

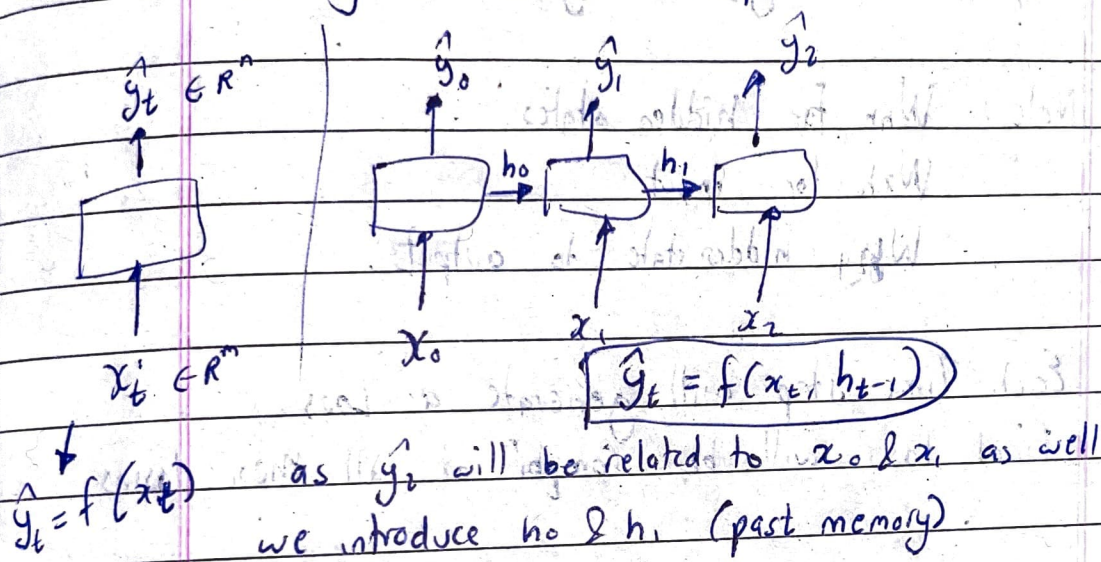
Sequential Data

Sequences

- Sentiment Classification → many to one
- Image Captioning → one to many
- Machine Translation → many to many

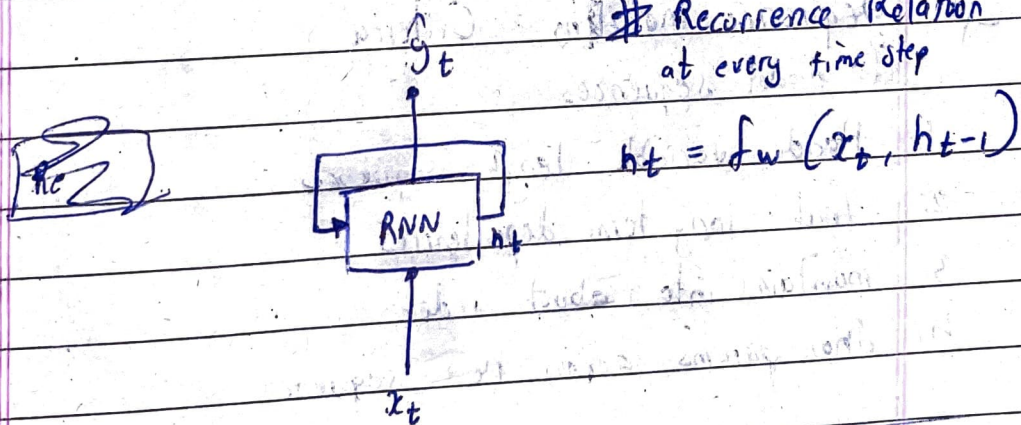
Neurons with Recurrence

Handling Individual Time Steps



Appls

Recurrence Relation at every time step



RNN State Output & Update

Input Vector x_t

Update hidden state $h_t = \tanh(W_{hh}^T h_{t-1} + W_{hx}^T x_t)$

Output Vector $\hat{y}_t = W_{hy}^T h_t$

Note: W_{hh} for hidden states

W_{hx} for input

W_{hy} hidden state to output

→ Each Time step will generate a loss

Final loss will be combiⁿ of all these losses

Sequence modelling : Criteria

To model sequences

1. Handle variable length sequences
2. Track long term dependencies
3. maintain info about order
4. Share params across the sequence

