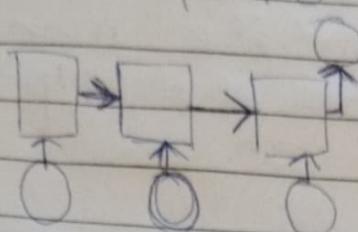
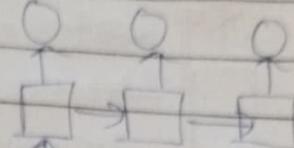
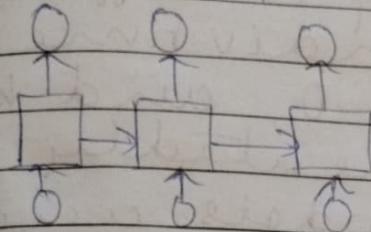


Recurrent NN, Transformers and Attention

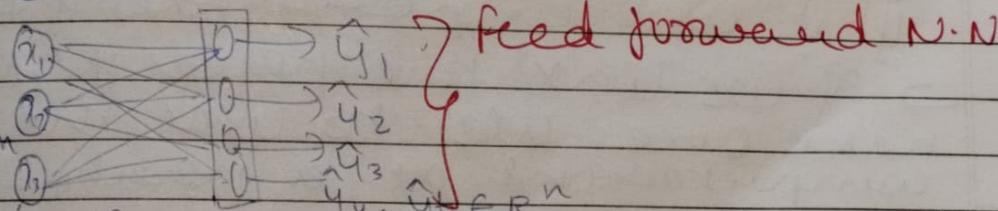
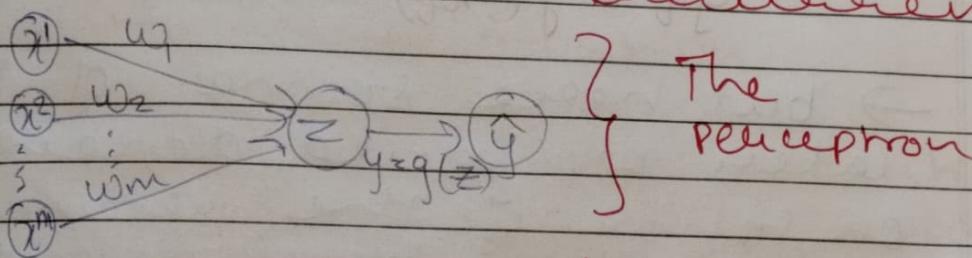
PAGE NO.
DATE

Sequence model applications

- 1) Many to one 
- 2) One to many 
- Many to one
eq - sentiment classification
- One to many
eq - captioning

- 3) Many to many 

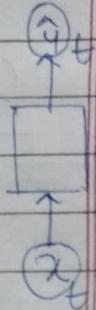
Neurons with recurrence.



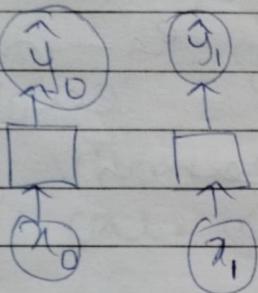
x_{t-1}^m \rightarrow Even if we stack perceptrons to form feed forward NN we still don't have notion of temporal inputs in sequential processing

Handling indiv. time steps

If we collapse the middle layer of stacked perceptrons into a single unit



- t denotes a single time step
- input at single time step generates single o/p using NN



