{y} = f(\omega_0 (\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3))$$

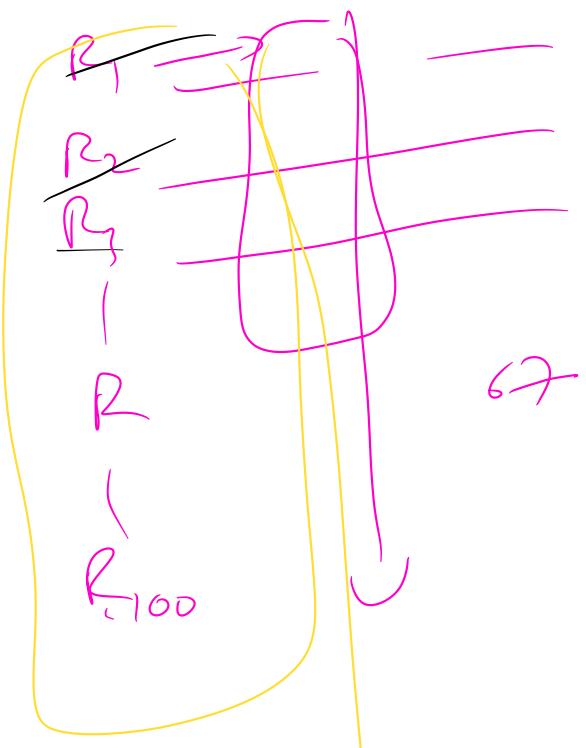




## Sentiment Analys

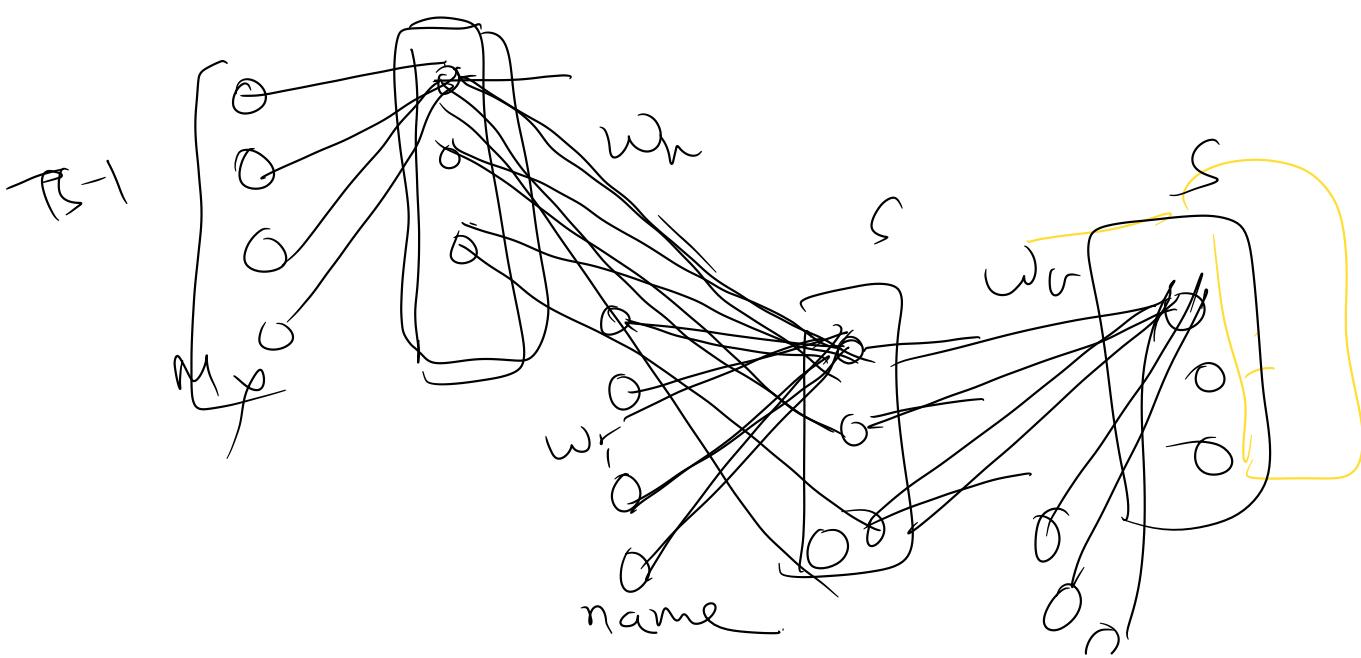
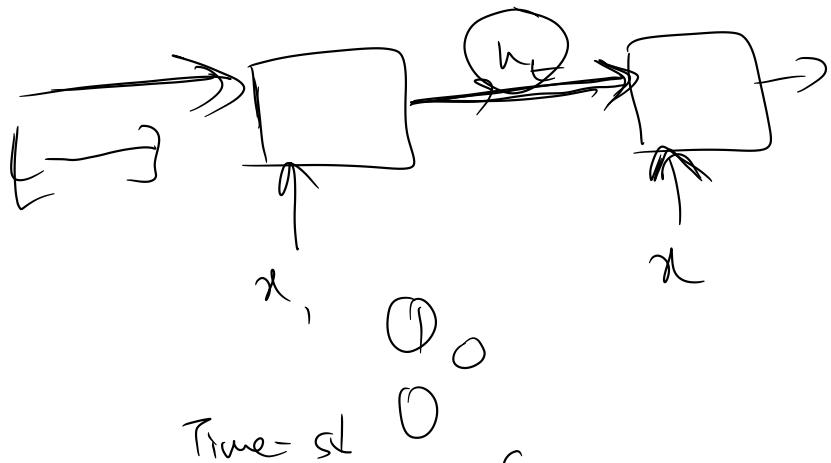
R1



