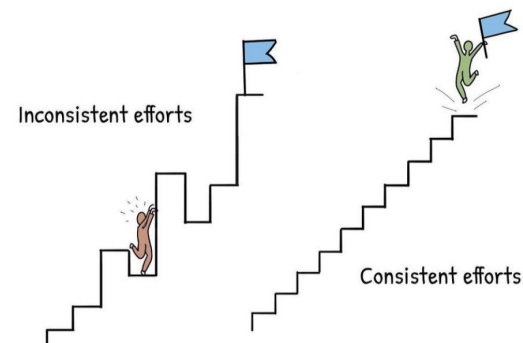# Deep Sequence Modelling

M.Tech. Data Science, Second Year, NMIMS
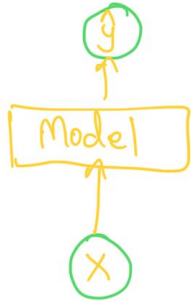
By,

Bilal Hungund, Data Scientist, Halliburton
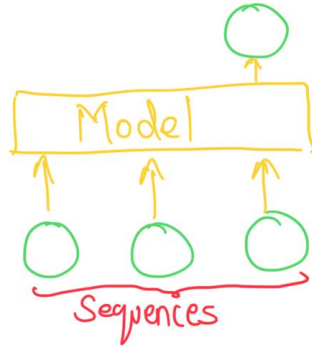
# Sequencing Modelling Applications

## One to one
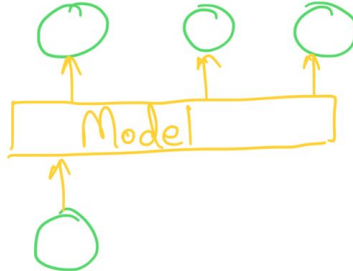(Classification)
Regression)

$y$

Model

$X$

## Many to One
(Sentimental Analysis)

Model

Sequences

## One to Many
(Image Captioning)

Model

Cat is sleeping

## Many to Many
(Language Modelling)

Model

Temporal Space

English to Hindi
Text

# Perceptrons Revisited



$W_0$

$W_1$

$W_2$

$W_3$

$g(z)$

$\hat{y}$

(No notion of time sequence yet)

# Feed Forward Models

$x_1$

$x_2$

$x_m$

$\hat{y}_1$

$\hat{y}_2$

$\hat{y}_n$

$x \in \mathbb{R}^m$

$\hat{y} \in \mathbb{R}^n$

# Handling with time steps

(How to relate network computations?)
Need some prior history

Output Vector $\hat{y}_t$

$\hat{y}_0$    $\hat{y}_1$    $\hat{y}_t$

Splitting t

$x_t$  Input Vector

$x_0$    $x_1$    ......    $x_t$

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

No interdependence or interconnectedness added yet

# Neurons with Recurrence

$\hat{y}_t$

Recurrent Cell

$h_t$ $\Rightarrow$

$x_t$

$\hat{y}_0$     $\hat{y}_1$     $\hat{y}_2$     $\hat{y}_b$

Linking the information and Computation of network

$h_1$    $h_2$    $h_3$

$x_0$     $x_1$     $x_2$     $x_t$

Recurrence relation

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

Output     Input     past memory

# Recurrence Neural Networks (RNNs)

$\longrightarrow$ RNNs have a state $h_t$, that is updated at each time steps as a sequence

i.e.,

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

- $h_t$ — cell state
- $f_w$ — weight
- $x_t$ — inputs
- $h_{t-1}$ — old state

| Input Vector | Update Hidden State | Output Vector |
|---|---|---|
| $x_t$ | $h_t = \tanh \left( W_{hh}^T h_{t-1} + W_{xh}^T x_t \right)$ | $\hat{y}_t = W_{hy}^T h_t$ |

# Simplifying RNNs

Losses

$L$

$L_0$     $L_1$     $L_2$     $L_3$

$\hat{y}_t$

$\hat{y}_0$     $\hat{y}_1$     $\hat{y}_2$     $\hat{y}_t$

RNN     $=$     $h_0$     $W_{hy}$     $h_1$     $W_{hh}$     $W_{hy}$     $h_2$     $W_{hh}$     $W_{hy}$     $b_3$     $W_{hh}$     $W_{hy}$

$W_{xh}$     $W_{xh}$     $W_{xh}$     $W_{xh}$

$x_t$     $x_0$     $x_1$     $x_2$     $x_t$

- - - $\rightarrow$ Backward Pass     $\longrightarrow$ Forward pass

# Gradient issues



$\longrightarrow$ Exploding gradients (values > 1)

$\longrightarrow$ Gradient Clipping : Scale big gradients

$\longrightarrow$ Vanishing gradients (values < 1)

$\longrightarrow$ Activation Function, Weight initialization, Network architecture

Vanishing gradients

$\longrightarrow$ Focusing on short-term dependencies and ignoring long term dependencies

Short term dependencies
$\underline{I}$ love neural $\underline{\quad ? \quad}$
(network)

$\xrightarrow{\text{Ignorance}}$

I studied data science and
I am fluent in $\underline{\quad ? \quad}$
(ML, DL, ...)

To alleviate the problem

$\longrightarrow$ Use of ReLU function as it prevents shrinking the gradients when $\alpha > 0$

$\longrightarrow$ Weight Initialization:

Initialize bias to zero and weights to identity matrix
$\underline{It}$ prevents the weights from shrinking to zero

$$I_n = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

$\longrightarrow$ Gated Cells : (LSTM)

Network Architecture
Use gates to add or remove information within each recurrent unit and to track information

# Design Criteria of RNNs

- Handle variable-length sequences
- Track long term dependencies
- Maintain order of information
- Sharing parameters

# Goals of Sequence Modelling

- Continuous stream
- Parallelization
- Long Memory

# Limitations of RNNs

- Encoding bottleneck
- Slow, no parallelization
- Not long memory

# NLP using RNN

Tokenization , pad Sequences and Embeddings

**Corpus/vocabulary**

F L O W
001 002 003 004

**Character Tokens**

W O L F
004 003 002 001

I love neural network
001 002 003 004

I love deep learning
001 002 005 006

Sequences
[[001, 002, 003, 004],
[001, 002, 005, 006]]'

Tokenization , pad Sequences and Embeddings

building vocabulary

I love neural network
I love deep learning
You love neural network
Do you think deep learning is good?

Tokens →

love - 1, i - 2, you - 3,
neural - 4, network - 5,
deep - 6, learning - 7,
do - 8, think - 9, good - 10
is - 11

Sequences

[2 1 4 5]
[2 1 6 7]
[3 1 4 5]
[8 3 9 6 7 11 10]

pad Sequences

$$\begin{bmatrix} 0 & 0 & 0 & 2 & 1 & 4 & 5 \\ 0 & 0 & 0 & 2 & 1 & 6 & 7 \\ 0 & 0 & 0 & 3 & 1 & 4 & 5 \\ 8 & 3 & 9 & 6 & 7 & 11 & 10 \end{bmatrix}$$

Pass to Embeddings