

Deep Sequence Modeling

- Audio is a kind of sequential data.
- Useful in analyzing medical signals.
-

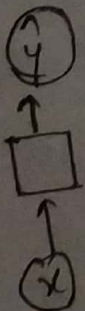
Design criteria:

- ↳ Handle variable length sequences
- ↳ Track long term dependencies
- ↳ Maintain information about order
- ↳ share parameters across the sequence

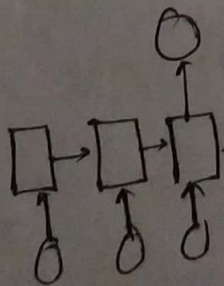
* ↳ RNN as an approach to sequence modeling problems.

RNN

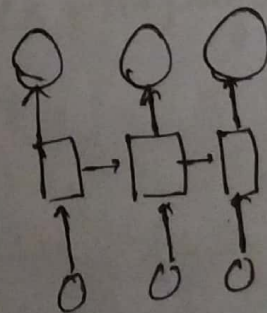
* standard feed forward neural network



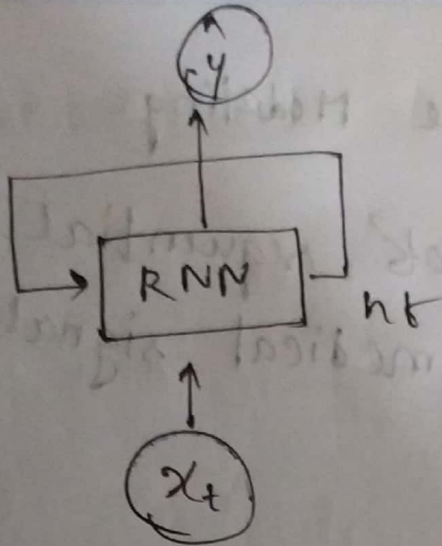
one to one



many to one
(sentiment)



Many to many



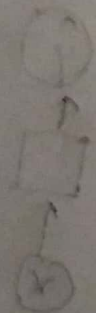
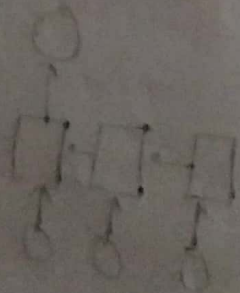
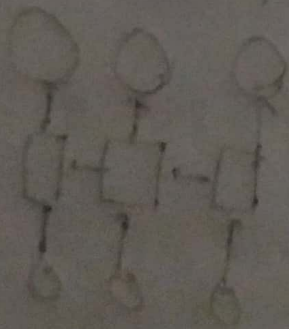
apply a recurrence relation at every time step to process.

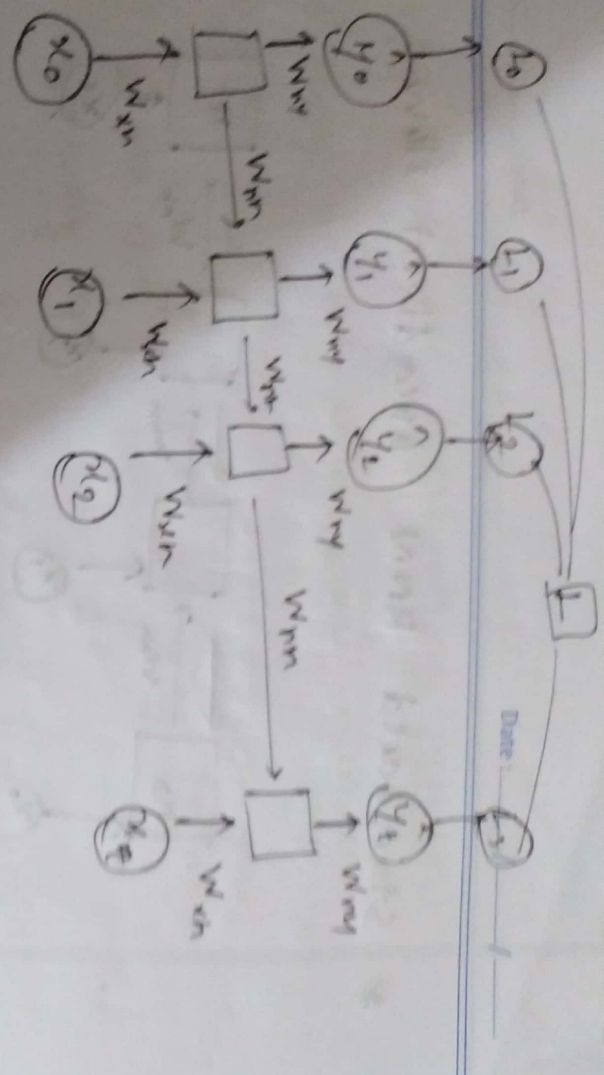
$$h_t = f_w(\underbrace{h_{t-1}}_{\text{old state}}, \underbrace{x_t}_{\text{input vector at } t})$$

cell state

function parameterized by w

$$\begin{aligned} \hat{y}_t &= W_{ny}^T h_t \\ h_t &= \tanh(W_{hn}^T h_{t-1} + W_{xn}^T x_t) \end{aligned}$$



