# Backpropagation in RNN

We gonna to take an example → Many to one RNN

| text | | | | x | | | N |
|------|------|------|---|---|---|---|---|
| cat | maat | crat | 1 | $x_1$ [100] [010] [001] | | | 1 |
| crat | rat | mat | 1 | → $x_2$ [001] [001] [010] | | | 1 |
| mat | mat | cat | 0 | $x_3$ [010] [010] [100] | | | 0 |

Vocab → Cat  Mat  Rat

$$[100] \qquad [010] \qquad [001]$$



Wi
(3,3)

Wh
(3×3)

3 bias

Wo (3,1)

1 bias



→ slot of zeros

$O_0 \xrightarrow{W_h} \square \xrightarrow{O_1} \square \xrightarrow{O_2} \square \to O_3$

$X_{11}$   $X_{12}$   $X_{13}$

$\hat{Y}$   Wo

$$O_1 = t(x_{11} w_i + O_0 w_h)$$

$$O_2 = t(x_{12} w_i + O_1 w_h)$$

$$O_3 = f(x_{13} w_i + O_2 w_h)$$

$$\hat{Y} = \sigma(O_3 w_0)$$

$$Loss = -y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y})$$

Loss calculate → minimize

using gradient ↙
descent
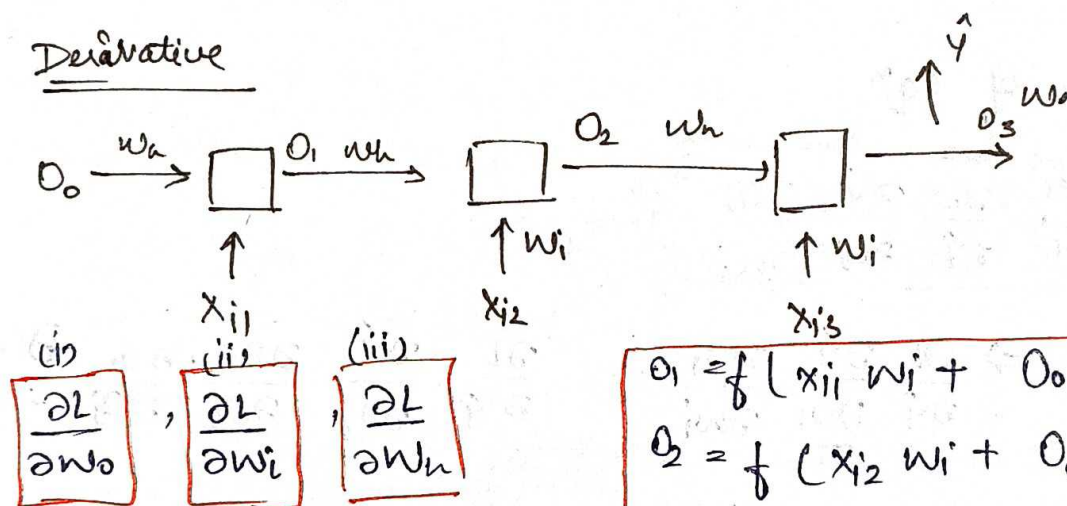
find the value of Wi wh Wo where Loss
should be minimum.

$$W_i = W_i - \eta \boxed{\dfrac{\partial L}{\partial W_i}} \qquad W_o = W_o - \eta \boxed{\dfrac{\partial L}{\partial W_o}}$$

$$W_h = W_h - \eta \boxed{\dfrac{\partial L}{\partial W_h}} \longleftarrow \text{Now, we have to find the derivative}$$

## Derivative



$$O_0 \xrightarrow{W_h} \square \xrightarrow{O_1 \; wh} \square \xrightarrow{O_2 \; wh} \square \xrightarrow{O_3 \; wo} \hat{y}$$

$$\uparrow W_i \qquad \uparrow W_i$$

$$X_{i1} \qquad X_{i2} \qquad X_{i3}$$

(i) $\boxed{\dfrac{\partial L}{\partial W_o}}$ , (ii) $\boxed{\dfrac{\partial L}{\partial W_i}}$ (iii) $\boxed{\dfrac{\partial L}{\partial W_h}}$

