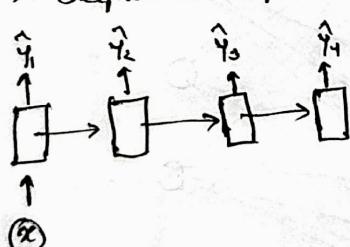


Recurrent Neural Network (RNN)

RNN

- It is a class of Neural network designed for sequential data like time series, text, audio etc.
- Unlike feedforward network, RNNs have memory - they use previous outputs as input for current processing.
- RNNs maintain a hidden state that carries information from previous time steps to influence the current prediction.

Types of RNN Architectures

- ① One to one → Regular neural network (e.g., image classification)
- ② One to many → Sequence output (e.g. image captioning), single input
- ③ Many to one → Single output (e.g. sentiment analysis), Sequence input
- ④ Many to Many → Sequence output (e.g. translation, speech recognition), Sequence input

(RNN) Recurrent Neural Network

RNN architecture

Review

	Sentiment
movie was good	1
movie was bad	0
movie was not good	0

Word to vector

unique!

$$\text{movie} = [10000]$$

$$\text{was} = [01000]$$

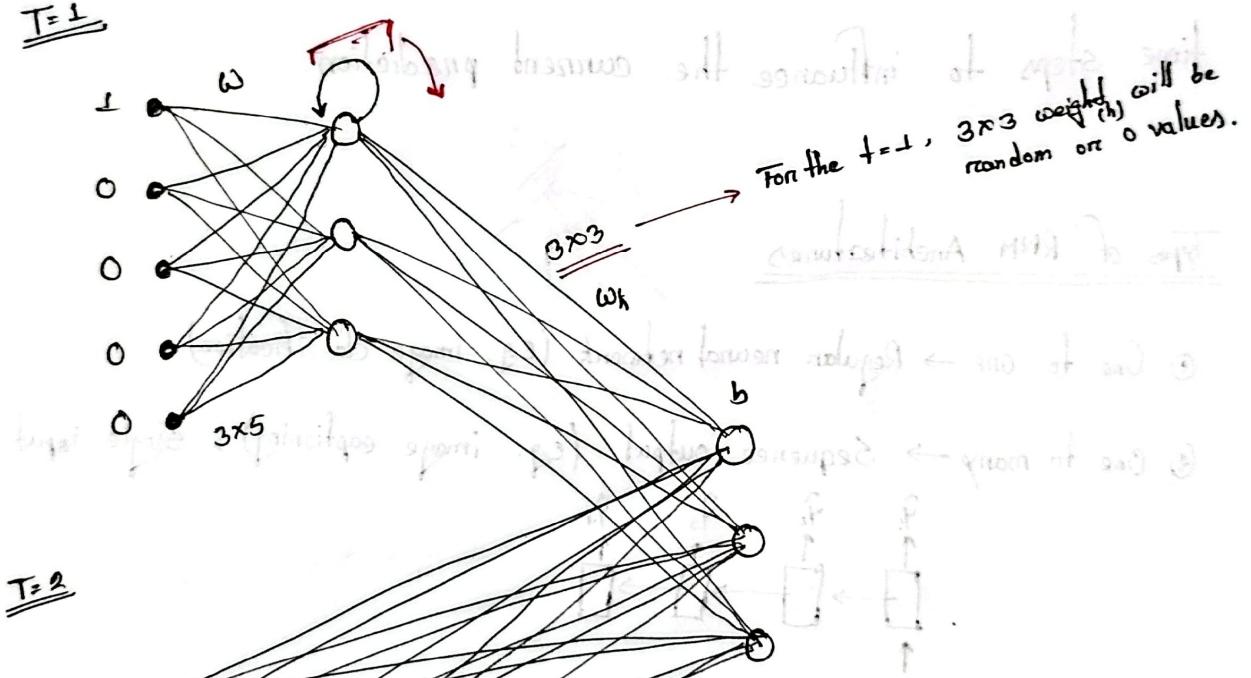
$$\text{good} = [00100]$$

$$\text{bad} = [00010]$$

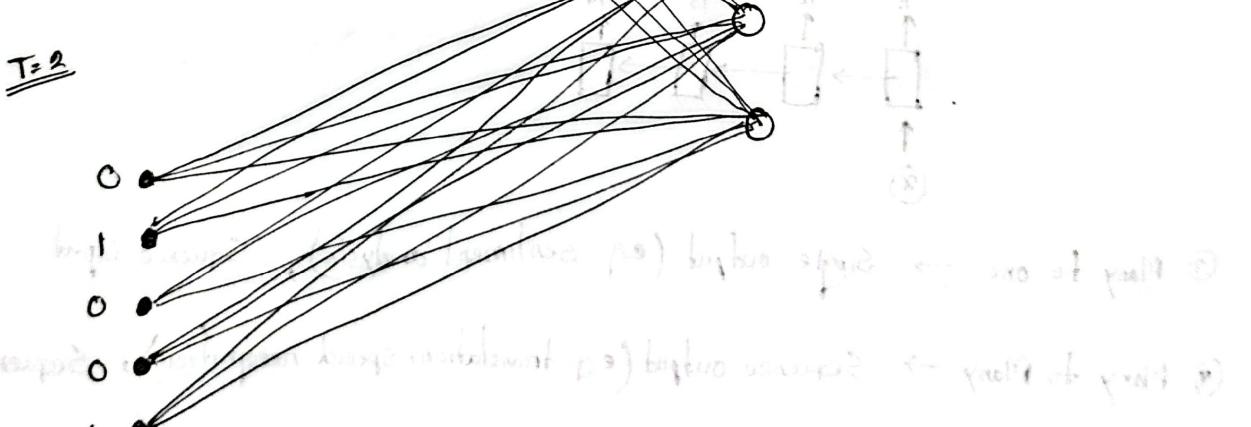
$$\text{not} = [00001]$$

answering next question comes after state update or hidden state.

$T=1$



$T=2$



Forward propagation in RNN

review

Sentiment

$$x_{11} \quad x_{12} \quad x_{13}$$

$$x_{21} \quad x_{22} \quad x_{23}$$

$$x_{31} \quad x_{32} \quad x_{33}$$

initial input state

loop [0] of first

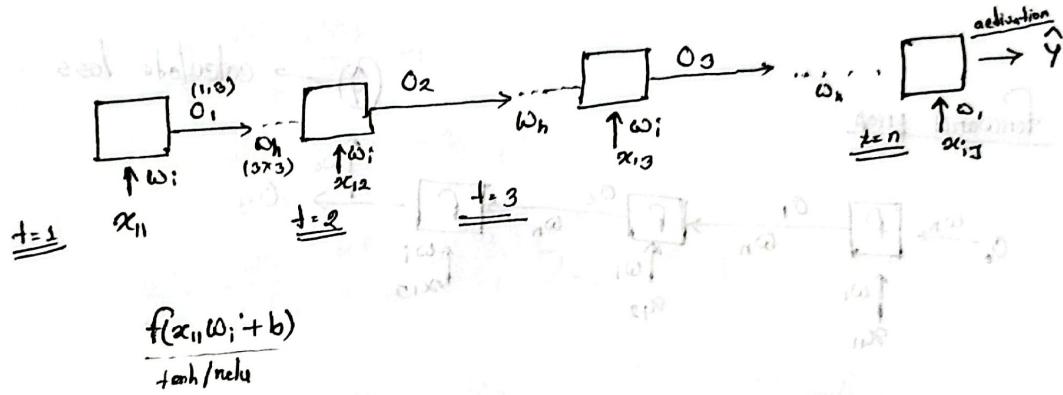
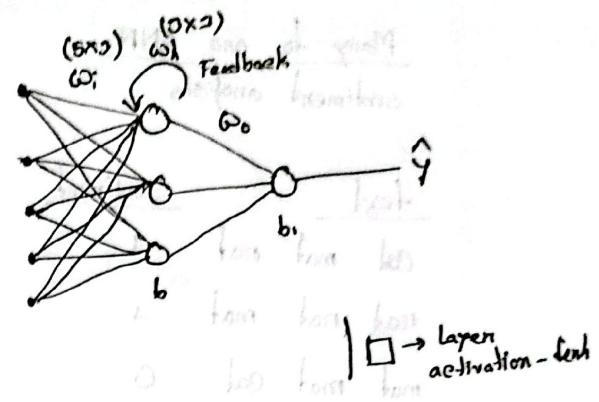
loop [0]

loop [0] of first

loop [0] of first

$$f(x_{12}w_i + \Delta w_h + b_i)$$

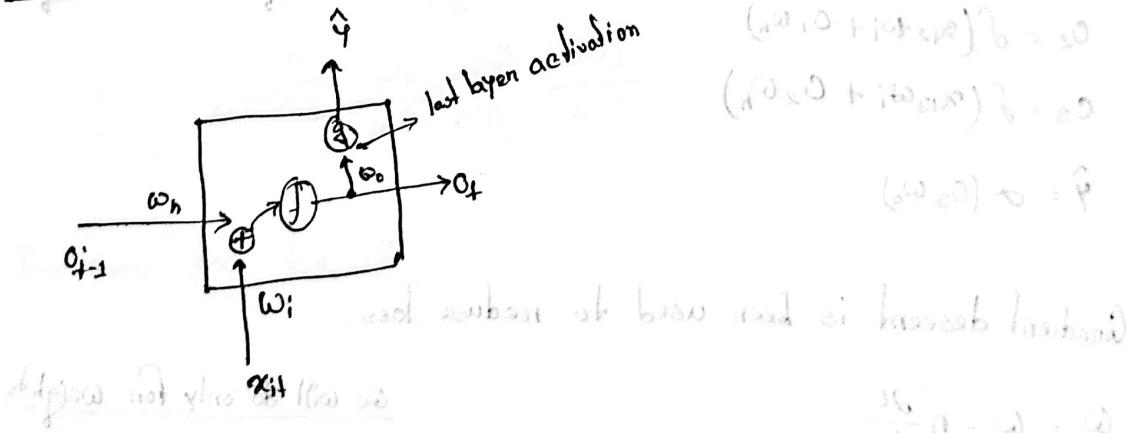
Forward \rightarrow unfolding through time



