Department of Mechanical Engineering
Gitam School of Technology
GITAM (Deemed to be University) Visakhapatnam

Category of course    Theory                    Semester    II
                      19EOE746
                      OPERATIONS RESEARCH

# Moodle - 3

①

## RNN

→ Recurrent Neural Nllw

→ Mainly used in Speech recognit$^n$ & NLP

remembers only previous data

→ follows sequential approach

→ it trained the process and convert a sequential data i/p ♦ into a specific sequential data o/p.
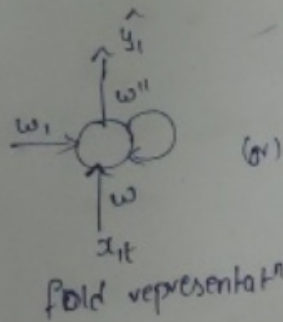
it is data,
such as words, sentences
or time series data.

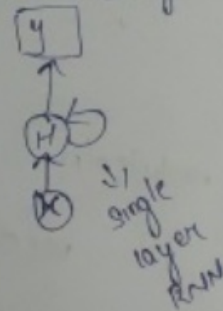Backpropogat$^n$ is used for weight adjustment
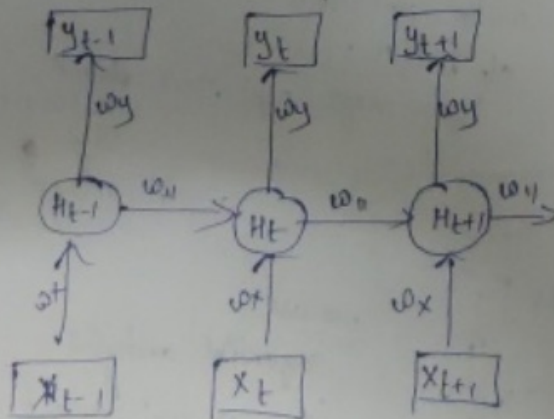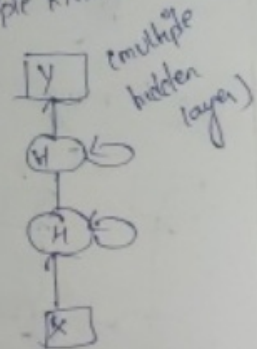
→ We use Backpropogat$^n$ algorithm.

Why RNN



fold representat$^n$

→ Multiple hidden layers untaiyee

(or)

single layer RNN

Multiple RNN

(multiple hidden layer)

Types of RNN

i) one to one
ii) one to many
iii) many to one
iv) Many to many

multiple i/p

$[y_1]$  $[y_2]$  $[y_3]$

$(H) \rightarrow (H) \rightarrow (H)$

$[X_1]$  $[X_2]$  $[X_3]$

multiple i/p

→ Best te...
is seq an improver
suited for seq
dencies

---

one to one :-

$[y]$ single o/p

$(H)$

single i/p $[X]$

→ this is called vanila neural n/w
→ used for regular to data

⊙ Global $\frac{\partial L}{\partial w_4} = \frac{\partial L}{\partial y}$ traditional RNN

⊙ $w_3 = \frac{\partial L}{\partial y}$ sysTMs address

⊙ weight adjust $\frac{\partial L}{\partial y}$ men that co

using chain rule $\frac{\partial L}{\partial c}$  STM n/w
$\frac{\partial L}{\partial y}$ data, wh
speech

⊙ $\frac{\partial L}{\partial k_1}$, $w_4$ $\frac{\partial L}{\partial g}$

⊙ $\frac{old}{new} = \frac{\partial L}{\partial y}$
old w
value

---

one to many

multiple o/p

$[y]$  $[y]$  $[y]$

$(H) \rightarrow (H) \rightarrow (H)$

$[X]$

single i/p

use case :-
image captioning
music generation
o/p is eq$^n$.
i/p — image o/p ·Text
description of
i/p of
of image

⊙

---

Many to one

single o/p $[Y]$

$(H) \rightarrow (H) \rightarrow (H)$

$[X_1]$  $[X_2]$  $[X_3]$

Multiple i/p

usecase
sequential analysis
i/p is seq$^n$ output is
fixed vector size

# LSTM

(3)

:- short term memory

- It is an improved vers$^n$ of rnn.
- Well-suited for seq$^n$ prediction tasks & excels in capturing long-term dependencies

- A traditional RNN has a single hidden state that passed through time which can take make it difficult for n/w to learn long-term dependencies

→ LSTMs address this problem by introducing a <u>memory cell</u>, which is a container that can hold information for an extended period.