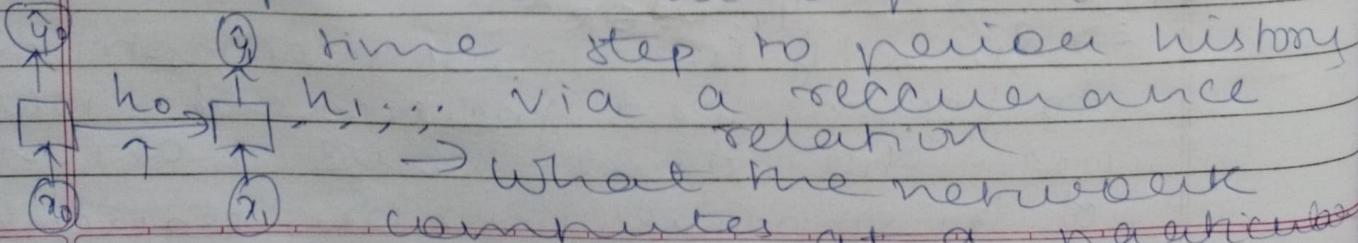
→ What if we do the operation multiple times to handle inputs fed in ~~one~~ corresponding to diff times

- We have an indiv time step starting with 0 and keep doing this repeatedly
- All these models are copies of each other ~~so~~

$$\hat{y}_t = f(x_t)$$

→ but here some of the later o/p would depend on some of the prev inputs if it is temporal dependence

→ If we link the prev NN with next one also to link the networks computations at a particular



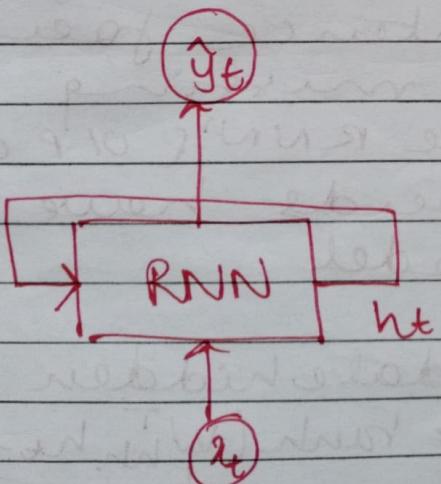
time step to prev history via a recurrence relation

→ What the network computes at a particular time step is passed to the later

→ we define this link by variable h (internal state) which we can think of as a memory team that's maintained by network & neurons of this state that's passed on from network to network as we process the seq. info

→ Networks O/P is not only a function of input x but also h which captures notion of memory that's being computed by the network & passed on over time $y_t = f(x_t, h_{t-1})$

→ We can also depict this relationship according to a loop -
→ This computation to this internal state variable h of t is being iteratively updated over time and fed back into neurons computation in the recurrence relation



RNN is maintaining the state & updating the state at each time step

* Apply recurrence relation at every time step to process a sequence

$$\text{at state } h_t = f_w(x_t, h_{t-1}) \text{ with weights } w$$

NOTE - Same function and set of parameters are used at every time step

RNN Intuition:- define RNN model

$\text{my_rnn} = \text{RNN}()$

$\text{hidden_state} = [0, 0, 0, 0]$ initialize hidden state

$\text{sentence} = ["I", "love", "recurrent"]$

Suppose "neutral"

we want to predict

the next word

of the sequence

which is our temporal sequence

for word in sentence:

prediction, hiddenstate =

$\text{my_rnn}(\text{word}, \text{hidden_state})$

next word prediction = prediction

we feed the word

to the RNN

→ This generates along with pred. for next word & prev hidden state updates the RNN's hidden state in turn

→ Finally our pred. for final word in sentence for word that we are missing is simply the RNN's output after all prior words have been fed in the model