$$f(x_{11}w_i + b)$$

tanh / relu

Simplified representation



$$(P \times 3 + P \times 3) \times 3 = 3$$

$$(3 \times 3 + 3 \times 3) \times 3 = 3$$

$$(3 \times 3 + 3 \times 3) \times 3 = 3$$

$$(3 \times 3) \times 3 = 3$$

$$\frac{d}{dt} \left(\frac{\partial L}{\partial h_t} \right) = \frac{\partial L}{\partial h_t} \cdot \frac{d}{dt} h_t$$

$$\text{backpropagation} \quad \frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \cdot \frac{d}{dt} h_{t+1}$$

$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \cdot \frac{d}{dt} h_{t+1}$$

Backpropagation in RNN

Many to one RNN

sentiment analysis

text	sentiment
cat mat rat	1
rat rat mat	1
mat mat cat	0

vector representation

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

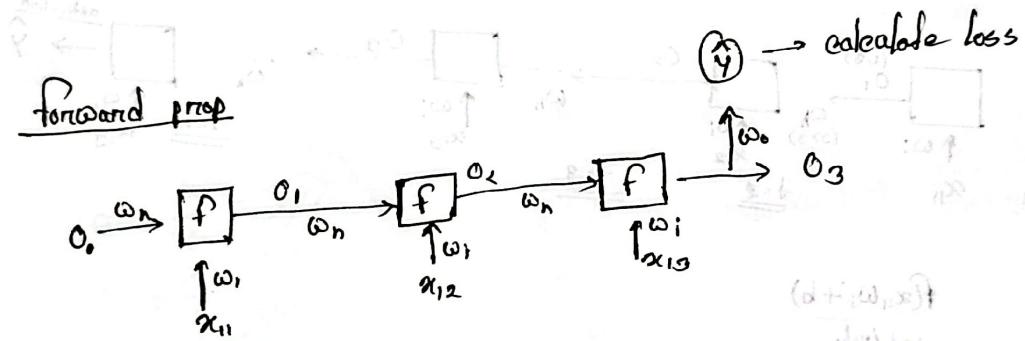
$$[0 \ 0 \ 1] [0 \ 0 \ 1] [0 \ 1 \ 0]$$

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

y

1

0



Here $\xrightarrow{\text{activation}}$

$$o_1 = f(x_{11}w_1 + o_0w_h)$$

$$o_2 = f(x_{12}w_1 + o_1w_h)$$

$$o_3 = f(x_{13}w_1 + o_2w_h)$$

$$\hat{y} = \sigma(o_3w_o)$$

Let, $\underline{\text{loss}} = -\varphi_i \log \hat{y}_i - (1-\hat{y}_i) \log (1-\hat{y}_i)$



Gradient descent is been used to reduce loss.

We will do only for weights

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

We have to find

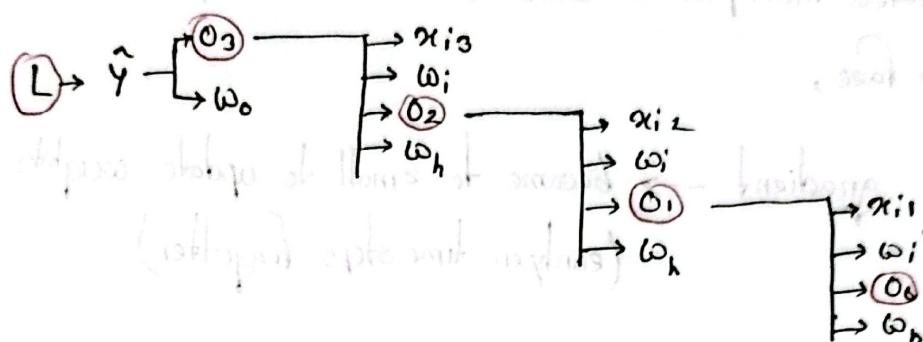
$$w_h = w_h - \eta \frac{\partial L}{\partial w_h}$$

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

If we unfold for 3 time,

the flow will look like,

[all dependency of Loss]



$$\text{For } i, \frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial w_i} +$$

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

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

$$\therefore \frac{\partial L}{\partial w_i} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_j} \frac{\partial O_j}{\partial w_i}$$

$$\text{For } h, \frac{\partial L}{\partial w_h} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_j} \frac{\partial O_j}{\partial w_h}$$

The same goes for $\frac{\partial L}{\partial w_0}$.