→ LSTM n/w are capable of learning <u>long-term dependencies in sequential</u> data, which makes them well-suited for tasks such as language translation, speech recognition & timeseries forecasting.
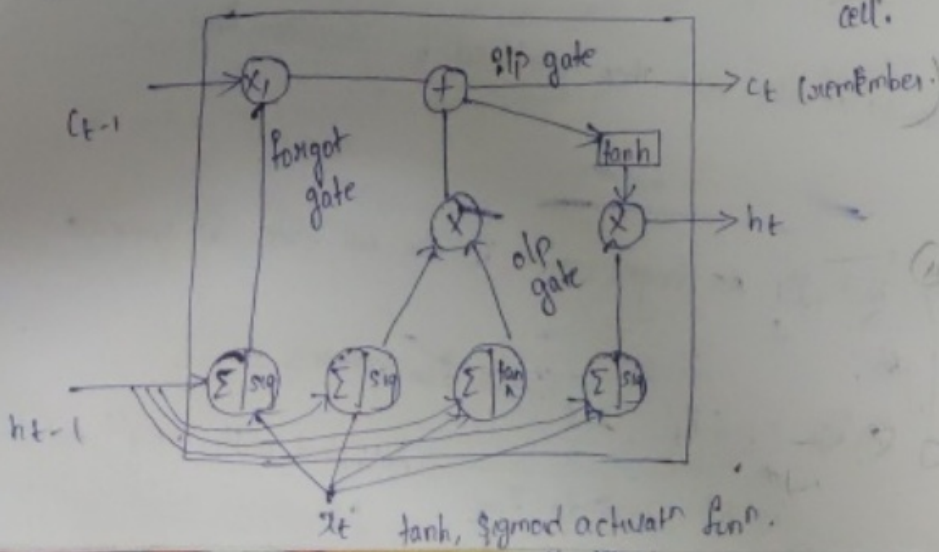
→ Memory cell is controlled by 3 gates :- i/p gate → ctrls what information is added to memory cell

ctrl what information is removed from memory cell.

forget gate

o/p gate.

↓

ctrls what information is o/p from memory cell.

LSTM.



ht-1

"It tanh, sigmoid activat$^n$ fun$^n$.

gradient

Types of RNN

Rnn form

$\rightarrow h_t$

tanh    tanh    tanh

$h_{t-1}$

$x_t$

## Bidirectional LSTM

$\rightarrow$ it is a RNN that is able to process sequential data in both forward & backward direction.

$\rightarrow$ allows Bidirectional LSTM to learn longer range dependencies in sequential data than traditional LSTMS which can only process sequential data in one direction

$\rightarrow$ Bi LSTM made up of two LSTM nlws.

    i) ilp seq$^n$ - in the forward direction

    ii) ilp seq$^n$ - in the backward direction

olp of two LSTM nlw are then combined to produce final olp.

forget Gate :-

    information that is no longer useful in the cell state is removed with forget gate.

    2 ilp $x_t$ & $h_{t-1}$ $\rightarrow$ olp cell from previous

    ilp at particular time $\rightarrow$ multipled with weight natrices followed by adding of bias

Department of Mechanical Engineering
Gitam School of Technology
AM ... University, Visakha...

Semester    II

...CH    Date    10.04.2024
90 Min

# GRU

(4)

$\rightarrow$ Recurrent unit.

$\rightarrow$ is a specialized variant of RNN

$\rightarrow$ applications :- NLP

Speech recognit$^n$

time series predict$^n$

$\rightarrow$ The GRU presents itself as an innovative sol$^n$ to the vanishing gradient problem in traditional RNN.

$\rightarrow$ The architecture of GRU is designed with two specific gates

    i) update gate $\rightarrow$ it recognizes long term connect$^n$

    ii) Reset gate.

