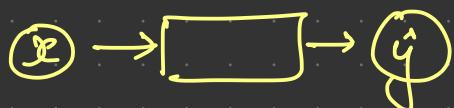


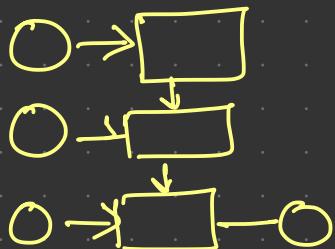
Recurrent Neural Networks

- Sequential modelling and data
- Sequencing model applications

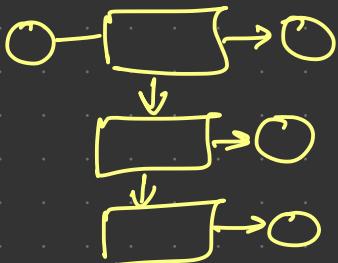
1) Binary Classification: (One to one)



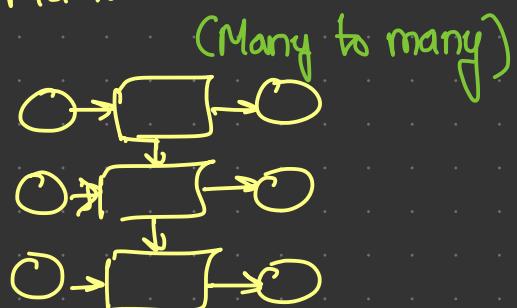
2) Sentiment classification (Many to one)



3) Image captioning



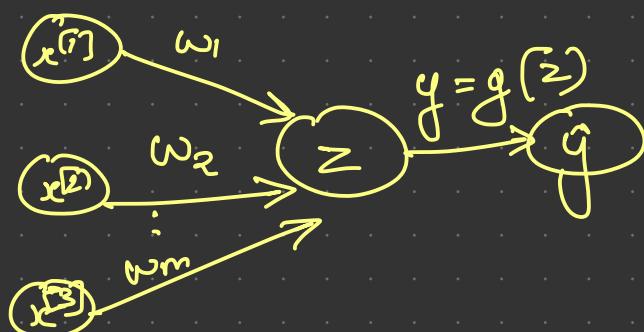
4) Machine translation



Sequences to another sequence

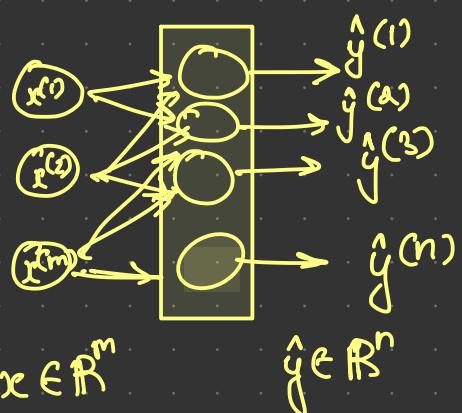
Neurons with recurrence:

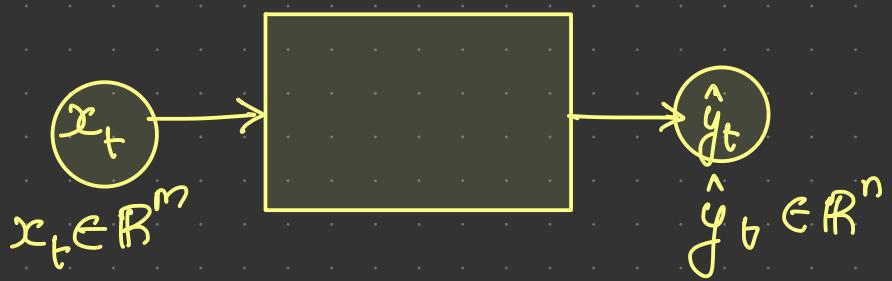
In perception,



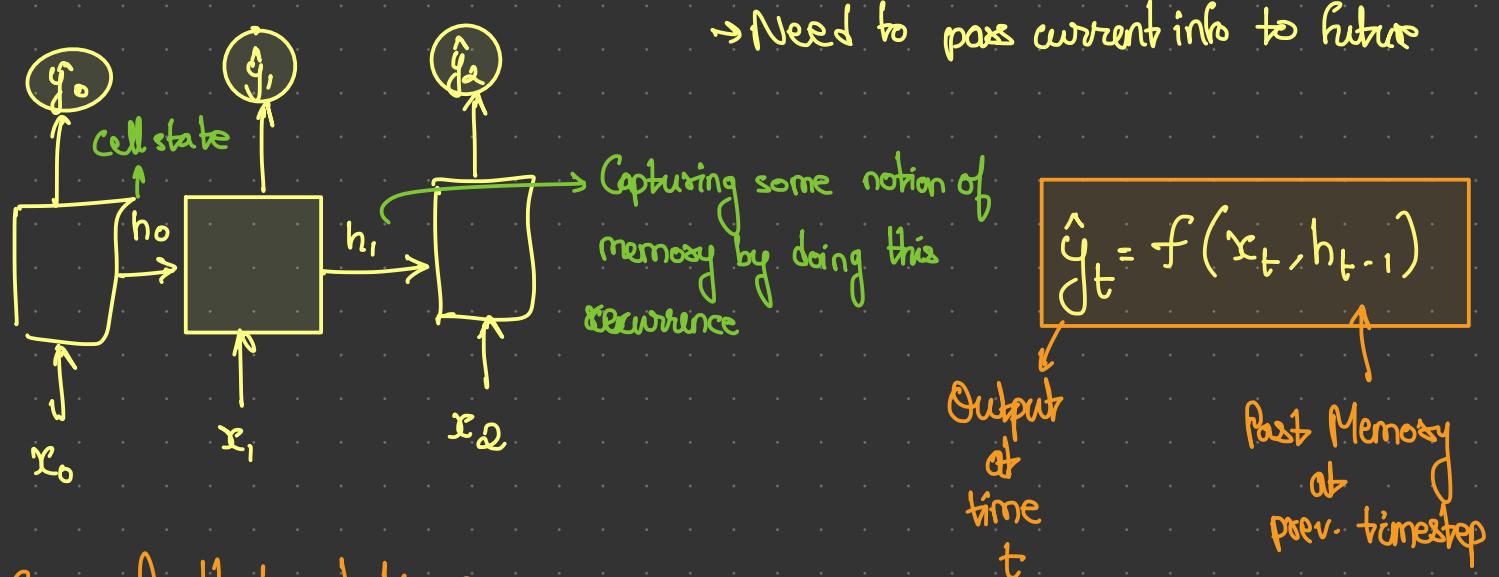
No notion of time sequences
in above

Forward propagation

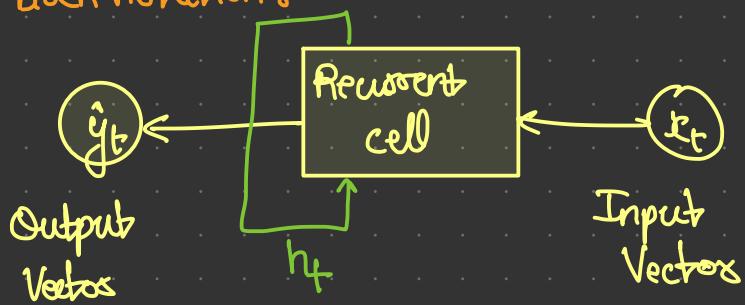




If we have multiple timesteps:



General block notation:



RNNs have internal state
 h_t , is updated at each time step

Recurrence relations

$$h_t = f_w(x_t, h_{t-1})$$

Cell state \downarrow
function with weights w

old state \downarrow
input \downarrow

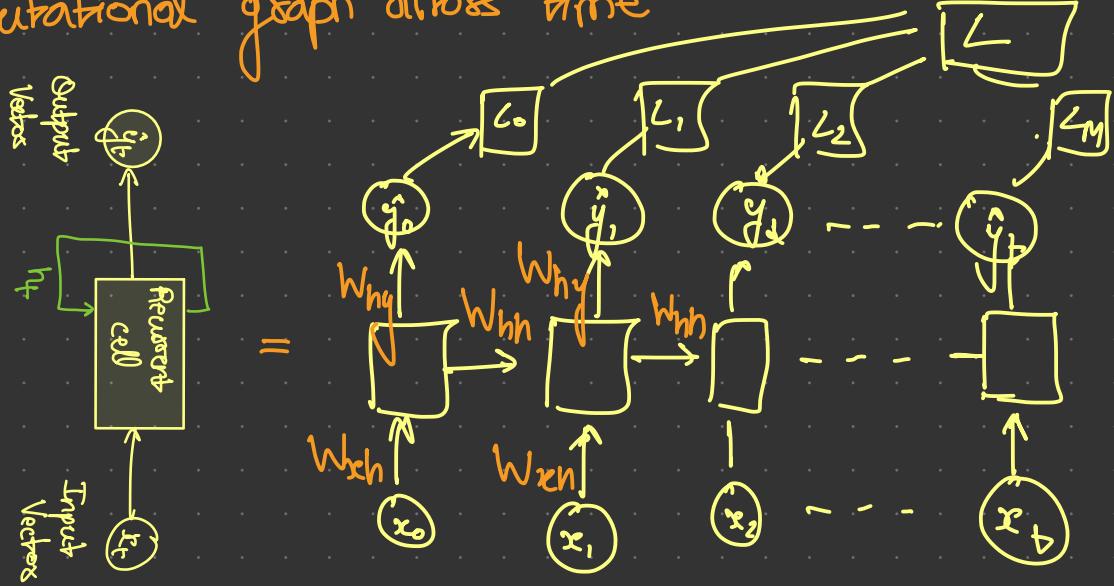
$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$

x_t \nwarrow
respective weight

$$\hat{y}_t = \hat{W}_{hy} h_t$$

NOTE: Same params,
functions are used every time step

Computational graph across time



Design Criteria for sequence modelling in general

- 1) Handle variable-length sequences
→ Ability to handle variable length of sentences etc
- 2) Track long-term dependencies
To have a notion of memory
- 3) Maintain information about order
How prev. affects next
- 4) Share parameters across sequences
RNN's have all these 4

Solved examples

Given Predict

Represent language as NN:



- Issues
- 1) Usually work with mathematical operations i.e. needed require it

Instead use embeddings

↳ Transform indexes into vectors of fixed size.

I. Vocabulary [Generate all possible words]

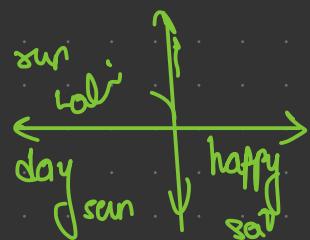
II. Indexing [$\begin{matrix} 1 \rightarrow 1 \\ 2 \rightarrow 2 \\ \vdots \\ n \rightarrow N \end{matrix}$]

III Embedding:

1) One shot embedding

i^{th} index is mapped
correspondingly.

2) Learned embedding



NOTE: feed forward loops can't handle variable lengths

1) RNNs have variable length

2), 3)

cell state dependent on prev.

BACK PROPAGATION THROUGH TIME [BPTT]