$$\boxed{\begin{aligned} O_1 &= f(X_{i1} W_i + O_0 W_h) \\ O_2 &= f(X_{i2} W_i + O_0 W_h) \\ O_3 &= f(X_{i3} W_i + O_0 W_h) \\ \hat{y} &= \sigma(O_3 W_o) \end{aligned}}$$
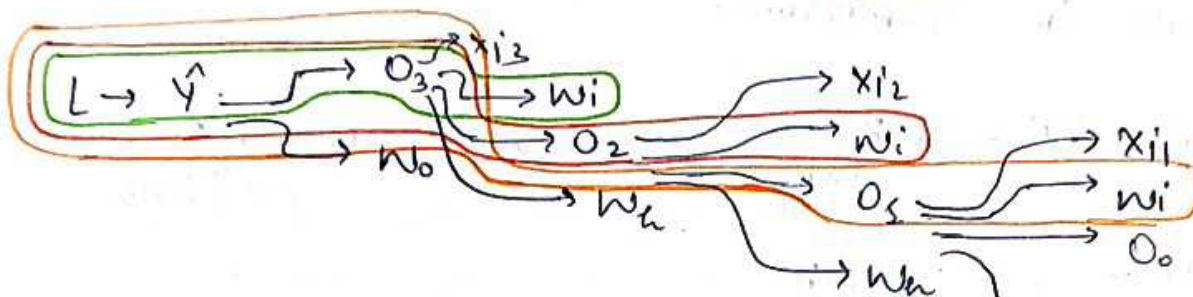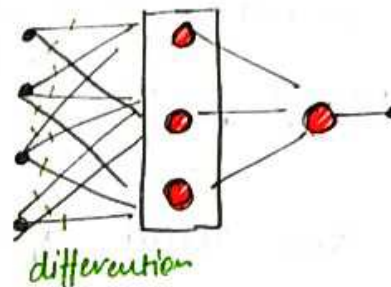
$$\dfrac{\partial L}{\partial W_o} \qquad L \to \hat{y} \begin{cases} O_3 \\ W_o \end{cases}$$

(i) $\boxed{\dfrac{\partial L}{\partial W_o} = \boxed{\dfrac{\partial L}{\partial \hat{y}}} \boxed{\dfrac{\partial \hat{y}}{\partial W_o}}} \longrightarrow L = Y_i \log \hat{Y}_i - (1 - Y_i) \log(1 - \hat{y})$

↓  ↑

$$\hat{y} = \sigma(O_3 W_o) \longleftarrow \boxed{\text{differentiation}}$$

(ii) $\dfrac{\partial L}{\partial W_i} \longrightarrow$  $\partial W_i \Rightarrow$ change

$L \Rightarrow$ how much  $L$ change when $W_i$ change

$$\frac{\partial L}{\partial w_i} \rightarrow \quad L \rightarrow \hat{Y} \rightarrow O_3 \begin{array}{l} \rightarrow X_{i3} \\ \rightarrow w_i' \\ \rightarrow O_2 \\ \rightarrow w_h \end{array}$$

$$\rightarrow W_0$$



differentiation



$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{Y}} \frac{\partial \hat{Y}}{\partial O_3} \frac{\partial O_3}{\partial w_i} + \frac{\partial L}{\partial \hat{Y}} \frac{\partial \hat{Y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial w_i} + \frac{\partial L}{\partial \hat{Y}} \frac{\partial \hat{Y}}{\partial O_3} \frac{\partial O_3}{\partial O_2}$$

$$\frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial w_i}$$

## Summarization of eq$^n$
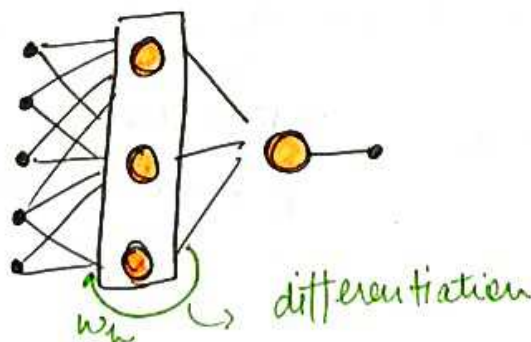
because we have 3 words

$$\boxed{\frac{\partial L}{\partial w_i} = \sum_{j=1}^{t=3} \frac{\partial L}{\partial \hat{Y}} \frac{\partial \hat{Y}}{\partial O_j} \frac{\partial O_j}{\partial w_i}}$$
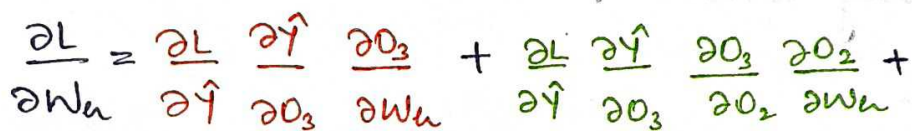
Proof ↗

