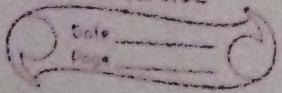


# RNNs, Transformers & Attention

## SEQUENTIAL DATA

classmate

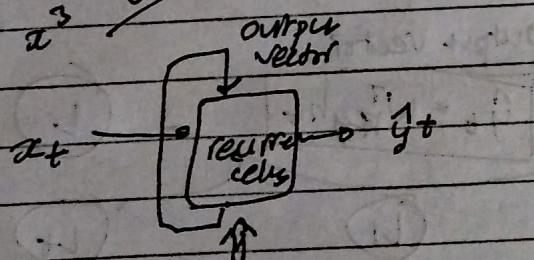
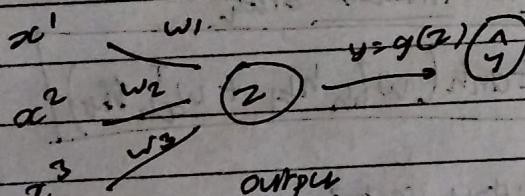


## Sequence Modelling Applications

- One to one      Many to one      One to many
- ▷ Binary Classification      ▷ Sentiment Analysis      ▷ Image captioning
- Many to many
- ▷ Machine Translation

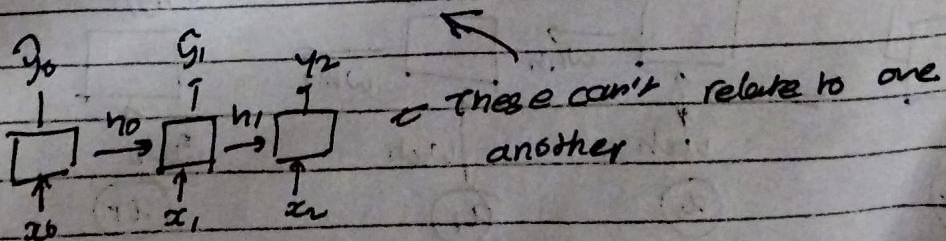
## Neurons with Recurrence

### Perceptron Revisited



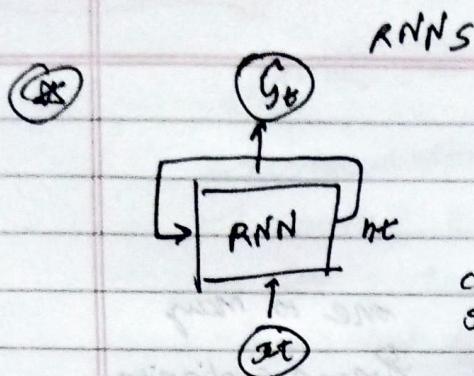
### Handling Individual Time Steps

$$g_t = f(x_t)$$



$$\hat{y}_t = f(x_t, h_{t-1})$$

output    input    past memory

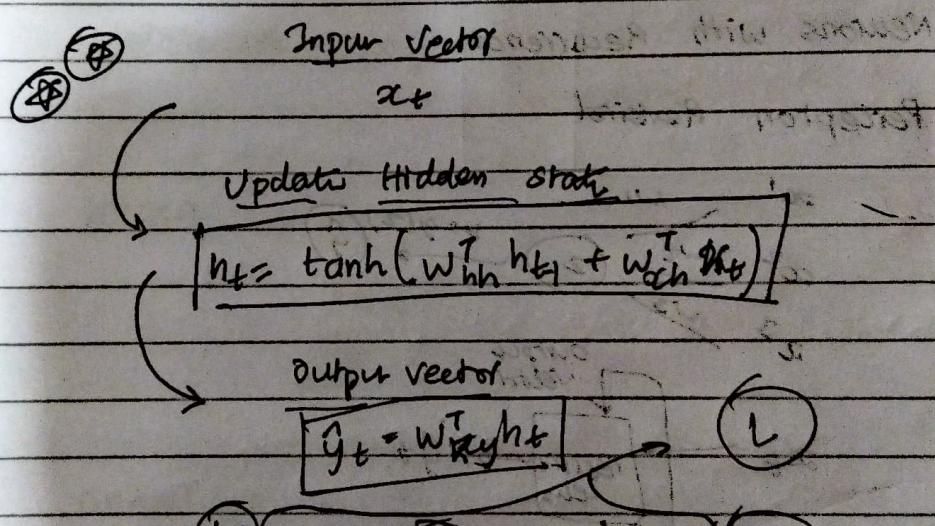


