

# CS 452

## Deep learning - Internals

MS : 25

ES : 40 {open notes}

Q1, Q2 : 10, 10 {want for

Term paper : 10 {fixe file.

Demo : 10 {Project)

[Attendance  
Compulsory]

Absolute grading.

90+ : A A B

80+ : B B

60+ : C C {lebanon, multithreaded code.

40+ : D D {Infered file

### # Two main paradigms of Processing:

1) Algorithmic based on given rule  
→ step-by-step instructions follow rule  
→ Eg. Processing course registration request, sorting numbers

2) Model based on learned lesson rule.  
→ Handles Uncertainty, make predictions  
→ Eg. Choosing your elective course  
predicting age of a person

Algorithms

Models

works for every input

they work for certain contexts

You can always explain how you particular at a result

\* "All models are wrong, but some are useful" - George Box

→ Models simplify reality.

→ we don't know the algorithm

→ we might not know if the algorithm still exists.

Then why we use it?

→ Even if we don't know the true algorithm, models can still be useful.

## \* Supervised learning

→ we have to gather data

→ we collect data with labels.

## Unsupervised learning

→ the model learns on its own.

→ it learns patterns without labels.

→ sometimes the distinction is blurred.

→ It's a myth - we always use some form of supervision (Supervised learning), directly or indirectly.

## # Algorithmic Solution

It might require too much data and too much time consuming.

- # Model Based Examples :-
  - Projectile Motion
  - Relational data model
  - Random Access Machine (RAM) model
    - one instruction at a time
    - each instruction takes equal amount of time.
  - used for data storage in tabular form
- Relational Databases :
  - forcefull conversion of bases into rows & columns.
- ↳ These models help abstract away from low-level details and focus on patterns.

- 
- # Statistical Models :-
  - Why do we build statistical models ?
  - ↳ Built when we can't write rules explicitly.
  - ↳ Try to model underlying patterns from the data.

# Explainable vs Non-Explainable Models :-

\* Model Parameters :  
(Statistical Models requires various Parameters)

→ Find the appropriate value of the Parameters (like weights in ML).

→ Process :

→ Initial Guess

→ Predict Output

→ Compare with actual output (calculate error)

→ Update guess

(using learning algorithms).

Example :-

Input	Output
1	2
2	6
3	10
10	30

Shows a pattern in data that we want a model to learn.

• linear Model

$$y = ax + b$$

• Parameters :  $a, b$

make initial guess :  $a=1, b=0$ ,  
use a cost function (like MSE)

to measure error.

- Cost function: To predict how better the guess is.
- Gradient descent: To update a & b.  
we need to reduce cost function but it is not the only thing.  
~~(It's like an exam)~~

~~Real world~~ ~~exam~~ ~~vs~~ Real world performance.

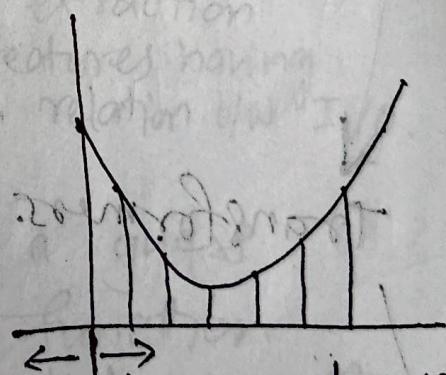
- Gradient Descent
  - Stochastic Gradient Descent (SGD)
  - Learning rate
    - Extending it to an ANN.
- Types of Gradient Descent
  - Batch
  - SGD.

\* Learning rate

→ Controls how much we move in the direction of minimum error.

→ Too small: slow learning.

→ Too large: may overshoot the optimum



How much we go right?

learning rate?  
(we will take baby steps).

(Ug, Ug))

gradient

gradient descent

gradient descent, gradient step, gradient

# Overview (course)

## Model Based Solutions (Paradigms)

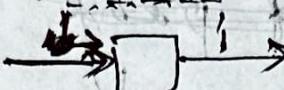
### Feed-forward Neural Network

Artificial



### Backpropagation

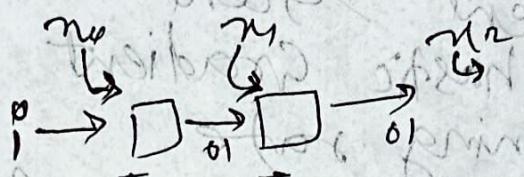
### Recurrent Neural Network [Vanilla,



LSTM

GRU

### Attention mechanisms



### RNN Configuration

History (Recurrent)

How to selectively  
use the  
history  
(Attention Mech.)