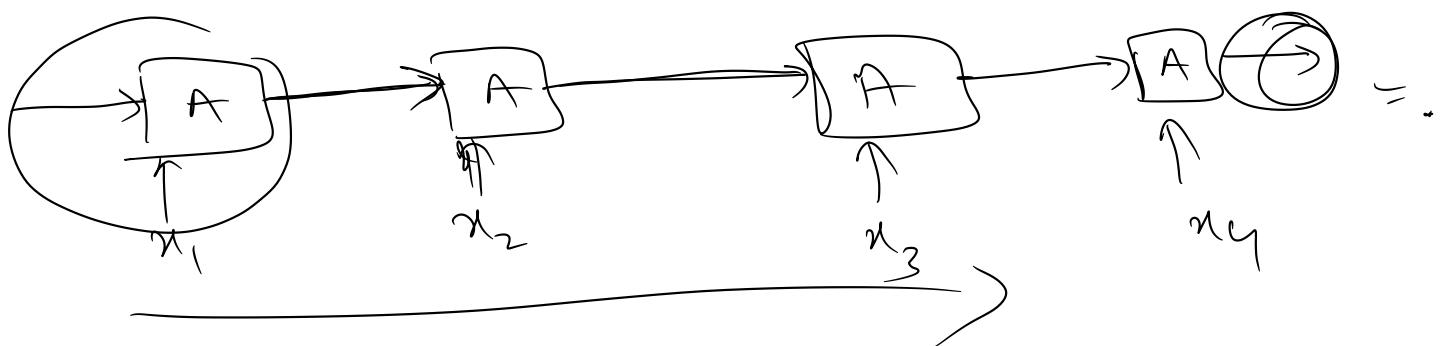
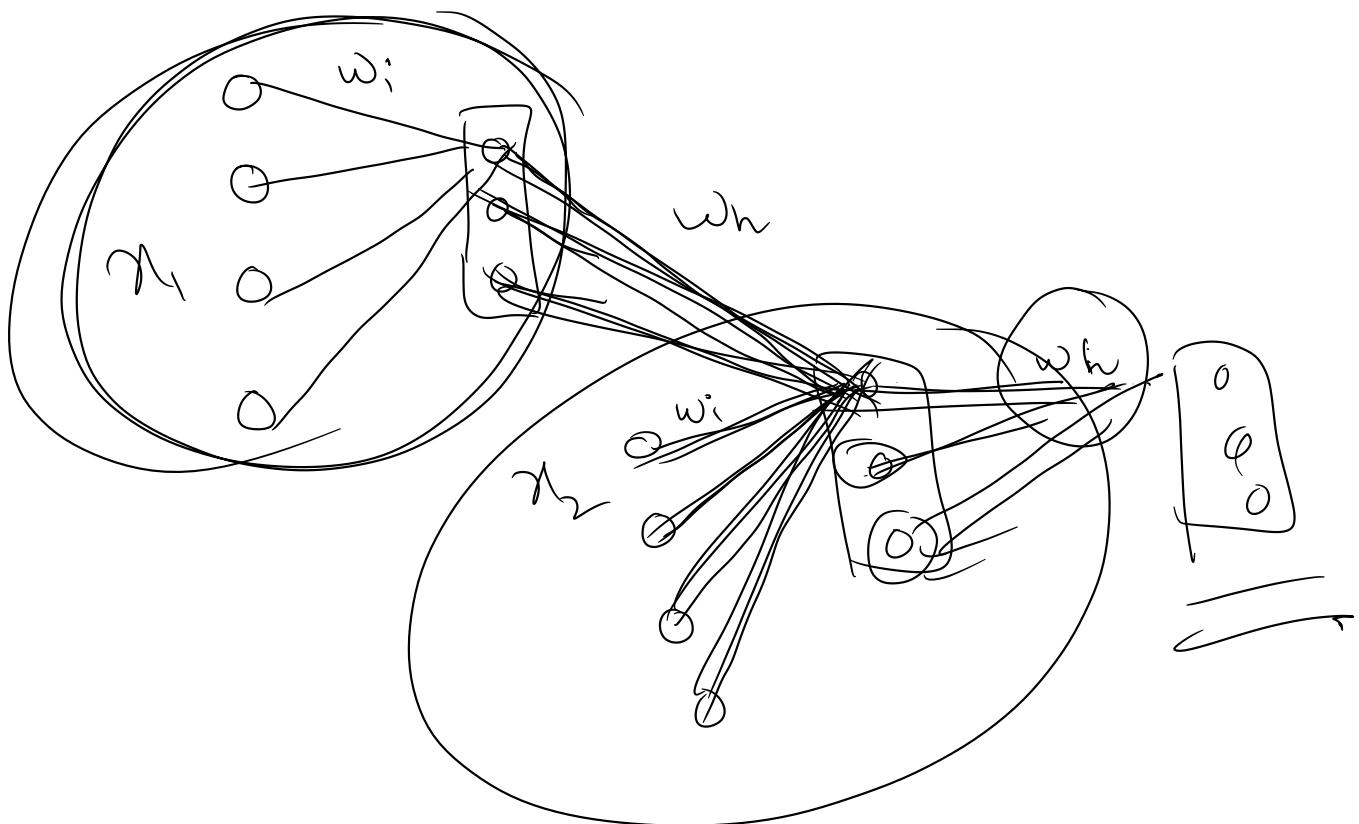