Apply a recurrence relation at every time step to process a sequence

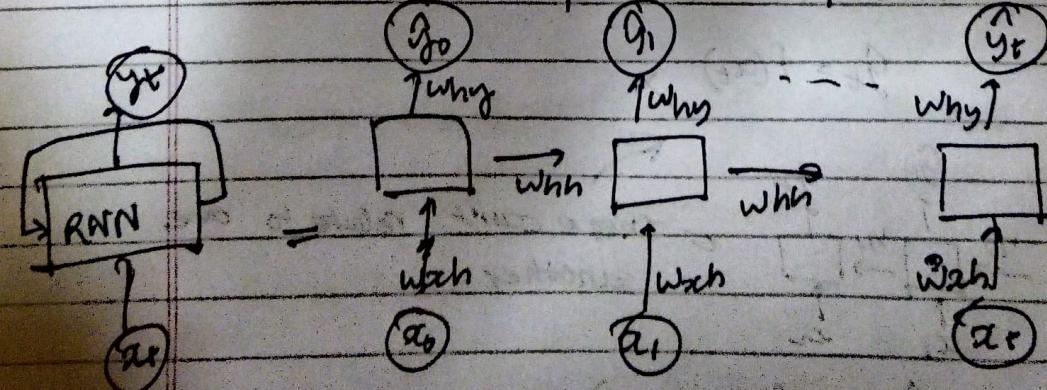
$$h_t = f_w(x_t, h_{t-1})$$

↑      ↓  
input    old state  
functions  
with weights  
 $w$

Same function & same set of parameters are used at every time step



RNNs: computational graph across time



## tf. keras.layers. SimpleRNN (mn-units)

RNN's for

### Design criteria - sequence modelling

1. Handle variable-length sequences
2. Track long-term dependencies
3. Maintain information about order
4. Share parameters across the sequence.

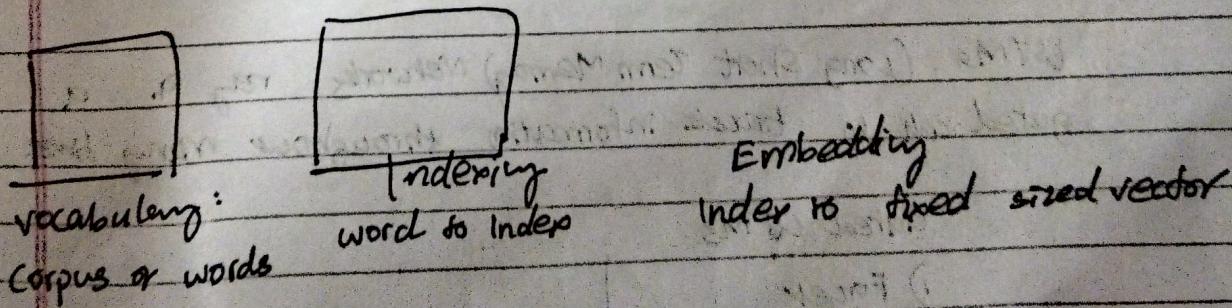
Predict the next word

This morning I took my car for a walk.