RNN - ~~Recurrent~~ Recurrence

Keep Attention

### Transformers.

### Encoder / Decoder

① Encoder Only (BERT)

② Decoder Only (LLMs)

### Hardware (CPU, GPU)

### various Open problems / issues

Inference, Fine Tuning, Bias, Ethics

- Problems where algos can't be used
- ① No algo's ~~exist~~ is possible
  - ② Algo. might exist, not ~~use~~.
  - ③ Algo. exists, but
    - i) Inefficient (uses too much resources, time, memory, network bandwidth)
    - ii) we want a diff. insight

∴ we abandon Algo. & uses model based solutions.

→ Backbox

Input | Output

ID | Yes/No

large training dataset

from input/output  
guess the logic  
behind the solution.

traditional way:  
features extraction  
find features having  
high relation b/w I/O.

$f(d)$

↓  
features

Embeddings

① Make a guess

② Cost function

(with the current  
model did I get the  
same output)

e.g.,  $\sum_{i=1}^{ID} (f_i - o_i)^2$

for distance  
checking

③ Change the parameters  
to reduce the error (cost  
function)

$$K \frac{\partial E}{\partial a} \text{ learning rate}$$

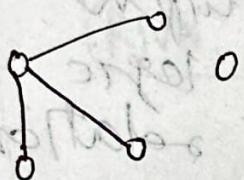
$$a = a + \Delta a$$

$$-\eta \frac{\partial E}{\partial a}$$

learning rate required so that we don't change the parameters abruptly, slow changes are required.

## # ANNs :-

Neurone in brains  $\leftrightarrow$  Nodes in graph connected to other nodes.



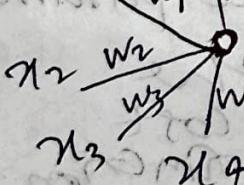
$$\frac{10^{11} \times 10^4}{2} = \frac{10^5}{2} = 5 \times 10^4$$

parameters in human brain

$10^3$  switching time human brain in

$10^9$  switching time for computer models.

=> Perceptron (not fake input) This neuron receives 4 inputs



Pass : 1  
Fail : -1

marks of student

what is thresholded perceptron

thresholded perceptron =  $\frac{\downarrow \text{Inputs}}{(+1, -1 \text{ OP})}$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

Parameters

$\bar{x} \cdot \bar{w} > 0$  then  $O/P = 1$  output of  
 $f(\bar{x}, \bar{w})$  else  $O/P = -1$  a perception

Training:  $w_i = w_i + \Delta w_i$ ,  $\Delta w_i = \eta(t - O)x_i$   
 (earlier days)

↳ guarantees convergence if data is linearly separable.

what if data is not linearly separable?

→ Delta rule

→ Gradient descent

→ Dot product of weighted vector & input vector

→ Unthresholded:  $a(\vec{x}) = \vec{w} \cdot \vec{x}$

(perception)

→ cost function  $E(\vec{w}) = \frac{1}{2} \sum (t_d - O_d)^2$

(error function)

→ gradient of cost function  $\nabla E = \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$

$$w_i = w_i + \Delta w_i$$

$$\Delta w_i = \eta \frac{\partial E}{\partial w_i}$$

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \sum_{d \in D} \frac{\partial (t_d - O_d)^2}{\partial w_i}$$

$$\begin{aligned} \frac{\partial (t_d - O_d)^2}{\partial w_i} &= 2(t_d - O_d) \cdot \frac{\partial (t_d - O_d)}{\partial w_i} \\ &= 2(t_d - O_d) \cdot -\frac{\partial O_d}{\partial w_i} \\ &= -2(t_d - O_d) \cdot \frac{\partial (\vec{w} \cdot \vec{x}_d)}{\partial w_i} \\ &= -2(t_d - O_d) \cdot \frac{\partial (\vec{w} \cdot \vec{x}_d)}{\partial w_i} \end{aligned}$$

$$O_d = \vec{w} \cdot \vec{x}_d$$

$$\vec{w} \cdot \vec{x}_d = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4$$

$$\frac{\partial w_i x_i}{\partial w_i} = x_i$$

$$\Rightarrow -\frac{\partial d}{\partial w_i} = x_i$$

if  $d$  is 1 then  $\vec{w} \cdot \vec{x}_d = 1$

$$\Delta w_i = w_i + \Delta w_i$$

$$-\frac{n \Delta E}{\Delta w_i} = n \left[ \frac{1}{n} \sum_{i=1}^n (t_d - o_d) x_i \right]$$

$\Sigma$  is also these

similar to the  
the  $\Delta w_i$  of earlier  
days  
(i.e. similar for  
linearly separable)

$$\Delta w_i = \frac{n (t_d - o_d) x_i}{n}$$

But here output is different  
(previously only two outputs +1, -1)

$$\frac{(t_d - o_d) x_i}{n} \rightarrow w'_0 \rightarrow w''_0 \rightarrow w'''_0$$

stop when  
the updates  
are not  
significant  
or when  
resources  
ends.