67

$(1, 0, 0, 0)$   $(0, 1, 0, 0)$   $[0, 0, 0, 1]$   
~~My~~ name is Baender  
 $x_1$   $x_2$   $x_3$   $x_4$   
 $(0, 0, 1, 0)$



my name is  
— T.S



## Sequential Neural Networks

Why to avoid ANNs for sequential data??

For Sentiment Analysis, let's say

the dataset has

Sentences like

$$\left. \begin{array}{l} S_1: \text{Bauder is studies at IITK} \\ S_2: \text{" " " " "} \\ S_3: \text{" " " " "} \end{array} \right\} \text{We assign vectors to all words.}$$

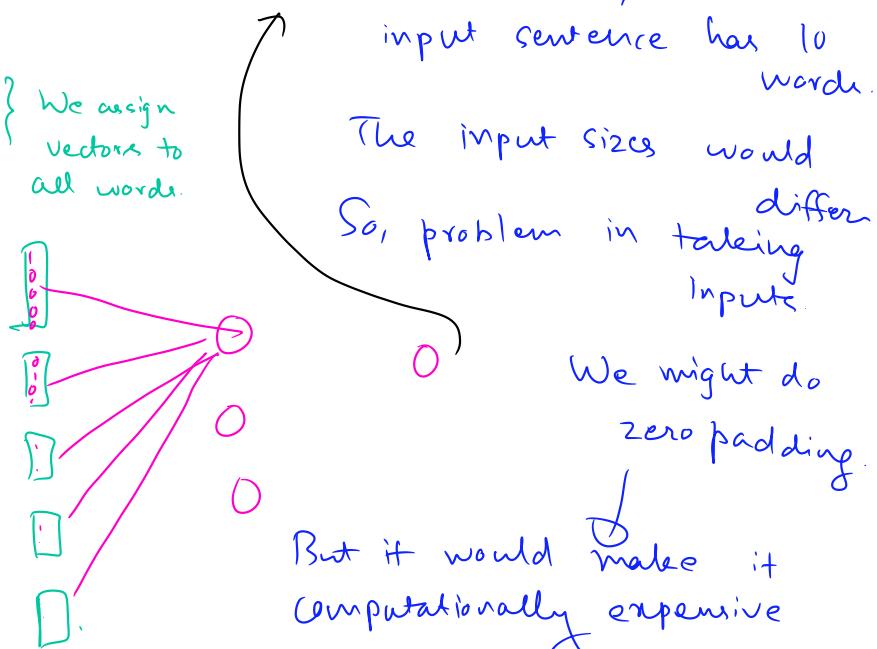
$$\text{Bauder} = [1 \ 0 \ 0 \ 0 \ 0]$$

$$\text{is} = [0 \ 1 \ 0 \ 0 \ 0]$$

$$\text{studies} = [0 \ 0 \ 1 \ 0 \ 0]$$

$$\text{at} = [0 \ 0 \ 0 \ 1 \ 0]$$

$$\text{IITK} = [0 \ 0 \ 0 \ 0 \ 1]$$



Why zero padding becomes an issue?? & sparsity in data increases

→ Dense layers expect dense inputs.

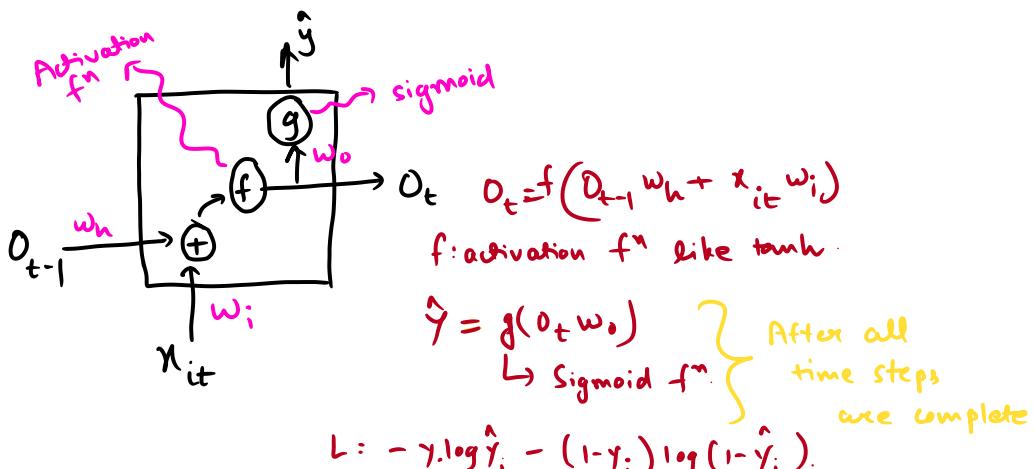
→ Sparsity of data ↑: reduces effective information per sample.

→ Zeros are treated as real data, not padding

→ Many weights never receive useful gradients.

(The corresponding weights barely get updated).

## Recurrent Neural Networks



" Recurrent" term is used because the hidden layers appear multiple times. The RNN "unfolds" through time.

" Back propagation through time" happens here : the algo is "unrolled" across time steps, allowing gradients (errors) to flow backward not just through layers but also through the sequence's history

## Why can't RNNs do machine translation?

An RNN can't produce an output, if we don't feed it an input at every time step.

Let's say we are doing English to Hindi  
Both the languages have different sentence structures.

So, we can't start translating right away at every time step.

We have to capture the meaning first.

So, we have to force our RNN to generate `<PAD>` or some dummy token till we feed the entire sentence into it..

So, the model doesn't have actual values to BPTT.

To force our RNN;

After final `<eos>` of English should output `<sos>` of French.

The model has to switch from Eng to French and we have to give the French sentence as input at every time step word by word.