Same

expand $j=1$ $\Rightarrow \dfrac{\partial L}{\partial \hat{Y}} \dfrac{\partial \hat{Y}}{\partial O_1} \dfrac{\partial O_1}{\partial w_i} \Rightarrow \dfrac{\partial L}{\partial \hat{Y}} \dfrac{\partial \hat{Y}}{\partial O_3} \dfrac{\partial O_3}{\partial O_2} \dfrac{\partial O_2}{\partial O_1} \dfrac{\partial O_1}{\partial w_i}$

Same

expand $j=2$ $\Rightarrow \dfrac{\partial L}{\partial \hat{Y}} \dfrac{\partial \hat{Y}}{\partial O_2} \dfrac{\partial O_2}{\partial w_i} \Rightarrow \dfrac{\partial L}{\partial \hat{Y}} \dfrac{\partial \hat{Y}}{\partial O_3} \dfrac{\partial O_3}{\partial O_2} \dfrac{\partial O_2}{\partial w_i}$

expand $j=3$ $\Rightarrow \boxed{\dfrac{\partial L}{\partial \hat{Y}} \dfrac{\partial \hat{Y}}{\partial O_3} \dfrac{\partial O_3}{\partial w_i}}$

(iii) $\quad \dfrac{\partial L}{\partial w_h}$



$w_h \quad \longrightarrow$ differentiation

$$\boxed{\frac{\partial L}{\partial W_a}}$$



$$\frac{\partial L}{\partial W_a} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial W_a} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial W_a} +$$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial W_a}$$

Summarize

$$\frac{\partial L}{\partial W_a} = \sum_{j=1}^{u} \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_j} \frac{\partial O_j}{\partial W_a}$$
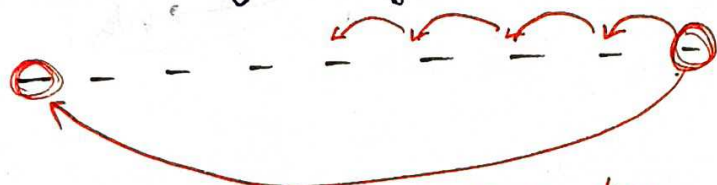
$u = $ timesteps

# Problem with RNN

RNN → Sequential data → textual, time, series
↓
Suffer ② major problem → LSTM
↳ problem of long term dependence
↳ unstable gradients / Stagnated Training

(i) Problem of long term dependence



\* short term memory loss.

This data not remember
Starting word (just like yayani)

g:-   Marathi is spoken in Maharashtra.
                                remember

Maharashtra is a beautiful place - went there
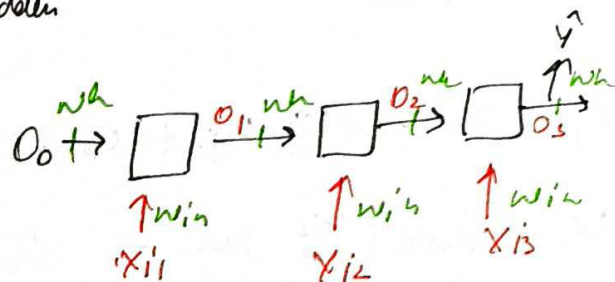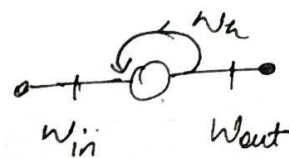last year but I could not enjoy properly
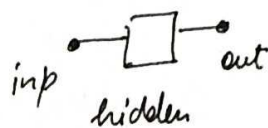because I don't understant marathi
              Not remember.