Pick one research paper  
published in 2024-2025

Learn  
Research

from → Conference

CS

Journal

Nature  
Science  
Cell

- ACL★  
- NeurIPS★  
- NATACL★  
- AACL★  
- AAAI

- EMNLP★  
- CVPR  
- ICLR

BUT no  
quality  
check

Preprints

arXiv

CORE

funding  
agency for  
CS research  
papers

Ratings

A+

B

C

other

Presentation

(Explain  
the paper)

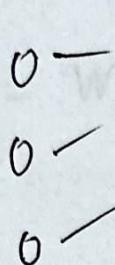
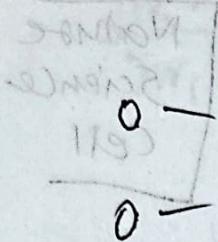
Before  
mid-sem

Demo → After  
Midsem

Research paper

- Title
- Abstract
- Intro.
- Our work
- Experiments
- Related work
- Conclusion
- Limitations

instead of taking one neuron  
take multiple neurons.



O —  $\times$   $n$  no

Similar  
to,  $y = a_1 n + b_1$

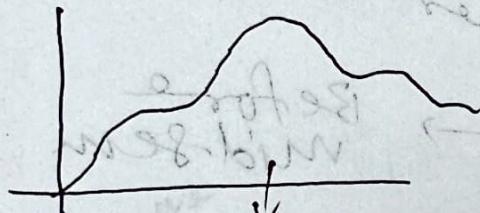
\* DATA  $a_2 n + b_2$

\* DATA  $a_3 n + b_3$

DATA  $a_4 n + b_4$

DATA  $a_5 n + b_5$

multiple linear functions will still give  
linear function  $\Rightarrow$  that's a problem  
 $\therefore$  Complete non-linear function



output required  
(from data)

linear functions  
can't give  
this type of  
non-linear outputs

$\therefore$  Use non-linear  
functions.

$f_1$   
O —  
 $f_2$   
 $f_3$   
 $f_4$

O —

Layer 1      Layer 2

$f_1(I/P)$

$f_2(f_1(I/P))$

$f_3(I/P)$

$f_4(f_3(I/P))$

I/P  $\rightarrow$   $L_1$

$L_2$

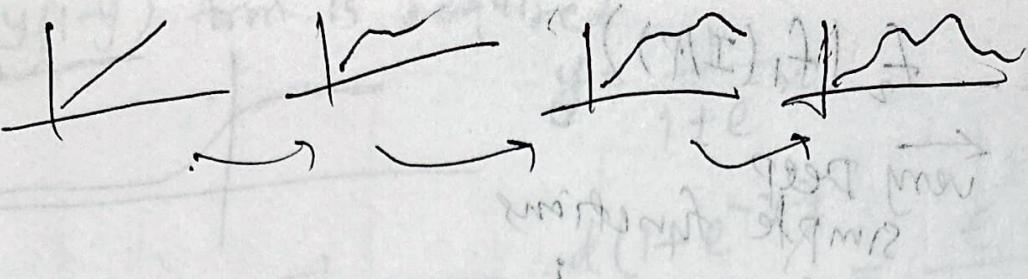
$L_3$

$f_1, f_2, f_3, f_4$  are  
non-linear functions

use multiple layers &  
neurons to get the  
output as close as the  
desired output.