Major Problems with RNN

Due to repeated multiplication while calculating the partial differentiation, we can face,

① Vanishing gradient → became too small to update weights
(earlier time steps forgotten)

② Exploding Gradient → became too large and destabilize training

Because of such problem, RNNs are been used very less.

For better performance we use,

⇒ LSTM

⇒ GRU

Long Short-Term Memory (LSTM)

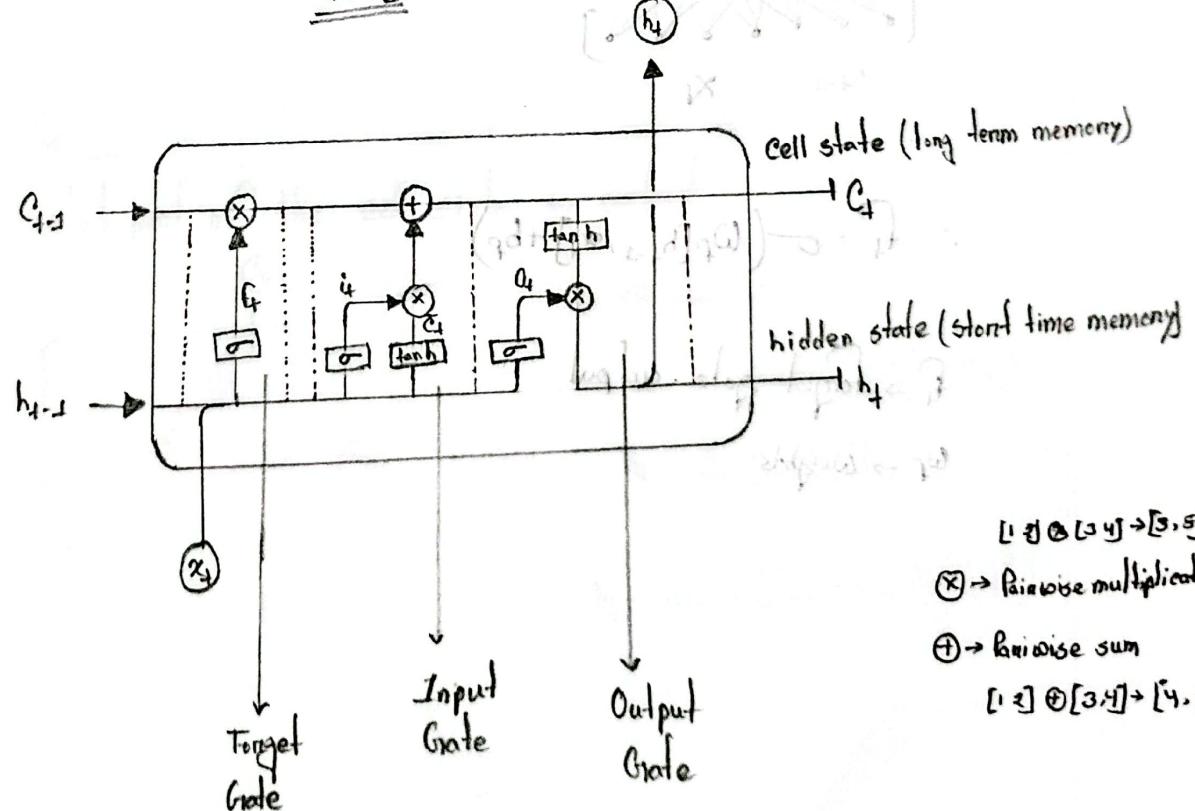
LSTM

It's a special kind of RNN capable of learning long term dependencies.
It is mainly designed to solve the vanishing Gradient problem of RNN.

What was the key idea?

- It uses gates to control how much past information should be remembered or updated.
- LSTM introduce memory cells, and gating mechanisms to regulate information flow.

