# Course 5:-
## RNN

Regular NN: A regular network takes out input and gives output and then forget it.

RNNs are different they remember things over time

Key tricks: An RNN has a layers that gets passed from one step to the next, carrying information from the past. (hidden state)

$$\text{tanh/Relu} \rightarrow a^{[1]} = g(W_{aa}a^{[0]} + W_{ax}x^{[1]} + b_a]$$

$$\hat{y}^{[1]} = g_a(W_{ya}a^{[1]} + b_y) \quad \text{sigmoid}$$

Normal equation:

$$a^{t} = g(W_y \, a^{(t-1)} + W_{ax}x^{[t]} + b_a)$$

$$y = g(W \, a^{[t]} + b_y)$$

$$a^{[t]} = g(W_a [a^{[t-1]}, x^{[t]}] + b_a)$$

## Backpropagation through time in RNN:

Forward propagation

$$a^{[0]} \longrightarrow a^{[1]} \longrightarrow a^{[2]} \longrightarrow a^{3} \ldots \ldots a^{[T]}$$

$$x^{[1]} \quad x^{[2]} \quad x^{[3]} \quad x$$

Loss

$$\mathcal{L}^{[t]}(\hat{y}^{<t>}, y^{<t>}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{<t>}) \log(1-\hat{y}^{(t)})$$

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1} \mathcal{L}^{[t]}(\hat{y}^{[t]}, y^{(t)})$$

## # Types of RNNs

\# One to one (One input, one output)

   e.g Image → class

☆ One to many (One input, many output) Music Generation

   ~~One in~~    Image ⟹ [Word 1], [Word 2], [word 3]

☆ Many to One (Many Inputs, one output)

   Word + Word + Word ⟹ [label]

* Many to many (Same length)

Word 1 → Word2 → Word3 ⟹ Tag1, Tag2, Tag3

* Many to many (different length)

Input: Word 1 → Word 2 → Word 3

| encoded Context |

↓

output

[Mot 1] → [Mot 2] → [Mot 3] → [Mot 4]

e.g

English to urdu translation

# LANGUAG MODEL ond RNNS.

Suppose: I said The apple and pear salad
Language model: what did he say?

$P(\text{The apple and pear salad}) = 3.2 \times 10^{-13}$

$P(\text{The apple and pair salad}) = 5.7 \times 10^{-10}$

Solution:

Language model select sentence with greater probability

How it works

Training set: large corpus of english text

tokenize

Cats average 15 hours of sleep a day.
(tokenize it)

RNN as a model language:-

A language model predicts the probability of the next token (word, subword or character) given the previous one.

Suppose: we have a large language model
How it will correct the sentences
Cats average 15 hours of sleep a day
RNN tokenize this senctence:
⇒ and measure probability of every token (word, subword)
⇒ It then arrange or predict words with best probabity.

like:

I → love = proba (0.9) ✓
→ kill = proba (0.5)
→ nonsense = proba (0.1)

then

| I Love | Answer.

Loss function:

$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = -\sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$\mathcal{L} = \sum_t \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

To Sample Novel Sequence:-

$$a^{<0>} \rightarrow \boxed{a^{<1>}} \rightarrow \boxed{a^{<2>}} \rightarrow \boxed{a^{<3>}} \rightarrow \boxed{a^{<T_y>}}$$

$$\hat{y}^{<T_y>}$$

$$y^{<T_x - 1>}$$

Generate a sentence

I    you    love

$$\hat{y} = I \qquad \hat{y} + \hat{y}^1 = love \qquad \hat{y}, \hat{y} + \hat{y}^{\hat{}} = I + love y$$

$$a \rightarrow a^{<1>} \rightarrow a^{[2]} \rightarrow a^{\{3\}}$$

0.9    0.2 → 0.5       love 0.1    I = 0.1
I      You  love       0.9  you

The RNN select the word with higher
probability. (Answer: I love ✓)

# Vanishing Gradients Problem:-

→ Vanishing gradient Problem is hard to deal

⇒ Exploding gradients can fixed by clipping

## Vanishing gradient problem:-

Happens during training when gradients (the signal used to update weights) become smaller and smaller as the are propagated backward through many layers.

### Consequence
- The RNN remembers short term patterns
- It forgets long-term context because the weights for earlier time steps hardly get updated.
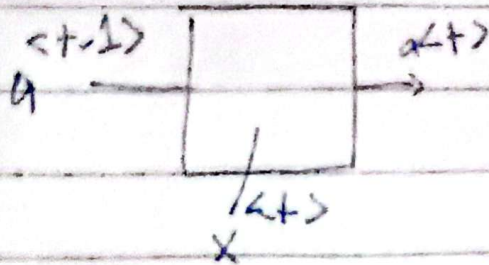
### fixing:-
★ GRU
★ Relu instead of tanh/sigmoid
★ Residual connection in deep RNNs.

## Gated Recurrent unit (GRU) :

Capturing long cell and treating vanishing
gradient

memory

$$a^{<t-1>} \longrightarrow \boxed{\phantom{xxxx}} \longrightarrow a^{<t>}$$

$$x^{<t>}$$

$C =$ memory cell , $\quad c^{<t>} = a^{<t>}$

$\tilde{c}^{<t>} = \tanh\left(W_c\left[c^{<t-1>}, x^{<t>}\right] + b_c\right.$

gamma $\Gamma_u =$ update gate $= (0, 1,)$

$$\Gamma_u = \sigma\left(W_u\left[c^{<t-1>}, x^{<t>}\right] * b_u\right)$$

$\Gamma_u$ used to find singular or plural.

$$\boxed{\text{softmax}} \longrightarrow \hat{y}^{<t>}$$

$c^{<t-1>} \longrightarrow \qquad \longrightarrow c^{<t>}$

$= a^{<t-1>}$

$\boxed{\text{tanh}} \quad \boxed{\sigma}$

$x^{<t>}$

## Equation

$$LSTM_5 = c^{<t>} = \tanh(W_c [\Gamma_r * c^{<t-1>}, u^{<t>}] + bc)$$

$$\Gamma_u = \sigma(W_u [c^{<t-1>}, u^{<t>}] + bu)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 + \Gamma_u) * c^{<t-1>}$$

## Detailed

GRUs: simplify LSTMs by merging some gates, making them faster and lighter while keeping long-term memory.