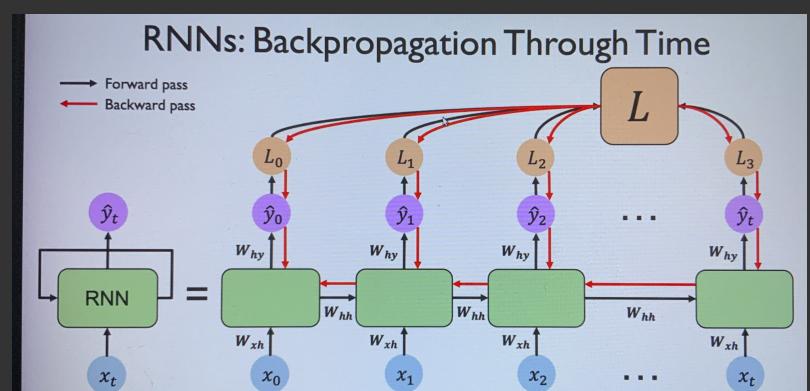
In feed forward algorithms

1) Take gradient of loss w.r.t. each parameter

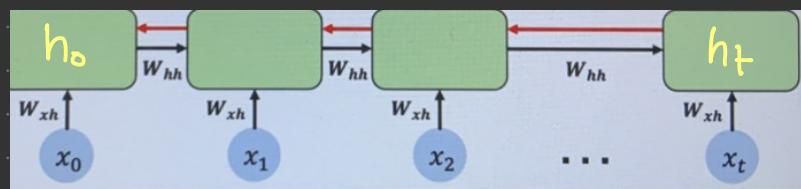
2) Shift parameters to minimise loss

In RNNs,

1) Forward pass through time



Standard RNN Gradient flows



Involves many factors to get w.r.t. to h_0

If values are too large
[exploding gradients]

Gradient clipped to scale
big gradients

If values are too small
[vanishing gradients]

- ① Activation fn
- ② Weights init
- ③ Network Architecture

Why is vanishing a problem?

→ No learning [can't propagate error]

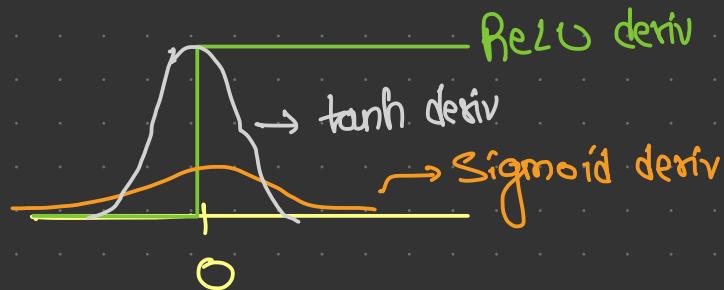
→ End up biasing in short term, won't work in long term

→ Won't learn well when discussing long sentences etc

↳ long term dependence

① Activation Functions:

ReLU prevents f' from shrinking the gradients when $x > 0$



② Parameters init

$I_n = \begin{bmatrix} 1 & 0 \\ 0 & \ddots & -1 \end{bmatrix}$

Initialise biases to 0 } Helps prevent weights from
Initialise weights to I } shrinking to zero

③ Gated Cells

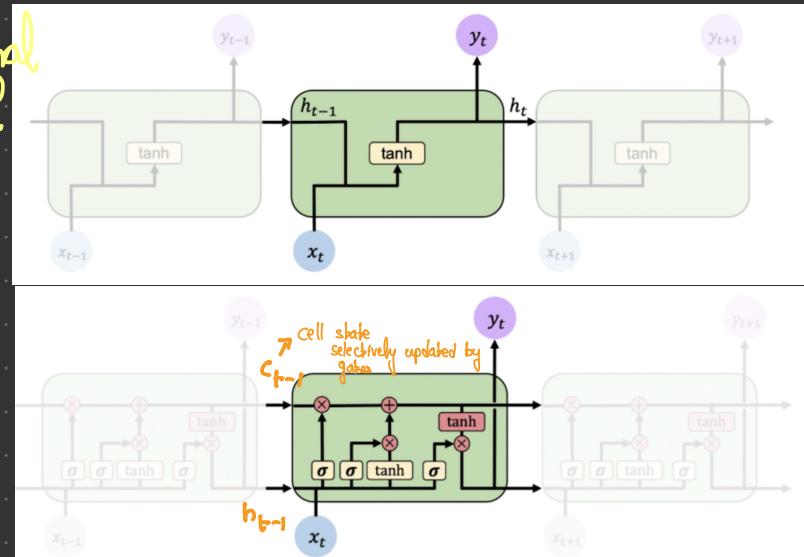
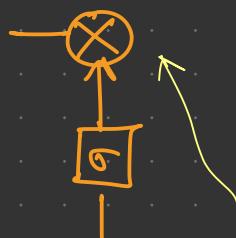
Using more complex recurrent unit with gates to control what information is passed through