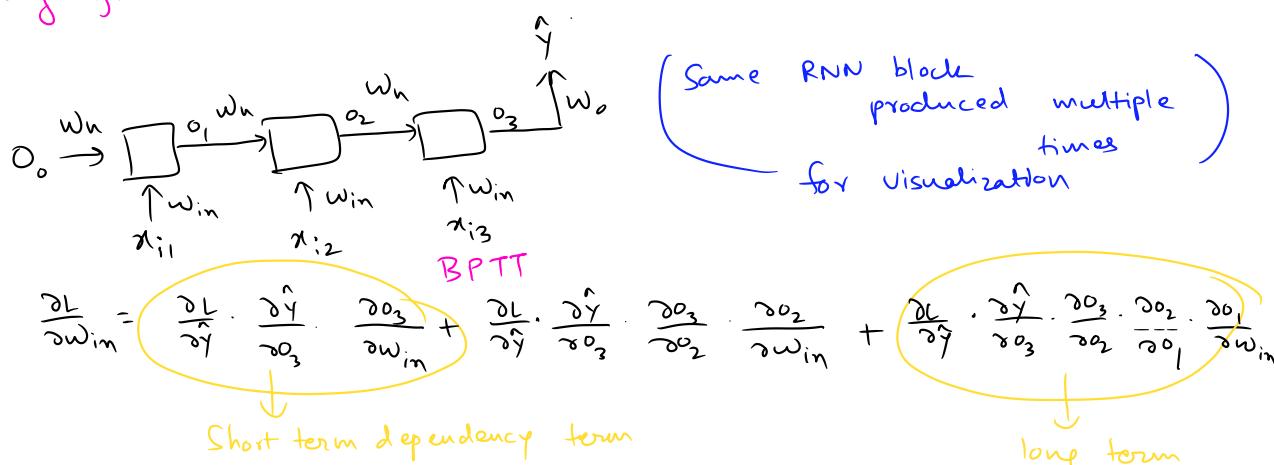
## Major Issues with RNNs

→ Problem of long term dependencies.

→ Stagnated training.

Long term dependencies arise due to Vanishing gradient problem.

↓  
Exploding gradients.