Input vector

x_t

Update hidden state

$$h_t = \tanh(W_h^T h_{t-1} + W_x^T x_t)$$

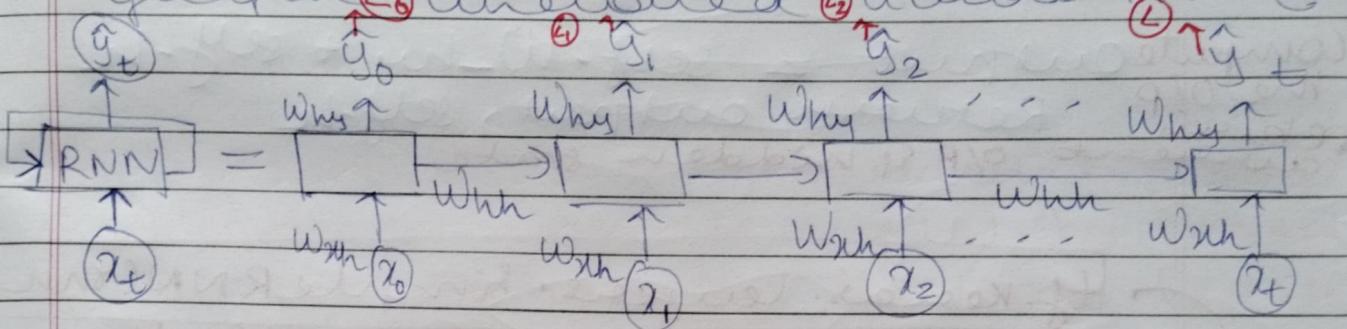
we have 2

separate weight matrices that the network is learning

Output vector = Why +
separate
wt matrix

Computational Graph Across time

→ Represent as computational graph (1) unrolled (2) across time (3)



→ prediction at an individual time step will simply amount to a computed loss at that particular time step

→ Now we can compare those predictions time step by time step to the true label and generate a loss value for those timesteps and finally get total loss by taking indiv. losses & summing & defining total loss for particular RNN

RNNs from scratch

class myRNNCell(tf.keras.layers.Layer)

def __init__(self, run_units, input_dim, output_dim):

super(myRNNCell, self).__init__()

{ self.w_xh = self.add_weight([run_units, input_dim])

self.w_hh = self.add_weight([run_units, run_units])

self.Why = self.add_weight([output_dim, run_units])

initialise
wt
matrix

Dr. Nitish

$\text{self_h} = \text{tf.zeros}([\text{num_units}, 1])$ hidden states to zeros

update def call (self, x):

the hidden $\text{self_h} = \text{tf.math.tanh}(\text{self}\cdot\text{W_hy} *$
state $\text{self}\cdot\text{h} + \text{self}\cdot\text{W_xh} * x)$

compute output = $\text{self}\cdot\text{W_hy} * \text{self}\cdot\text{h}$
return Review output, self.h
current O/P & hidden state

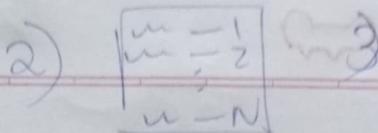
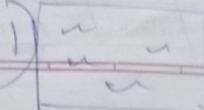
C If keras.layers.SimpleRNN(num_units)
tensorlayer abstracts the
operation of defining RNN
so we can define simple
RNN by this call

Sequence modelling: Design Criteria
Eg diff length sentences same dependencies far away from each other
ie same wts used in diff time steps will give a meaningful prediction
Represents language to NN
→ NN cannot interpret words and require numerical inputs

Encoding language for NN

C Embedding: transforms indexes into vector of fixed size
transforms indices or something that can be rep. as an index into a numerical vector of a given size

Embedding can be used for words as:-



PAGE NO.	111
DATE	

Vocabulary indexing
(all words) word index

are not embedding
 $w_{cat} = [0, 1, 0, 0, 1]$
7 in index

Instance of me
word corresponds to a vector
at its index

This is a very sparse method
simply based on count index

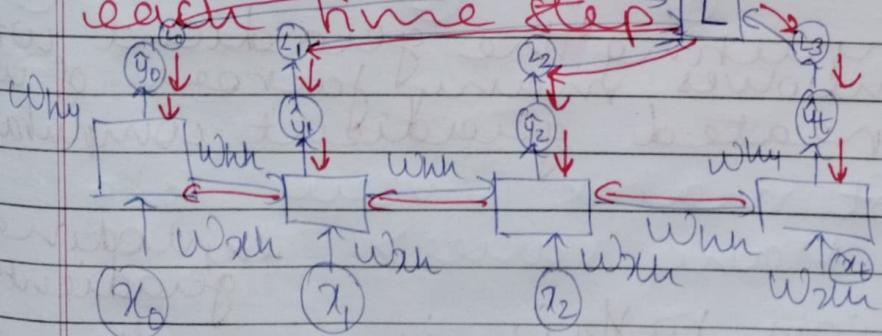
⇒ Alternatively, commonly what is done is to use a NN to learn in embedding. We can learn a NN that captures some meaning in the S/P data & maps related words/ inputs close in this embedding