Let Predict the next word

→ NN can't interpret words

∴ Embedding: transform indexes into a vector of fixed size

1. Vocab

corpus words

2. Indexing

word to index

3. Embeddings

i. one hot embedding

cat = [0, 1, 0, 0, ...]

↑
index

ii. learned embedding

run walk	dog cat
day sun	happy sad

similar words
close

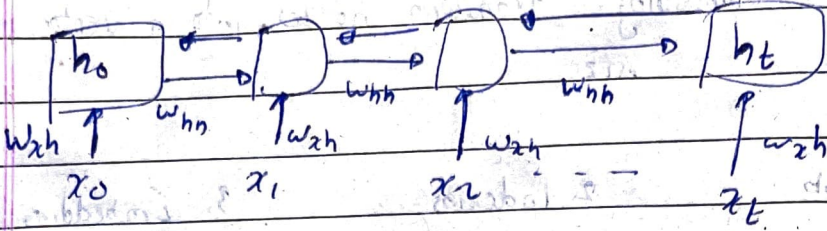
RNN can handle long variable sequence length
so different sequence length are just different
number of time steps

Long term dependency met by RNN by using past
hidden state

Capture differences in sequence order

Back Prop through time (BPTT)

as we go through horizon for grad which is time



gradient wrt h_0 involves many factors of w_{hh} & repeated gradient computation

Problems →

- ① many values > 1 :
exploding gradient
grads clipping to
scale big grads
- ② Many values < 1
vanishing grads
 1. Activation func
 2. wt initialization
 3. Network architecture

Errors due to further
back time have smaller &
smaller grads

→ Bias params to
capture short term
dependencies

Trick 1 ReLU prevents shrinking

Trick 2 Init cut to Identity matrix & biases to zero
helps preventing cuts to shrink to zero

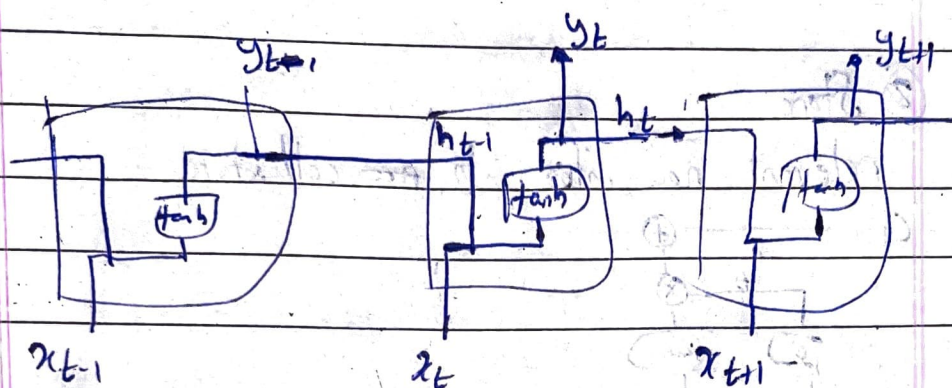
Trick 3 →

Trich 3 Gated Cells (LSTM, GRU)

use a more complex recurrent unit with gates control what info passes

Long short term memory (LSTM)