$f_3 f_2 f_1 (I/P)$

$f_4(f_3 f_2 f_1 (I/P))$



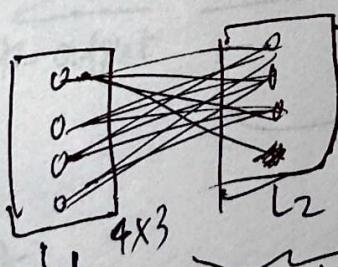
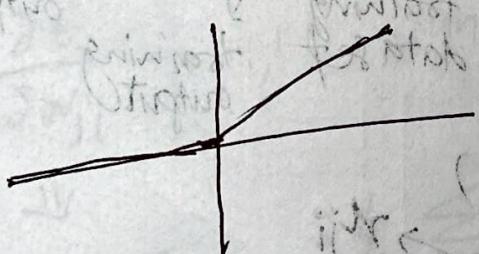
$\rightarrow$  non differentiable  
 $\therefore$  has a problem.  
 Used in earliest models  
 (Perception threshold)  
 Sgn ( $n$ )  
 Sigmoid function

3. Use sigmoid function

$$f(x) = \frac{1}{1 + e^{-x}}$$

$\rightarrow$  is differentiable  
 so we can use gradient descent, etc  
 $\Rightarrow$  Activation function

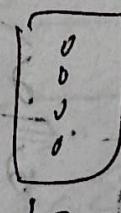
Or ReLU (Rectified



(Input layer)

multiple hidden layers

completely connected

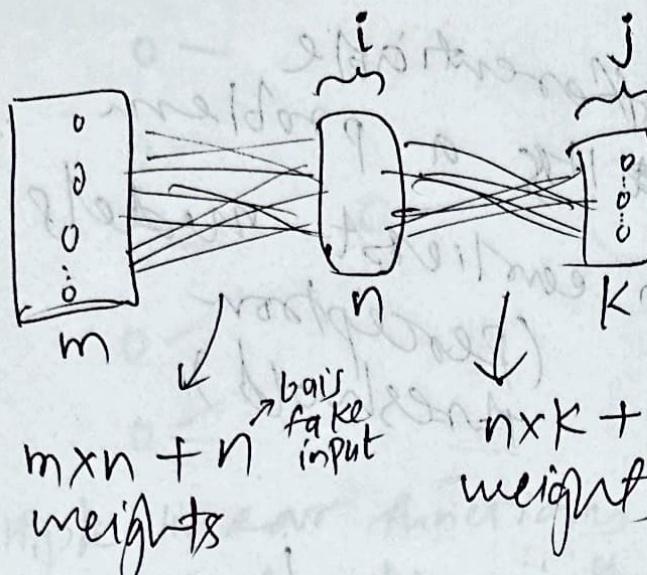


(Output layer)

(All neurons of one layer connected to all the neurons of next layer)

$$f_{\text{out}} \rightarrow f_{200} \rightarrow f_2(f_1(I/P))$$

very Deep simple functions



net input received at neuron  
 $\Rightarrow \text{net}_j$

weight  
 $w_{ji}$

output  
 $o_j$

$x_{ji}$   
the input  $i$  received

$$w_{ji} = w_{ji} + \Delta w_{ji}$$

$$-\eta \frac{\partial E}{\partial w_{ji}} \rightarrow \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

training data set  
model output  
training output

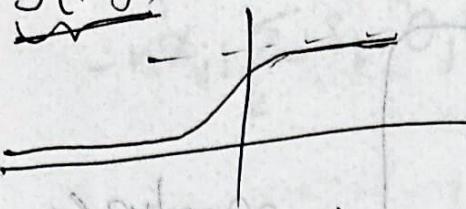
$$\frac{\partial E}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_{ji}}$$

$$\frac{\partial o_j}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{ji}}$$

$$o_j = \text{activation}(\text{net}_j)$$

desired  
 $\Delta w_{ji}$

$y(1-y)$  form is required



$$y = \frac{1}{1+e^{-x}}$$

$$\frac{dy}{dx} = \frac{-1}{(1+e^{-x})^2} \cdot \frac{d(1+e^{-x})}{dx} = \frac{e^{-x}}{(1+e^{-x})^2}$$

$$\therefore y(1-y) = y - y^2 = \frac{1+e^{-x}}{(1+e^{-x})^2} - \frac{1}{(1+e^{-x})}$$

$$\frac{\partial o_j}{\partial \text{net}_j} = o_j \cdot (1-o_j)$$

$$\text{net}_j = \sum_{\substack{\text{all} \\ \text{dimensions}}} x_m w_m \rightarrow \frac{\partial \text{net}_j}{\partial w_{ji}} = x_j$$

$$\frac{1}{2} \sum_{k=\text{output}} \sum_{d \in D} (t_{dk} - o_{dk})^2$$

→ go through all neurons  
and the training dataset

$$\frac{\partial E}{\partial w_{ji}}$$

$$\sum_{j=\text{output}} \frac{\partial \frac{1}{2} \sum_{d \in D} (t_{dj} - o_{dj})^2}{\partial o_j}$$

$$\frac{1}{2} \sum_d 2(t_{dj} - o_{dj}) \cdot \frac{\partial (t_{dj} - o_{dj})}{\partial o_j}$$