Handle Variable sequence length

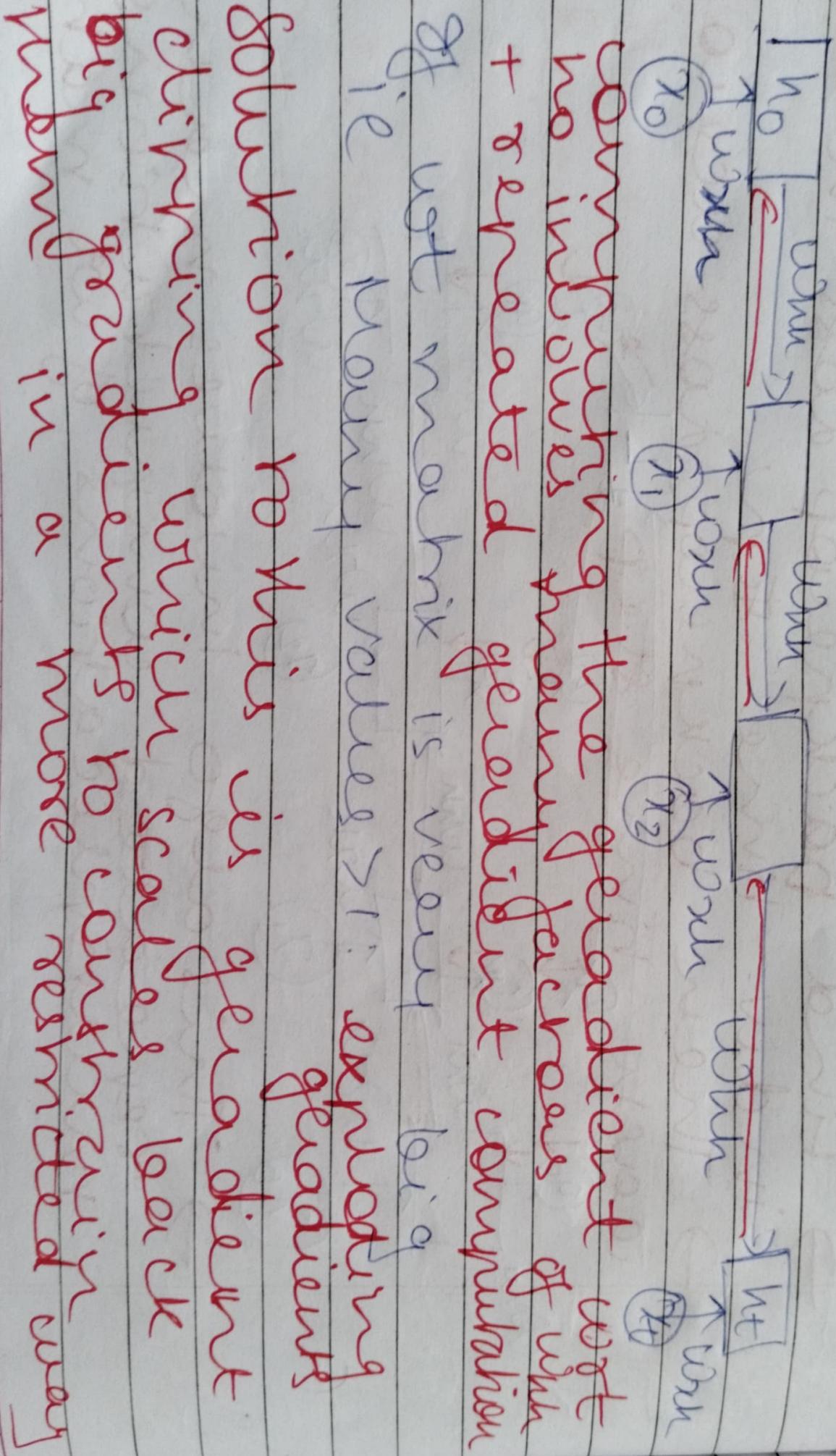
Backpropagation through time (BPTT)