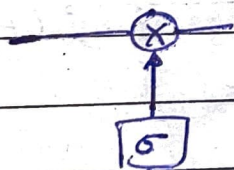
→ Standard RNN



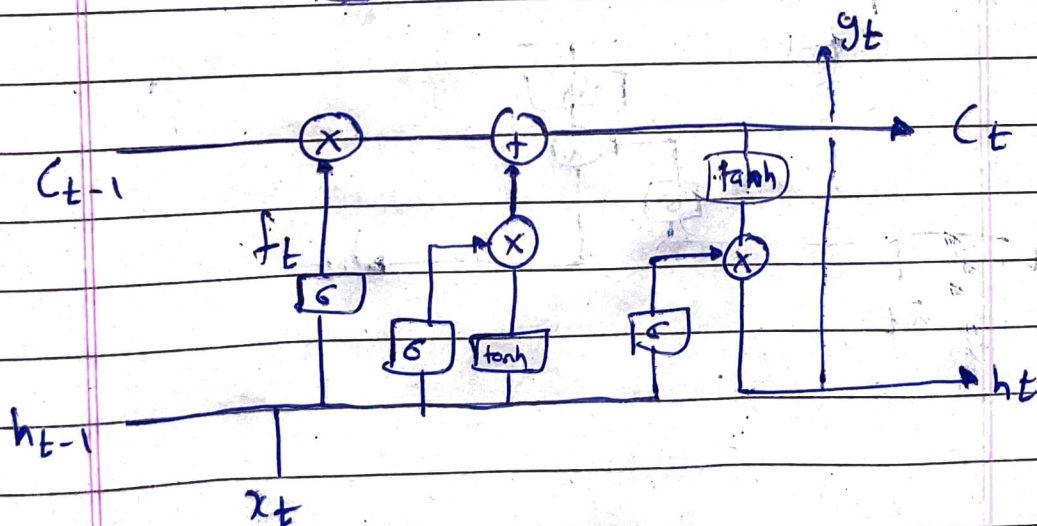
→ LSTM contain computational blocks that control information flow

alot more complex

→ Info added or removed using gates



eg. sigmoid layer & pointwise multi"

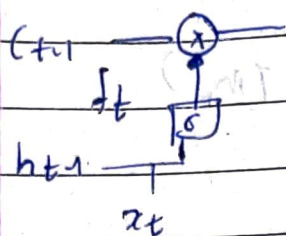


#

LSTM

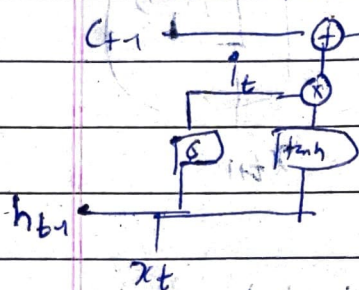
① Forget

irrelevant parts of previous state



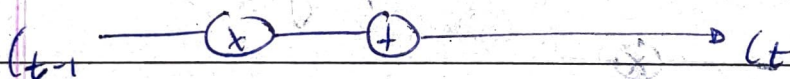
② Store

relevant new info into new cell state

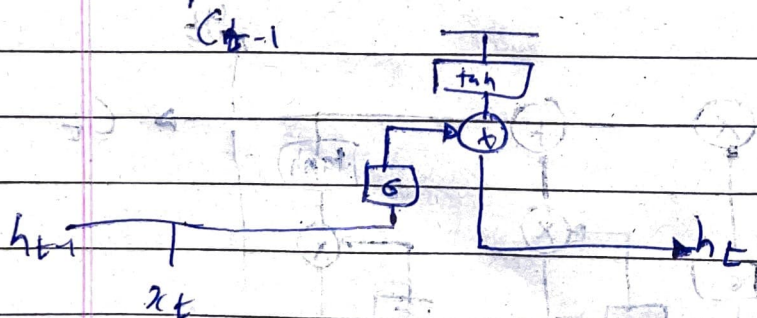


③ Update

selectively cell state



④ Output



→ maintaining a separate cell state allows uninterrupted gradient flow

Key concepts

- 1- maintain separate cell state from output
- 2- Use gates to control info flow
 - forget irrelevant
 - store relevant
 - selectively update cell state
 - output
- 3 BPTT with ~~unint~~ uninterrupted grad flow

Problems

1. Encoding bottleneck
2. Slow, no parallelisation
BPTT very expensive
3. Not long-memory

Attention model provide learnable memory access
back prop not needed