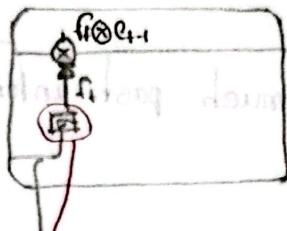
LSTM Architecture



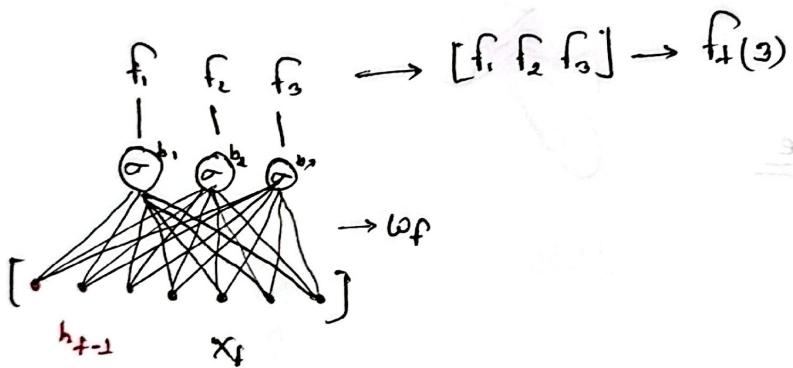
(HIDDEN) space || first - last part

Components of an LSTM Cell

① Forget gate



This gate decides what information to remove from the cell state.

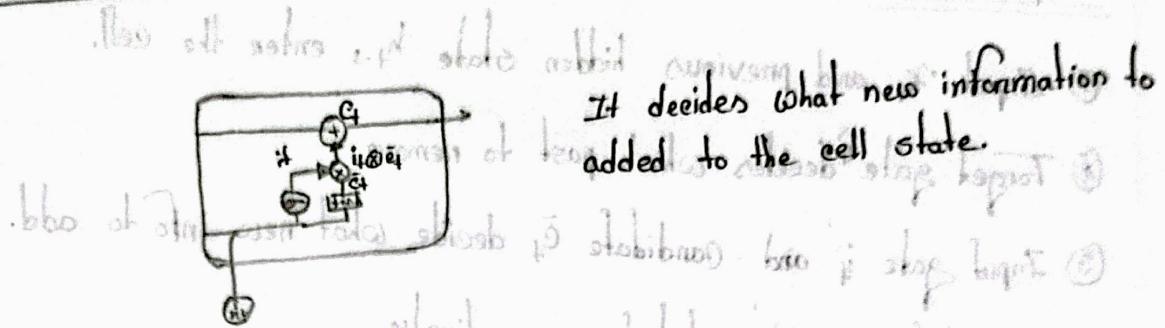


$$\therefore f_t = \sigma(w_f[h_{t-1}, x_t] + b_f)$$

f_t → forget gate output

w_f → weights

② Input Gate



$$c_t = f_t \otimes c_{t-1} + i_t \otimes \bar{c}_t$$

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

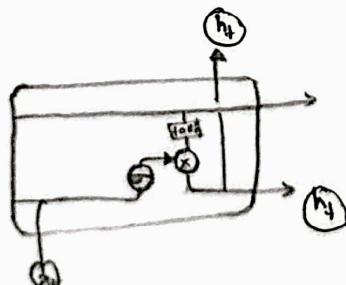
$$\bar{c}_t = \tanh(\omega_c [h_{t-1}, x_t] + b_c)$$

i_t = input gate activation (Controls how much to update) (filter)

\bar{c}_t = Candidate cell state (new potential memory)

③ Output Gate

Decides what part of the cell state to output.



$$\rightarrow \tanh(-1, 1)$$

$$o_t = \sigma(\omega_o [h_{t-1}, x_t] + b_o)$$

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

h_t = new hidden state (output)

Complete Flow

short logic ①

- ① Input x_t and previous hidden state h_{t-1} enter the cell.
- ② Forget gate f_t decides which past to remove.
- ③ Input gate i_t and candidate \tilde{c}_t decide what new info to add.
- ④ Cell state c_t is updated accordingly.
- ⑤ Output gate o_t decides what to output as h_t .

$$(id + [W_i, b_i] \cdot v_t) \odot = p_i$$

$$(sd + [W_f, b_f] \cdot v_t) \text{danh} = \tilde{p}$$

(forget) (forget of don't want old) no forget short logic - p
 (input) (input of want new) states like old - p

short logic ②

but two of state has set to bring back original

$$(i, \tilde{c}_t) \text{danh} = p$$

$$(sd + [W_i, b_i] \cdot v_t) \odot = p$$

$$(f) \text{danh} \otimes p = d$$

(forget) state didn't care - p

Gated Recurrent Unit (GRU)