→ find backprop loss in each indiv. time step

→ Then backprop loss across each time step



→ This algo involves many repeated computations & multiplications of wt matrices



Conversely the wt matrices could be very very small and we end up with vanishing gradients

Many values < 1 :

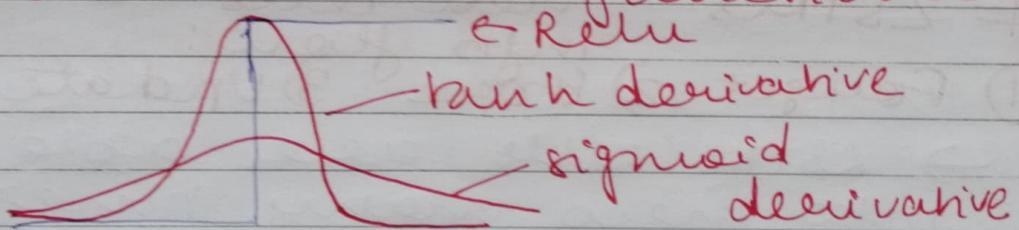
vanishing gradients

Tools to mitigate vanishing grad
1) Activation func 2) wt initialization 3) Network architecture

Why are vanishing gradients a problem? (Problem of long term dependencies)
Multiplying many small numbers together \rightarrow Errors due to further back time steps have smaller & smaller gradients \rightarrow Bias parameters to capture short term dependencies \rightarrow RNN becomes unable to connect the dots & establish long term dependencies

Trick 1 Activation functions

\rightarrow ReLU



ReLU

tanh derivative

sigmoid

derivative

Using ReLU prevents J from shrinking gradients when $x > 0$ as its derivative is 1

Trick 2 Parameter initialization

\rightarrow Initialize weights to identity matrix

Initialize biases to 0

$$J_n = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \end{bmatrix}$$

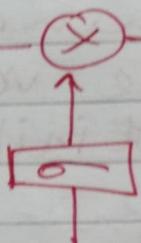
This helps prevent wts from shrinking to 0

Trick 3 gated cells - most robust solution

Idea: To use gates to selectively add or remove info within each recurrent unit with

Pointwise multiplicative

Sigmoid neural net layer



gated cell
[LSTM, GRU, etc]

Gates optionally let info thru the cell

→ We can filter out what's not important while maintaining what's imp.

→ The most popular ^{type of} Recurrent unit that achieves this gated is LSTM (long short term memory) - networks rely on gated cell to track info throughout many time steps

LSTM's These cells control info flow:

- 1) Forget
- 2) Store
- 3) Update
- 4) O/P

→ LSTM cells are able to track info thru many timesteps

[tf.keras.layers.LSTM(num units)]

Key concepts 1) Maintaining a cell state - just like a std RNN & that cell state is independent from what is ^{directly} outputted

2) use gates to control flow of info (for updating cell states)
• Forget gate gets rid of irrelevant info

- Share relevant info from current \hat{y}_t
- Selectively update cell state
- Output gate receives a filtered version of the cell

③ We can train the LSTM using backprop thru time algo. Mathematics of how LSTM is defined allows for uninterrupted flow of gradients which completely eliminates vanishing grad. problem

RNN applications & limitations

Eg 1 - sentiment classification
 $I/P \rightarrow \text{seq. of words}$ O/P - Prob of having the sentiment
 $\text{loss} = \text{tf.nn.softmax_cross_entropy_with_logits}(y, \text{predicted})$