If we use activation  $f^n$  like

tanh or sigmoid

Vanishing gradient → weights may remain near their initialization

Sol<sup>n</sup>: Relu or Leaky Relu  
Better weight init

If activation  $f^n$  like

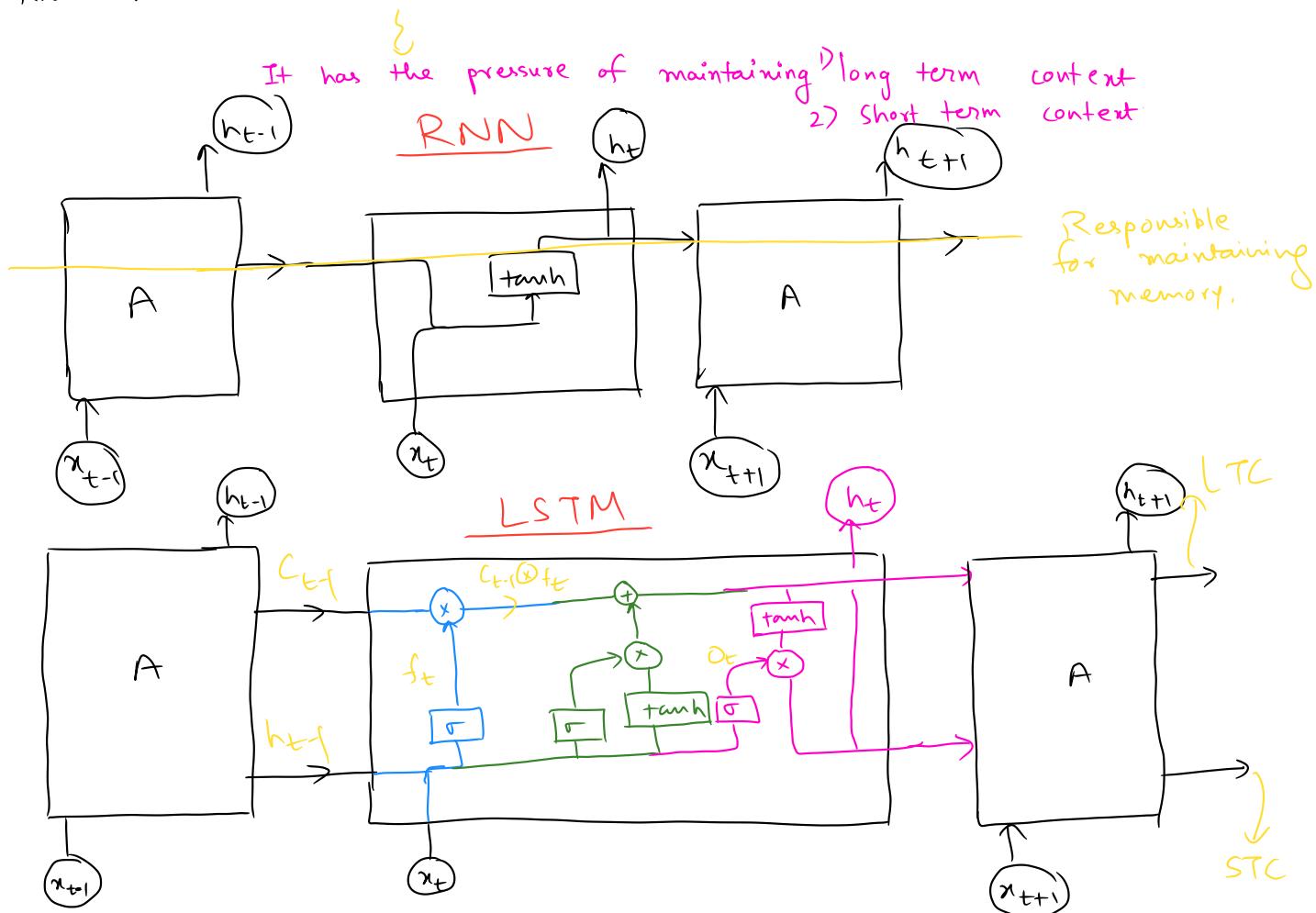
relu or leaky relu

exploding gradient → chaotic gradient descent, weights become very large

Sol<sup>n</sup>: Gradient Clipping (we set the max<sup>m</sup> gradient)  
Controlled Learning Rate.

## Long Short-Term Memory (LSTM)

RNN maintains memory only by the hidden state line.



### Forget gate

What to remove  
from LTM based  
on STC & current  
input

### Input gate

What to add to  
LTM

### Output gate

Based on LTM  
& current  
input gives  
output

also creates  
STM

LTC :  $c_{t-1}$  Cell State

Total inputs

STC :  $h_{t-1}$  hidden state

Input :  $x_t$

given  $c_t$  &  $h_t$   
updated created

Terms

$x_t$        $c_{t-1}$        $h_{t-1}$   
one word converted  
into vector

vectors with same dimension.



(X) : point wise multiplication