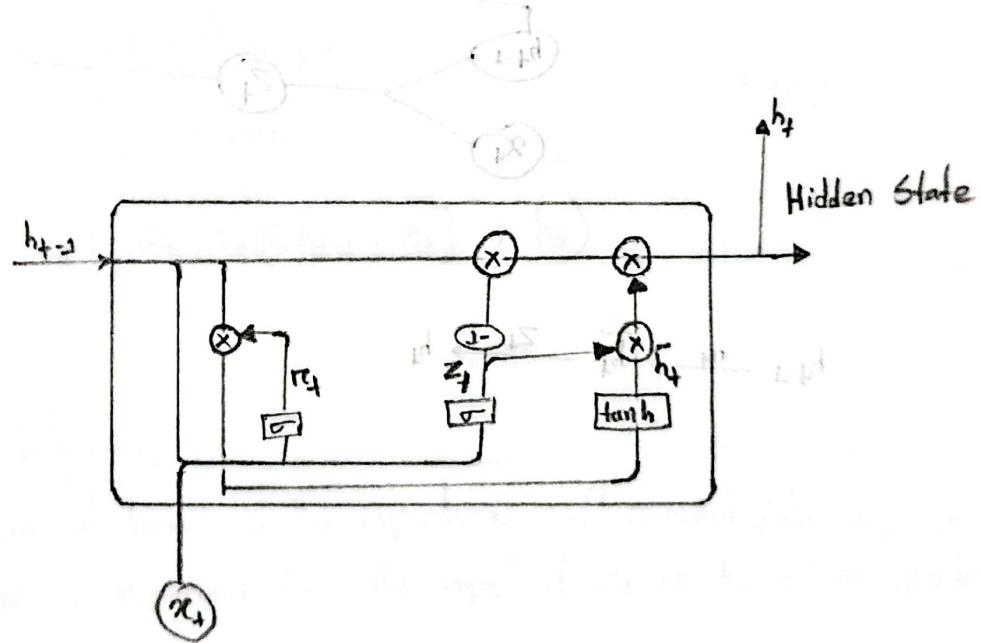
GRU

It's a variant of LSTM, designed to make training simpler and faster, while still solving vanishing gradient problem.

Key idea

GRU combines the forget gate and input gate of LSTM into a single update gate, reducing the number of parameters.

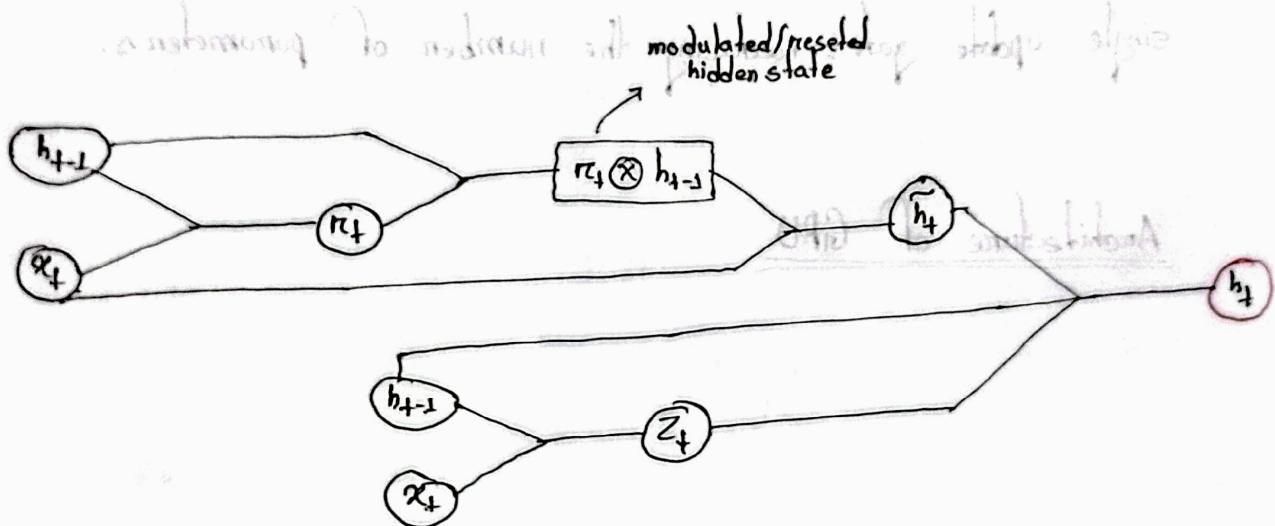
Architecture of GRU



(iii) Forward pass

Flow of calculating h_t

- ① Calculate reset gate (r_t)
- ② Calculate candidate hidden state (\tilde{h}_t)
- ③ Calculate Update gate (z_t)
- ④ Calculate current hidden state (h_t)



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

① Update gate (z_t)

Controls how much past information to carry forward.

$$z_t = \sigma(\omega_z [h_{t-1}, x_t] + b_z)$$

② Reset gate (r_t)

Controls how much past information to forget or re-set.

$$r_t = \sigma(\omega_r [h_{t-1}, x_t] + b_r)$$

③ Candidate hidden state (\tilde{h}_t)

Creates new memory content.

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

④ Final hidden state (h_t)

Here it balances out based on the importance of current data (z_t). If current information is important then the impact of z_t on h_t will be greater or vice versa.