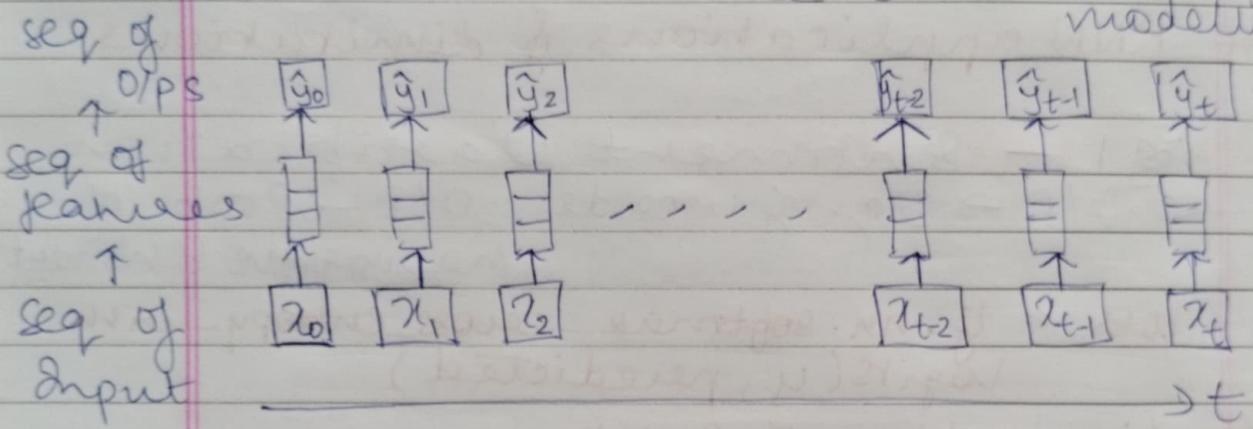
limitations of RNNs

1) Encoding bottleneck - RNN require that seq. info is taken in & processed time step by time step which imposes an encoding bottleneck - we're trying to encode a lot of content into a single O/P that may be just at the very last time step. We can't ensure that all that info leading up to that time step was properly maintained & encoded

- 2) RNN's can be quite slow, there is no easy way to parallelize that computation
- 3) Capacity of RNN & LSTM is not that long so we can't really handle data of tens of thousands or even beyond that effectively to learn the complete info & patterns present within a rich data src

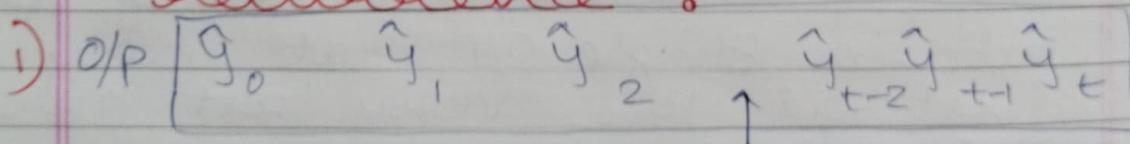
Goal of sequence modelling

RNN - recurrence to model sequence modelling



Desired capabilities - 1) continuous stream 2) Parallelization 3) long memory

Can we eliminate need for recurrence?



feature vectors
concatenate
feature vec
of all time steps

disadv - not scalable, no order, no long memory

2) Attention - foundational mechanism of transformer architecture

* Inhibition behind self-attention - attention to the most imp. parts of an input

- 1) Identify which parts to attend to
- 2) Extract features with high attention

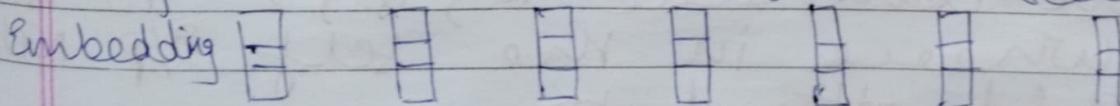
STEP 1 compute attention mask : how similar is each key to desired query

STEP 2 : Extracting values based on attention : Return me values of highest attention