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \otimes \begin{bmatrix} 3 & 4 & 5 \end{bmatrix} = \begin{bmatrix} 3 & 8 & 15 \end{bmatrix}$$

(+) : point wise addition

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \otimes \begin{bmatrix} 3 & 4 & 5 \end{bmatrix} = \begin{bmatrix} 4 & 6 & 8 \end{bmatrix}$$

(tanh) : tanh appl^n

$$[\tanh(1), \tanh(2), \tanh(3)]$$

$\sigma$  : Neural Networks with sigmoid activation  $f^n$

} Number of nodes same throughout LSTM

$\tanh$  : Neural N with tanh activation  $f^n$

Number of nodes, hyperparameter

Number of nodes in hidden layers = dimensions of  $c_t, h_{t-1}, f_t, \bar{c}_t, i_t, o_t$

### Forget Gate

Decided by  $x_t$  &  $h_{t-1}$

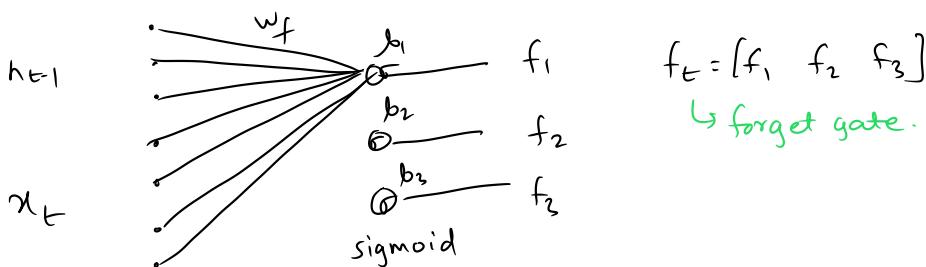
Let's say our  $x_t$  has dim=4

Let's No. of nodes' dim=3

Say dim  $(h_{t-1}) = 3$

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

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



### Input Gate

- 1) Cal<sup>c</sup> candidate cell state  $\tilde{c}_t$
- 2) Cal<sup>i</sup>  $i_t$   $\hookrightarrow$  decides which info of  $\tilde{c}_t$  should be added to  $c_t$
- 3)  $c_t$