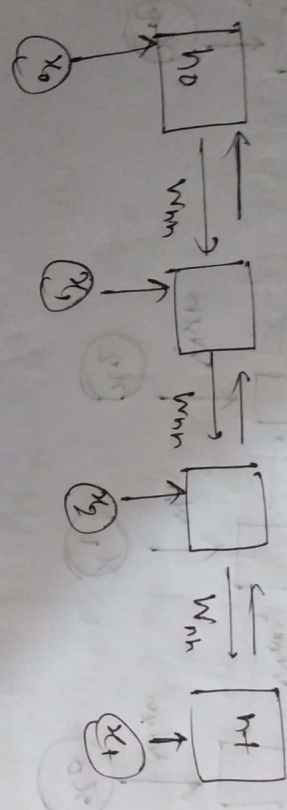


Backpropagation Through Time (BPTT)

- Take the derivative of the loss w.r. to each parameter
- shift parameters to minimize the loss

In RNN, errors are backpropagated in each individual time step

* Standard RNN gradient flow:



Computing h_0 involves many factors of w_{hn} + repeated gradient computation

- repetitive use of gradient can be problematic

Many values > 1 : exploding gradient problem

→ gradient clipping to scale big gradients

Many values < 1 : vanishing gradient problem

- ① Activation function
- ② Weight initialization
- ③ Network architecture

Why vanishing gradients a problem?

→ ~~many~~ multiply many small numbers together

↓
as the number will shrink, it will eventually vanish

↓
errors due to further back time steps have smaller gradients

↓
biasing the network to capture short term dependencies

Solution

1. Activation function:

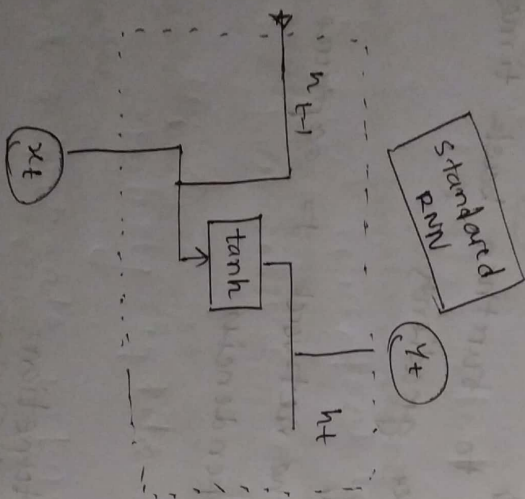
Using ReLU prevents it from shrinking the gradients when $x > 0$

2. Parameter initialization:

initialize weights to the identity matrix helps to prevent them to shrinking to zero

3. graded cells: use more complex recurrent unit with gates to control what information is passed through

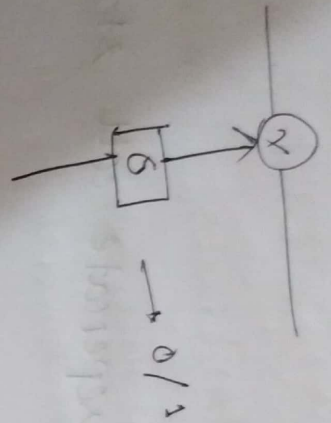
ex - LSTM, GRU



LSTM:

— contains computational blocks that control information flow.

— LSTM cells are able to track information throughout many time steps

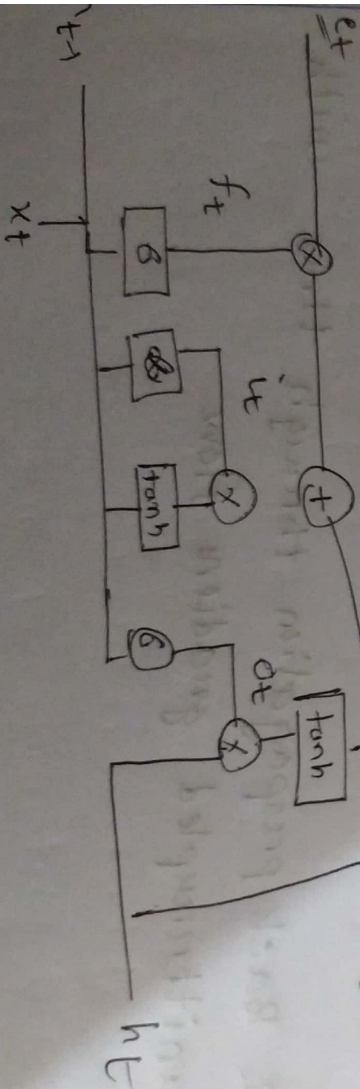


— Information is added or removed through structures called gates.

consists of a neural net layer
i.e. sigmoid

How LSTM works?

- (i) Forget f_t
- (ii) Store i_t
- (iii) Update o_t
- (iv) Output a_t



* ~~the~~ Uninterrupted gradient flow.

Key concepts:

1. Maintain a separate cell states from what is outputted
2. Use gates to control information flow
 - ↳ Forget irrelevant information
 - ↳ Store information from current input
 - ↳ selectively update cells
 - ↳ output gate returns a filtered version of the cell state
3. Back propagation through time with uninterrupted gradient flow.

RNN applications

1. Music generation
2. Sentiment classification
3. Machine translation
└ attention mechanism
4. Environmental modeling