        $\downarrow$
    identifies short term relat$^n$ships

The various components of architecture are :-

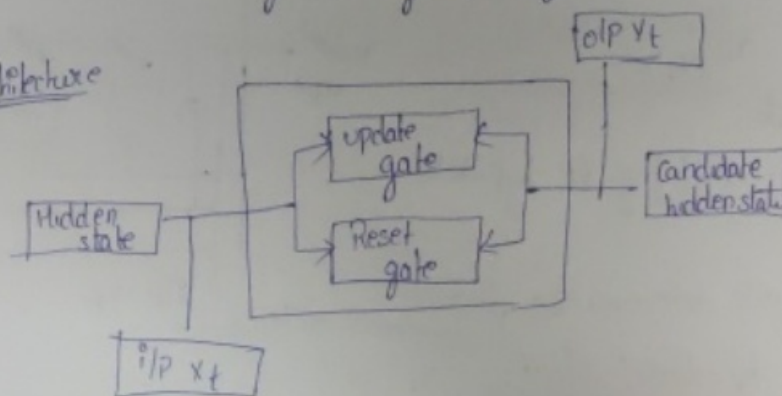    i) update gate $(Z)$ $\rightarrow$ determines the degree of past informat$^n$ forward to future

    ii) Reset gate $(R)$ $\rightarrow$ Decides the amount of past informat$^n$ to discard

    iii) Candidate hidden state $(H')$ $\rightarrow$ creates new representat$^n$ considering both i/p & past hidden state.

    iv) final hidden state $(H)$

        $\downarrow$ A blend of new & old memories governed by update gate.

architecture

update gate:-

→ 1st step of GRU

→ It employs the current i/p & previous hidden state, to decide how previous hidden state should be update

→ sigmoid fun$^n$ is used.

$$Z_t = \sigma\;(W_z \cdot [h_{t-1},\, x_t] + b_z)$$

update gate vector

weight matrix — hidden state — current i/p — bias

$$z_t = \sigma\,(\omega z \cdot [h_{t-1}, x_t] + b_z)$$
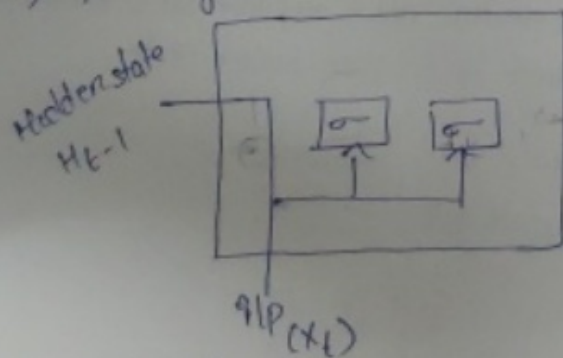
Reset gate:-

sigmoid fun$^n$ is used
similar to update gate
it identifies the volume of past informat$^n$ to be discarded.

$$r_t = \sigma\,(W_r \cdot [h_{t-1},\, x_t] + b_r)$$

reset gate vector

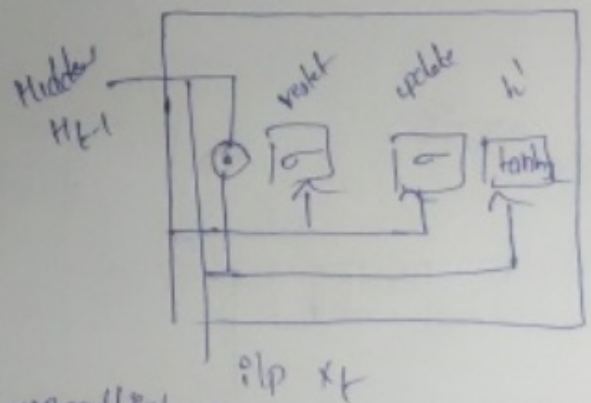weight matrix — s/bias

Reset & update gate follows



i/p $(x_t)$

Candiate hidden state

it is computed employing the hyperbolic tangent fun$^n$ (tanh).

value of reset gate determines the previous hidden state

$$h'_t = \tanh\,(W \cdot [r_t \times h_{t-1},\, x_t] + b]$$

→ elementwise multiplica$^n$
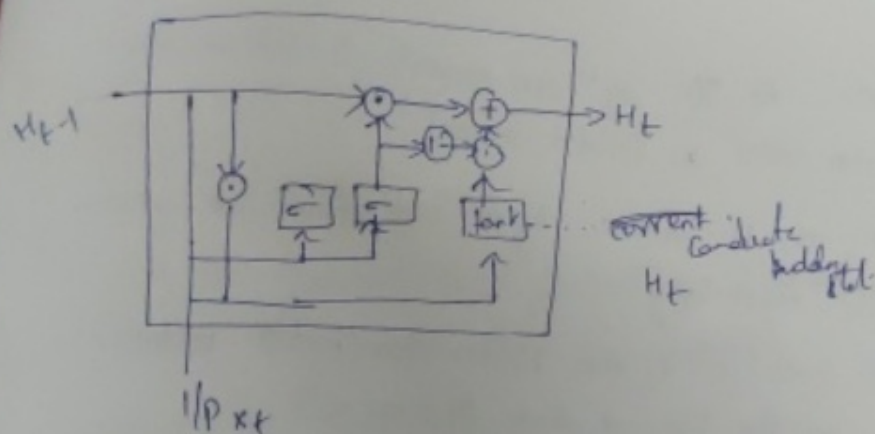


Hidden H$_{t-1}$   reset   update   h'