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

LSTM vs GRU

	<u>LSTM</u>	<u>GRU</u>
<u>Feature</u>	Has 3 gates at each step of sequence (Input, forget, output)	Has 2nd hidden state
<u>Gate</u>	Yes (Ex) [add]	No (only hidden state)
<u>Memory cell</u>	None	Fewer
<u>Parameter</u>	More	Faster
<u>Training Speed</u>	Slower	Nearly same
<u>Performance</u>	Slightly better for long dependency	(add + [add ⊕ add ⊕ add]) / 3 = add
<u>Complexity</u>	High	Simpler
<u>Suitable for</u>	Complex, long term sequence, big data	Simpler, small data

W.L.G. lets focus on unidirectional LSTM because it's simple
we can just add as many cells as we want with backpropagation or without backpropagation

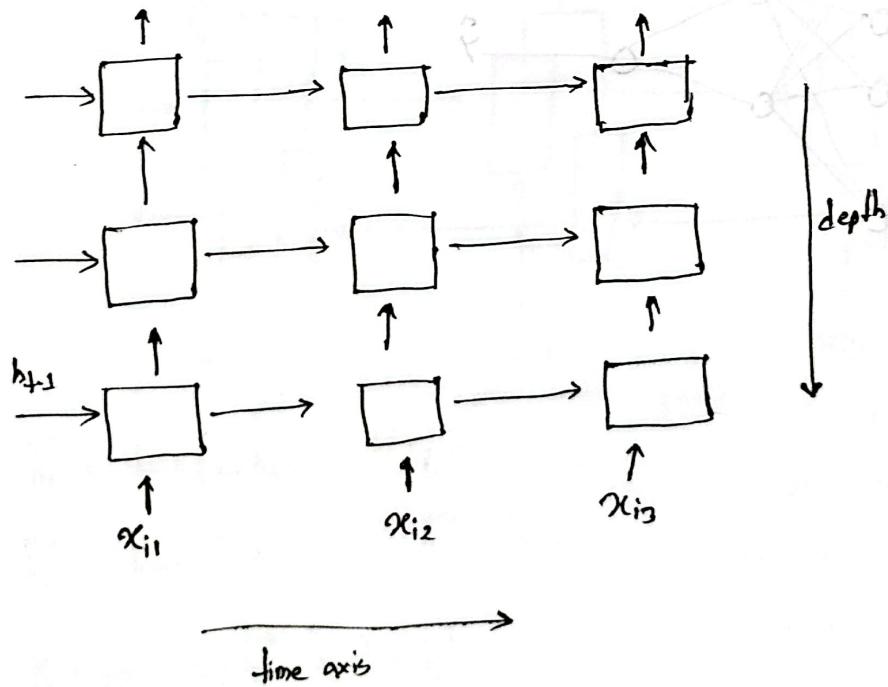
$$h_t = \text{tanh}(Wx_t + b + h_{t-1})$$

Deep RNN

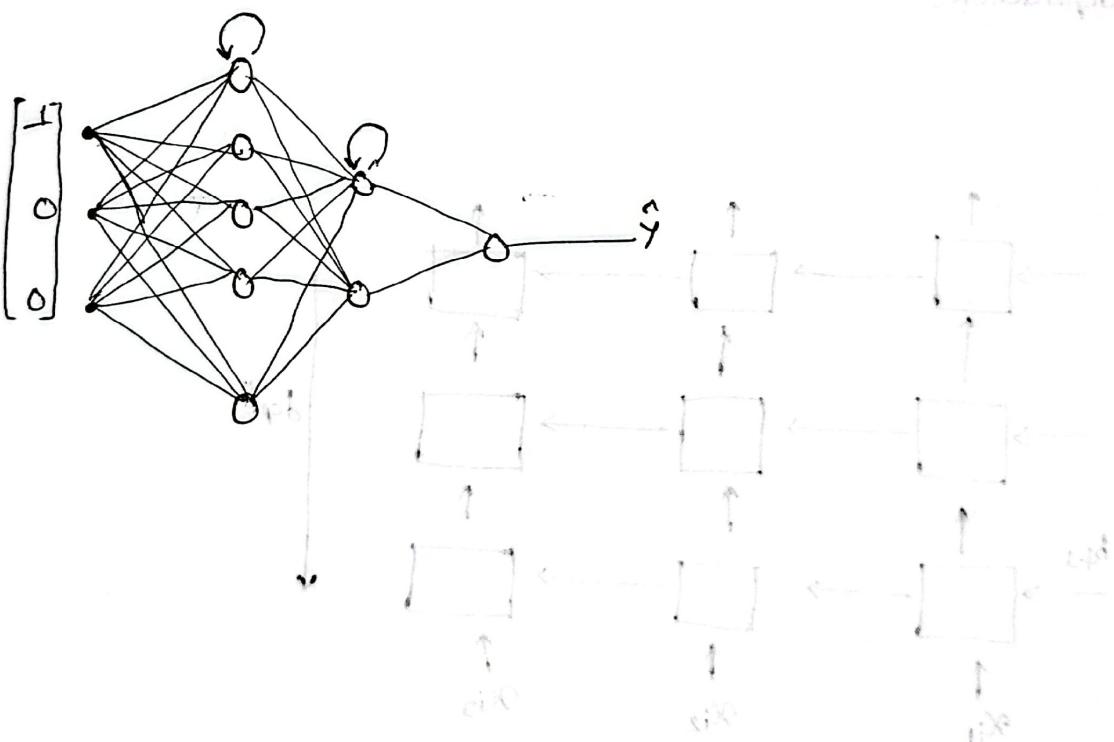
- ⇒ An extension of a standard RNN where multiple recurrent layers are stacked on top of each other.
- ⇒ Each layer processes the sequential data and passes its hidden states to the next layer.