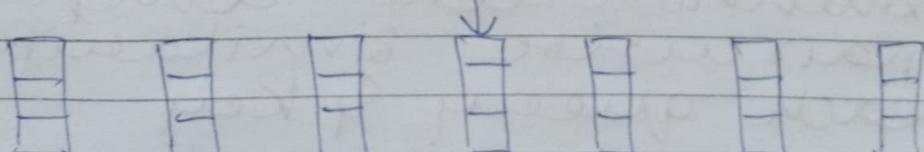
Learning seq attention with NN
goal - identify & attend to imp features in input

i) Encode position info

Eg We need the tennis ball to score



Position info



Position - aware encoding

* Data is fed in all at once
so need to encode position info to understand order

2) Extract query, key, value for each

$$\begin{array}{c} \text{Positional embedding} \\ \times \end{array} \begin{array}{c} \text{linear layer 1} \\ = \end{array} \begin{array}{c} \text{Query} \\ \text{o/p} \end{array}$$

$$\begin{array}{c} \cancel{\text{Same}} \\ \text{Positional embedding} \\ \times \end{array} \begin{array}{c} \text{linear layer 2} \\ = \end{array} \begin{array}{c} \text{K} \\ \text{Key} \end{array}$$

(These are
a diff set
of weights)

$$\begin{array}{c} \text{Positional embedding} \\ \times \end{array} \begin{array}{c} \text{linear layer 2} \\ = \end{array} \begin{array}{c} \text{V} \\ \text{value} \end{array}$$

(also
diff set of
wts)

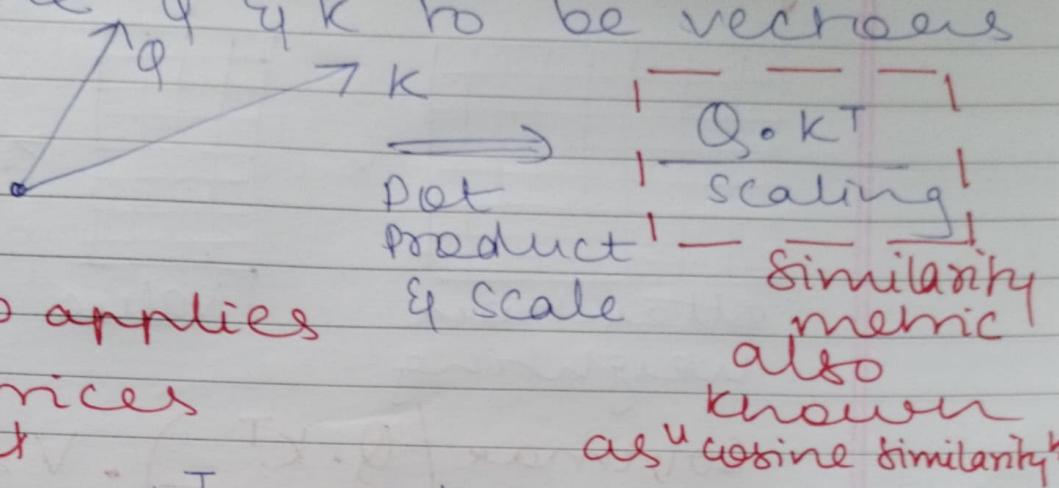
3) Compute attention weighting

We compare K, Q, V to see
where in the self o/p Network
shld attend to

Attention Score: Compute
pairwise similarity b/w
each query & key

To compute similarity b/w 2 sets of features:

→ assume Q & K to be vectors



This also applies
to matrices

$$Q \cdot K^T$$

query key

$$\frac{\text{Dot Product}}{\text{Scale}} \Rightarrow \frac{Q \cdot K^T}{\text{scaling}}$$

Captures how similar these vectors are and whether or not they point in the same direction.