→ LSTM: Long short term memory networks rely on a gate cell to track information through many time steps

LSTM Networks

- Contains computational blocks that control info flow
- Able to track info throughout time step

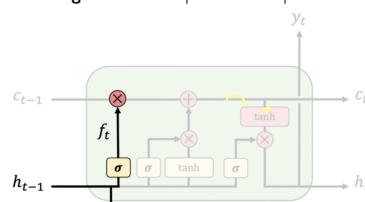


Info is added (or) subtracted through gates
Optionally let info through

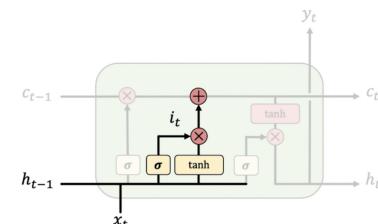
- Steps: 1) forget
2) store
3) update
4) output

→ Regulates info slow

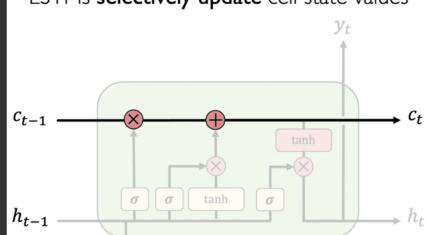
1) Forget 2) Store 3) Update 4) Output
LSTMs forget irrelevant parts of the previous state



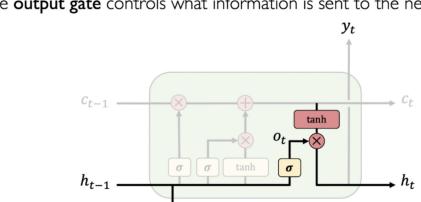
1) Forget 2) Store 3) Update 4) Output
LSTMs store relevant new information into the cell state



1) Forget 2) Store 3) Update 4) Output
LSTMs selectively update cell state values

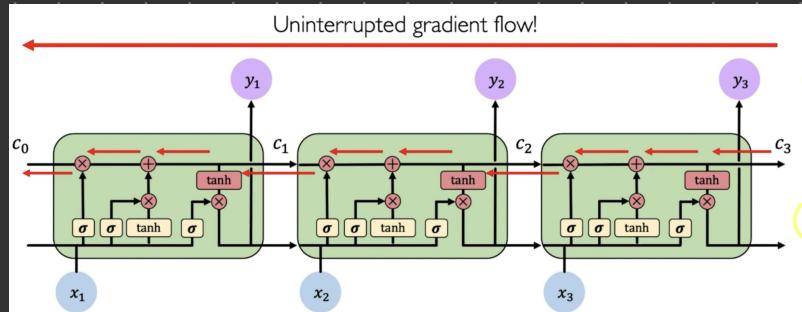


1) Forget 2) Store 3) Update 4) Output
The output gate controls what information is sent to the next time step



Allows uninterrupted gradient flow:

key ideas of LSTM:



- ① Maintain separate cell state from outputted
- ② Use gates to control flows

→ forget: Gets rid of irrelevant info
→ Stores relevant info from current input
→ Selectively UPDATE cell state
→ OUTPUT gate outputs filtered version of cell state

- ③ Backprop with uninterrupted gradient flow

Practical Applications:

- 1) Music generation
- 2) Sentiment classification
- 3) Machine Translation
- 4) Weather model fitting

Issues of RNNs:

- 1) Encoding bottleneck
- 2) Not efficient as its sequential [no parallelisation]
- 3) Not long memory [limited capacity]

Solutions

ATTENTION

block is to learn from input which points and states to attend to [very efficient]
→ provides reasonable memory access
→ Only one pass needed (no back prop)