GRU Gates:-

Update gate $\Rightarrow$ decide how info of past to keep.
Reset gate $\Rightarrow$ decide how much of info to forget when computing new info.

Step 1 Update Gate: $z_t = \sigma(W_z u_t + U_z h_{t-1} + b_z)$

$z \in (0, 1)$

if close to 1 $\Rightarrow$ keep more old memory
if close to 0 $\Rightarrow$ update more with new info

Step 2 : Reset Gate $\Rightarrow r_t = \sigma(w_r n_t + V_r h_{t-1} + b_r)$

if $0 \Rightarrow$ forget old memory

if $1 \Rightarrow$ keep most of old memory

Step 3 : activation (new memory)

$$\underline{c}^{<t>} = \tan h(\ )$$

In

GRU and long-short-term-Memory (LSTM

Correction of equation

$$a^{<t>} = \Gamma_o * \tan h(c^{<t>})$$

# GRU

$$\tilde{c}^{<t>} = \tanh\left(W_c\left[\Gamma_c * c^{<t-1>}, x^{<t>}\right] + b_c\right)$$

Update gate

$$\Gamma_u = \sigma\left(W_u\left[c^{<t-1>}, x^{<t>}\right] + b_u\right)$$

$$\Gamma_r = \sigma\left(W_r\left[c^{<t-1>}, x^{<t>}\right] + b_r\right)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

# LSTM

$$\tilde{c}^{<t>} = \tanh\left(W_c\left[a^{<t-1>}, x^{<t>}\right] + b_c\right)$$

Update gate

$$\Gamma_u = \sigma\left(W_u\left[a^{<t-1>}, x^{<t>}\right] + b_u\right)$$

$$\Gamma_f = \sigma\left(W_f\left[a^{<t-1>}, x^{<t>}\right] + b_f\right)$$

$$\Gamma_o = \sigma\left(W_o\left[a^{<t-1>}, x^{<t>}\right] + b_o\right)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$

## Bidirectional Infor

Take info from next and previous or present layer :-



We can take info from future and past neron cell.

☆ By taking info once in forward time, and then in backward, this way it take info from future and past.
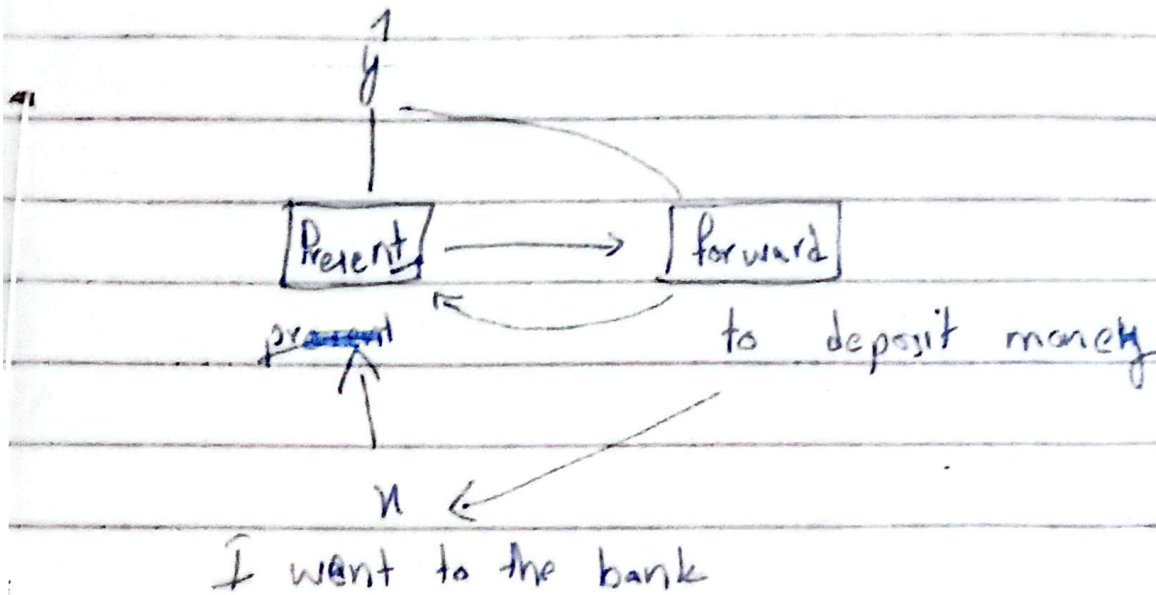
## Working

Trained Model

We want to check meaning

→ I went to the bank (financial or river bank)

step 2: To check this we will go forward and backward in time.

Forward RNN: I went to the bank ( itt a noun)

Back ward RNN
see "to deposit money" ( confirm it's a
financial bank not a river bank).



$$\hat{y} \uparrow 1$$

Present $\longrightarrow$ forward

present

to deposit money

$n$

I went to the bank

# Deep RNNs
Stacking multiple RNN layers on top of
each other like a sandwiched to learn complex
thing ( NLP's, GAN)