Attention weighting: where to attend to! How similar is k to q

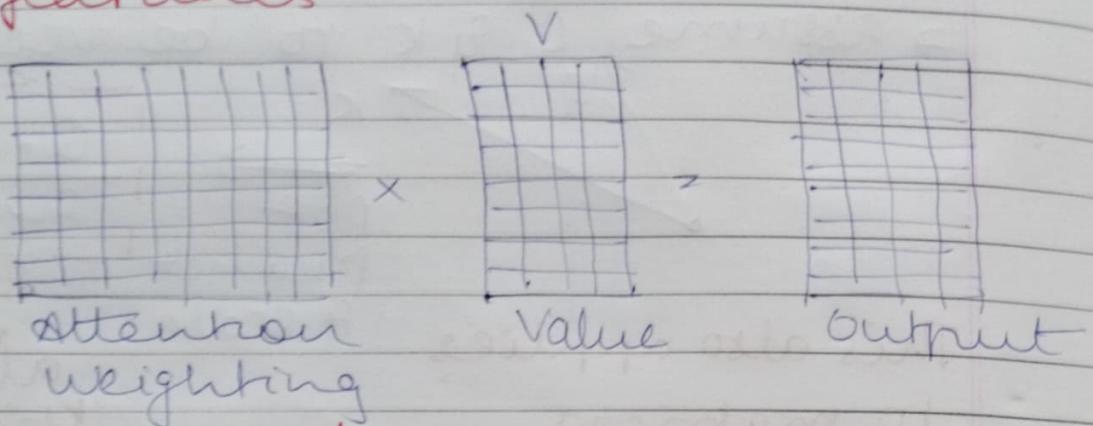
re	tossed
the	tennis
ball	serve
no	

$$\text{softmax } \frac{Q \cdot K^T}{(\text{scaling})}$$

attention weighting

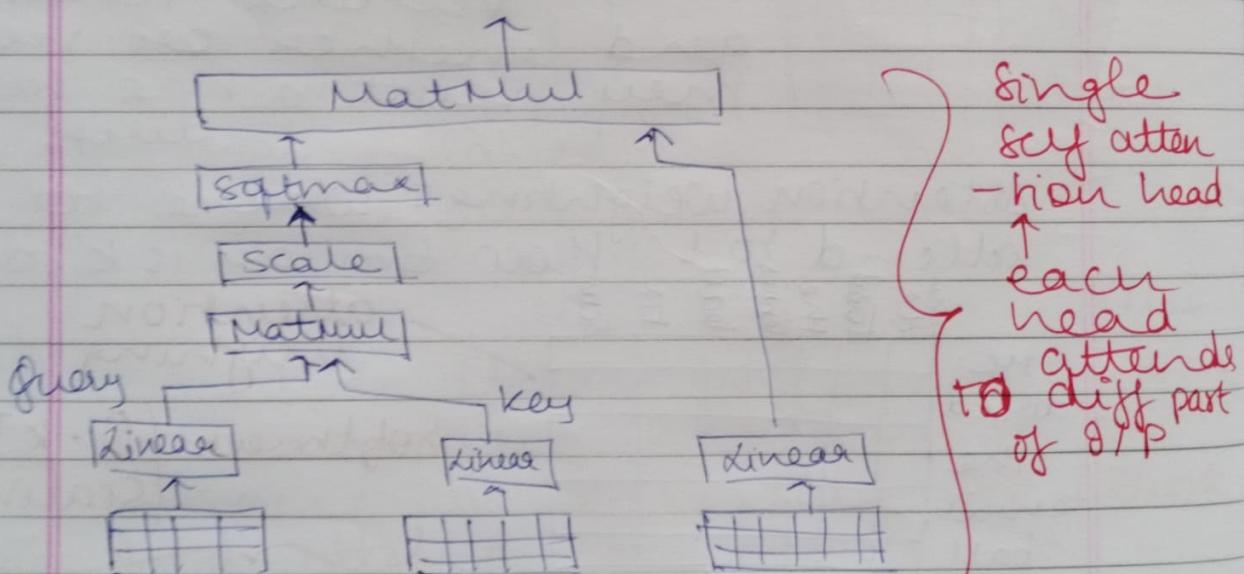
words related to each other in the sequence shld have high attention wts

4) Extract features with high attention: self-attend to extract features



$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right) \cdot V = A(Q, K, V)$$

Summary



These heads can plug into a larger network