X "deep" X  $\rightarrow$  [ ]  $\rightarrow$  "learning"

Neural networks cannot interpret  
words

Embedding: transform indexes into vector or fixed size



## Backpropagation Through Time (BPPT)

### Backpropagation algorithm

- 1) Take the derivative w.r.t to each parameter
- 2) Shift parameters in order to minimize loss

Many values $\rightarrow$ Exploding gradients	many values $< 1$ : vanishing gradients
Gradient clipping to scale big gradients	1. Activation function 2. Weight initialization 3. Network Architecture

#1

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

#2 Initialise weights to identity matrix

Initialise biases to zero

#3 Gated cells

LSTMs (Long Short Term Memory) Networks use a gated cell to track information throughout many time steps

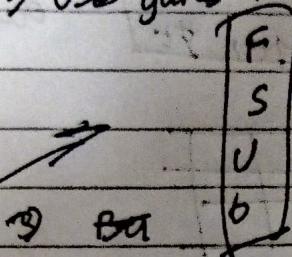
#### Gated LSTMs

- 1) Forget
- 2) Store
- 3) Update
- 4) Output

able to track through information throughout many timestamps

## LSTM Key concepts

- 1) Maintain a cell state
- 2) Use gates to control flow of information



- 3) Backpropagation through time will perturb uninterrupted gradient flow

## RNN Application & Limitations

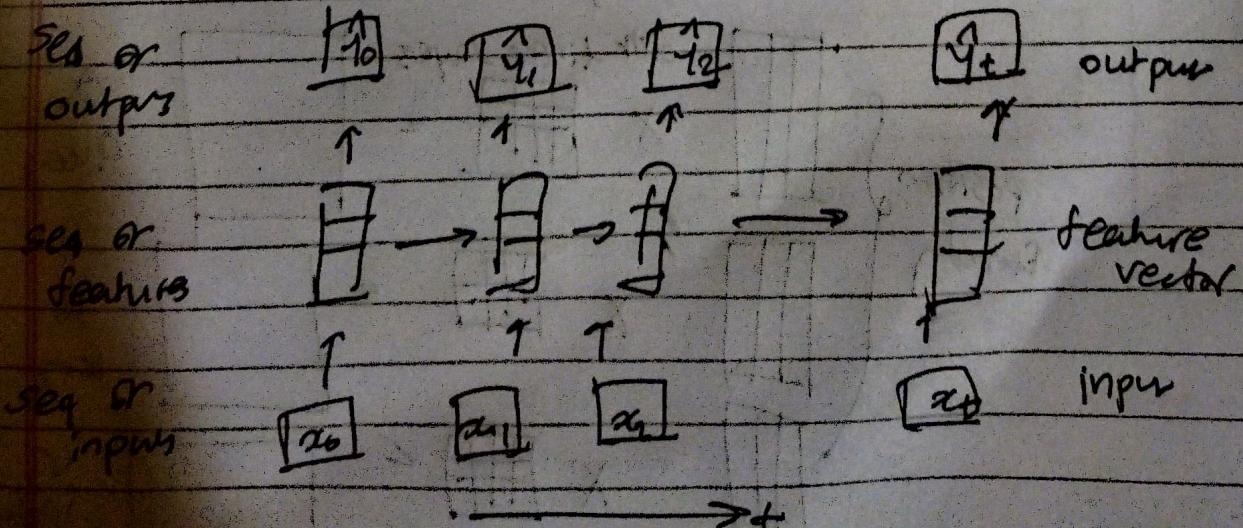
- 1) Music Generation
  - Input: sheet music
  - Output: next character in sheet music
- 2) Sentiment classification