i/p $x_t$

den state:-

known as current hidden state

subsequently defined through linear interpolation

It involves the previous hidden states
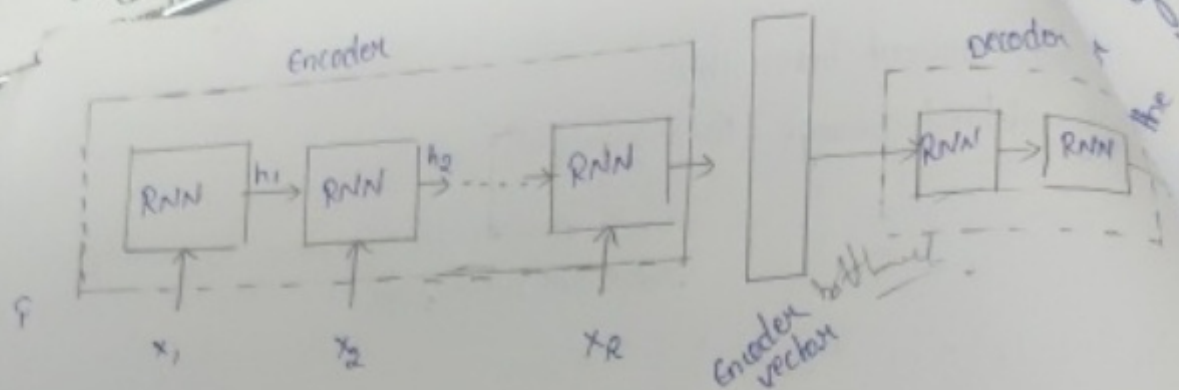Prospective hidden state.

$$h_t = (1 - Z_t) * h_{t-1} + Z_t * h_t'$$



$H_{t-1}$      $H_t$

current
conduct
$H_t$   hidden std.

I/p $x_t$

# Encoder-Decoder architecture:-

→ The encoder and decoder architecture is a powerful framework used for various seqⁿ to seqⁿ tasks in d.L

→ useful for tasks like m/c translatⁿ, where i/p & o/p seqⁿ have different lengths

→ architecture consists of two components :- i) encoder
                                     ii) Decoder.

Encoder

Encoder vector

$x_1$     $x_2$     $x_R$

## Encoder:-

→ The encoder processes the i/p seq$^n$ and converts it into a single dimensional vector (also called as hidden vector)

→ Multiple RNN cells can be stacked together to form the encoder.

→ RNN reads each i/p sequencially.

→ For everytime step (each i/p) $t$, the hidden state (hidden vector) h is updated according to the i/p of that timestep $x[i]$.

→ After all the i/ps are read by encoder model, the final hidden state of model represents the context/summary of whole i/p seq$^n$

## Encoder vector:-

→ This is the final hidden state produced from the encoder part of model.

→ The vector aims to encapsulate the information for all i/p elements in order to help the decoder make accurate predict$^n$

→ it acts as the initial hidden state of the decoder part of model

## Decoder:- (take encoder o/p vector as i/p)

→ The decoder takes the hidden vector produced by the encoder & generates the o/p seq$^n$ (translate the sentence in another language)

generates o/p seq$^n$ by predicting the next o/p $y_t$

the hidden state $h_t$

for the decoder is the final vector obtained at the end of encoder model.

each layer will have three o/ps hidden vector from previous layer $h_{t-1}$ —

and previous layer o/p $y_{t-1}$, original hidden vector h.

## o/p layer

use softmax activat$^n$ fun$^n$

used to produce the probability distribut$^n$ from a vector of values with the target class of high probability

o/p $y_t$ at time step t is computed using

$$y_t = softmax(W^s h_t)$$

## Application

text summarizat$^n$

speech recognition

Time series Application

Google's mlc Translat$^n$