$$\tilde{c}_t = \tanh(w_c [h_{t-1}, x_t] + b_c)$$

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

$$\tilde{c}_t \otimes i_t = \tilde{c}_t^* \quad (\text{filtered candidate state})$$

This solves Vanishing Gradient Problem

If LSTM feels the initial information is imp; then it would set

$$f_t = [1 \ 1 \ 1] \quad \text{example dim} = 3$$

$$i_t = [0 \ 0 \ 0] - \frac{\epsilon_{\text{max}}}{3}$$

$$\text{So, } c_{t-1} = c_t$$

$$c_t = f_t \otimes c_{t-1} + i_t \otimes \tilde{c}_t \rightarrow \text{Cell State}$$

### Output Gate

$$h_t = o_t \otimes \tanh(c_t)$$

$\tanh$  keeps a safe range of  $(-1, 1]$

- $\tanh$  is used because it keeps the hidden activations bounded
- stabilizes training by preventing exploding gradients.
- allows gates to control information flow
- exposes only a safe, filtered version of internal memory.

### Gated Recurrent Units

LSTM: complex architecture, more number of parameters.

↳ 3 gates (Training Time more)

GRU: simpler architecture, lesser parameters.

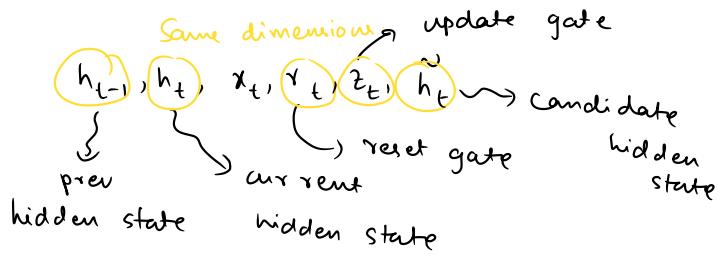
2 gates ↳ reduces training time

LSTM: STC & LTC is maintained separately

GRU: Only one hidden state

↳ Reset Gate & Update Gate.

We give  $h_{t-1}$  &  $x_t$  as input.

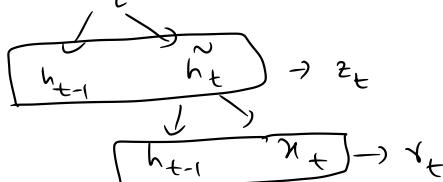


$\sigma$ ,  $\tanh$  have same number of nodes = dim ( $h_t, r_t \dots$ )

We take text, we use text vectorization techniques

- Steps
- 1) calc.  $r_t$
  - 2)  $\tilde{h}_t$
  - 3)  $z_t$
  - 4)  $h_t$

We need  $h_t$



from prev. hidden state  $(h_{t-1})$  & current input  $(x_t)$ : we form candidate  $(\tilde{h}_t)$  memory

$\tilde{h}_t$  is inclined towards to current input

$$h_{t-1}, \tilde{h}_t \text{ (Balanced)} \rightarrow h_t$$

$$h_{t-1} \xrightarrow{r_t} \tilde{h}_t \xrightarrow{z_t} h_t$$

Overall flowchart

$$\text{reset gate} \\ r_t = \sigma(w_r [h_{t-1}, x_t] + b_r)$$

$$\tilde{h}_t = \tanh(w_c [h_{t-1} \otimes r_t, x_t] + b_c)$$

$$h_t = (1 - z_t) \odot h_{t-1} \oplus z_t \odot \tilde{h}_t$$

$z_t$  is high it favours candidate hidden state

$z_t$  is low: favours prev. hidden state

### Encoder - Decoder

For NLP tasks like machine translation

where input & output have varying lengths.

Let's say English to Hindi

English: I want to ~~really~~ back home *+ tokens*

Hindi: मुझे घर जाना है <start> <end>

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

6 tokens

(We calc error across all time steps)

for BProp