### Limitations of RNNs

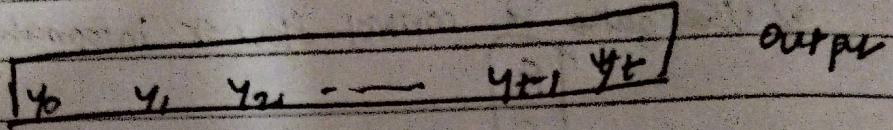
Encoding bottleneck

Slow, no parallelisation

No long memory



## Naive Approach

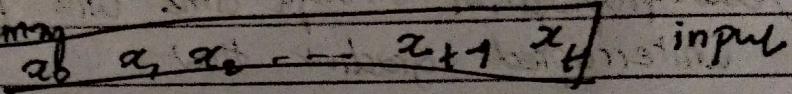


✓ No recurrence

✗ Not scalable

✗ No order

✗ No long term memory



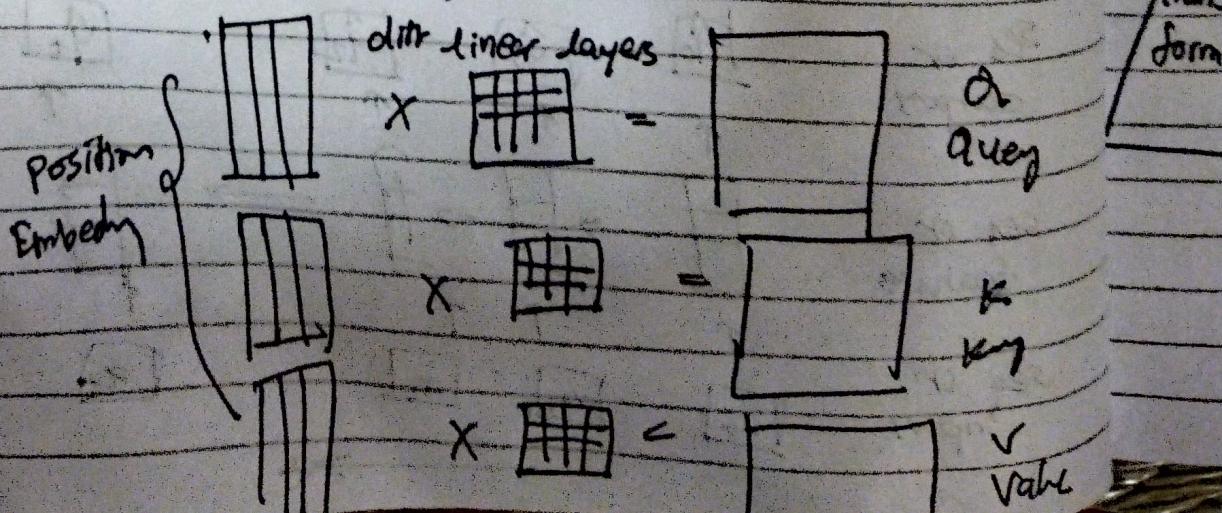
## Attention

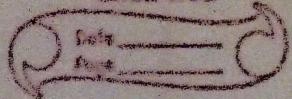
↳ foundational mechanism of transformer architecture

## Learning self Attention with Neural Networks

Goal: Identify and attend to most important features in inputs

- 1) Encode position information
- 2) Extract query, key, value for search
- 3) Compute attention weighing
- 4) Extract features with high attention

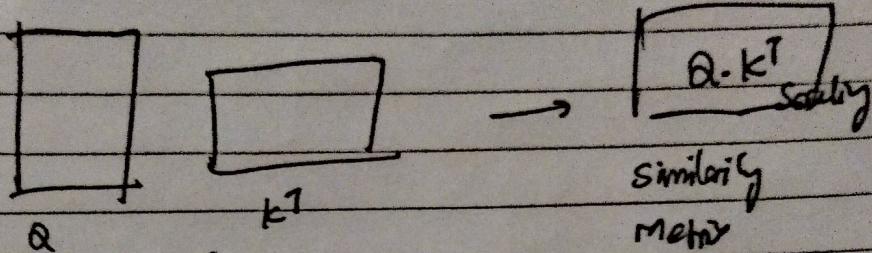




computing pairwise similarity bw query & key

$$\begin{matrix} Q \\ K \end{matrix} \xrightarrow{\text{Dot Prod}} \underline{Q \cdot K^T}$$

Scaling



Query                  Key                  Softmax  $\left( \frac{Q \cdot K^T}{\text{Scaling}} \right)$

Attention weighting

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

↑ value matrix