In Auto-suggestion (keypad) →   Suggest short sentence
                                but not long sentence

\* 2nd sentence RNN not remember "Maharashtra"
word. So, model not predict "Marathi" word
at last.      (Vanishing gradient problem)

# Problem #1 Problem of long term dependency → vanishing

| Input | | | out |
|---|---|---|---|
| 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 |

inp ▢ out
hidden

$W_{in}$ — $\circlearrowleft W_h$ — $W_{out}$

$O_0 \xrightarrow{W_h} \boxed{\phantom{x}} \xrightarrow[O_1]{W_h} \boxed{\phantom{x}} \xrightarrow[O_2]{W_h} \boxed{\phantom{x}} \xrightarrow[O_3]{W_h} \hat{y}$

$\uparrow W_{in}$    $\uparrow W_{in}$    $\uparrow W_{in}$
$x_{i1}$    $x_{i2}$    $x_{i3}$

3 time stamps

Loss → minimum → gradient descent

short term dependency

$$\frac{\partial L}{\partial W_{in}} = \boxed{\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial W_{in}}} +$$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial W_{in}} +$$

$$\boxed{\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial W_{in}}} \rightarrow \text{Long term dependency}$$

$$W_{in} = W_{in} - \eta \boxed{\frac{\partial L}{\partial W_{in}}}$$

$$W_{out} = W_{out} - \eta \boxed{\frac{\partial L}{\partial W_{out}}}$$

$$W_h = W_h - \eta \boxed{\frac{\partial L}{\partial W_h}}$$

This is only for 3 timesteps

# 100 timesteps
long term dependency

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{100}} \frac{\partial O_{100}}{\partial O_{99}} - - - - - - - - \frac{\partial O_2}{\partial O_2} \frac{\partial O_1}{\partial W_{in}}$$

During the calculating gradient descent, → value of long term dependency will be small

value of short term dependency will be large.

So, short term dependency contribute more as compare to long term memory to find gradient descent.

time stamps $\uparrow\uparrow$ $\Leftrightarrow$ long term dependency $\downarrow\downarrow$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{100}} \boxed{\frac{\partial O_{100}}{\partial O_{99}} \quad \text{-----} \quad \frac{\partial O_2}{\partial O_1}} \frac{\partial O_1}{\partial W_{in}}$$

$O_1 = \tanh (X_{11} W_{in} + O_0 W_h)$

$O_t = \tanh (X_{it} W_{in} + O_{t-1} W_h)$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{100}} \boxed{\prod_{t=2}^{100} \left(\frac{\partial O_t}{\partial O_{t-1}}\right)} \frac{\partial O_1}{\partial W_{in}}$$

$\dfrac{\partial O_t}{\partial O_{t-1}} = \tanh (X_{it} W_{in} + O_{t-1} W_h) W_h$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_{100}} \prod_{t=2}^{100} (\underbrace{\tanh' u}_{\text{bet}^n 0\rightarrow 1} \quad \underbrace{W_0}_{\text{bet}^n 0-1}) \frac{\partial O_1}{\partial W_{in}} \rightarrow$ Very close to zero

Sol $\rightarrow$ Different activation function $\rightarrow$ relu / leaky

not bet$^n$ 0-1

$\rightarrow$ better weight initialization

$\rightarrow$ Skip RNNs (dry) $\rightarrow$ self study topic

$\rightarrow$ LSTM

#2 Problem $\rightarrow$ Unstable Training (Exploding problem)

Long term dependency $\leftrightarrow$ very large

$\searrow \approx$ infinite

* if long term dep. is $\approx$ infinite then Gradient will be infinite then weights also infinite And Model not train.

for example → (i) Using relu and weight intialize
is 1. So, weights multiply and get
large number. Long Term dep → dominate and
short term dep. not contribute in gradient update.

(ii) If Learning Rate is ↑↑ (very heigher).

Solution : 1) <mark>Gradient Clipply</mark> → Search (self study)

2) Controlled Learning Rate

3) LSTM use.

## LSTM (Long Short Term Memory)

The what

LSTM core idea

Now we have to decide this        Story
story is good/Bad. decision
depend on geo-Area. Hero
of the story.

Pratapgarh

Vikram ——Fought→ enemy
died

Jr. Vikram
died

Super Junior Vikram
died

Super Jr.
vikram
xyz
LTC

Vikram
xyz
STC