Key idea

Deep RNNs learn hierarchical representational temporal representations, where lower layers capture short term patterns and higher layers capture long term dependencies.



model = Sequential([
 Embedding(10000, 32, input_length=100),
 SimpleRNN(5, return_sequences=True), // This will be true except the last layer
 SimpleRNN(2),
 Dense(1, activation='sigmoid')
])



Bidirectional RNN (BiRNN)

This is an extension of RNN that can process data in both forward and backward direction.

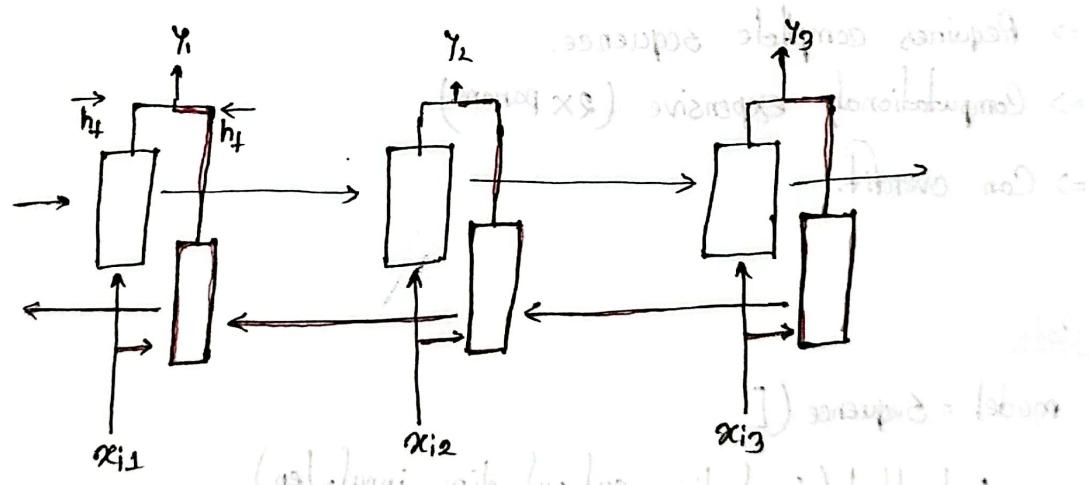
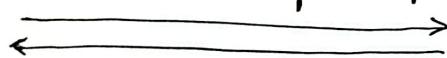
⇒ Useful for the tasks where the entire input sequence is available before prediction.

⇒ When the current prediction can depend on the future. like,

Amazon, a great website.

Amazon, a big river.

Amazon, a huge jungle.



$$\vec{h}_t = \tanh(\vec{W}\vec{h}_{t-1} + \vec{U}\vec{x}_t + \vec{b})$$

$$\bar{h}_t = \tanh(\bar{W}\bar{h}_{t+1} + \bar{U}\bar{x}_t + \bar{b})$$

$$y_t = \sigma(v[\vec{h}_t, \bar{h}_t] + b)$$

Mostly used with LSTM (BiLSTM) and (BiGRU)

⇒ Used LSTM and GRU for both direction and to overcome the issue

⇒ In BiLSTM, handles long term dependency better

⇒ In BiGRU, faster and computationally lighter

Adv

⇒ Use both past and future context.

⇒ Improve accuracy in sequence pred. tasks.

⇒ Works well in NLP and speech recognition.

Disadv

⇒ Requires complete sequence.

⇒ Computationally expensive ($2 \times$ param.)

⇒ Can overfit.



Code

```
model = Sequence([
```

```
    Embedding(input_dim, output_dim, input_len),
```

```
    Bidirectional(SimpleRNN(5, return_sequences=True)),
```

```
    Bidirectional(SimpleRNN(2)),
```

```
    Dense(1, activation='sigmoid')
```

])