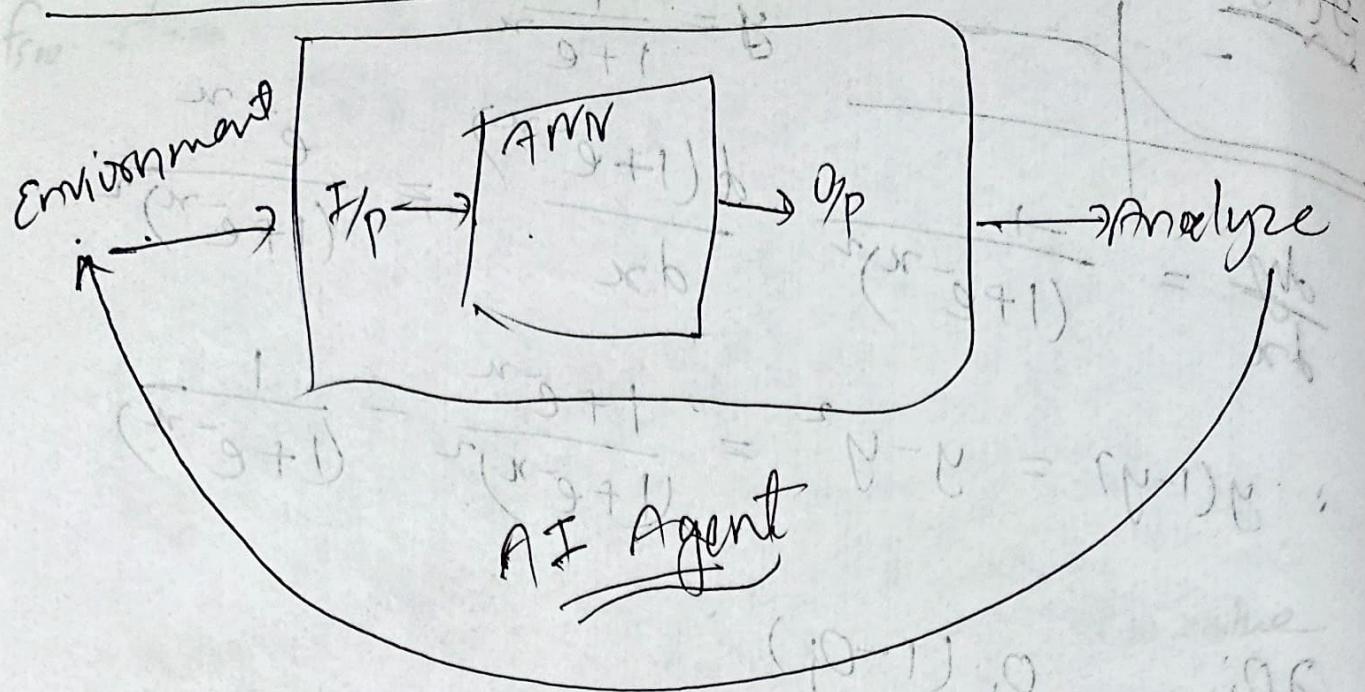
-ve of  
derivative of error  
w.r.t. net  $j$

$$\Rightarrow \frac{\partial E}{\partial \text{net}_j}$$

Backpropagation

$$\Delta w_{ji} = -\eta \cdot \delta_j \cdot x_{ji}$$

For output layer  
learning rate



① Frozen LLM Agent

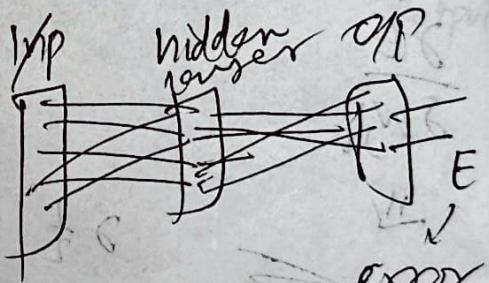
② RLLM Agent

$$\Delta w_{ji} = n \delta_j x_{ji} \quad \left\{ \text{for output layer} \right.$$

$$\delta_j = -\frac{\partial E}{\partial \text{net}_j}$$

$\Rightarrow$  ~~for output layer~~  $\Rightarrow$  ~~for hidden layers~~

$$\Rightarrow -n \frac{\partial E}{\partial w_{ji}}$$



$$\frac{\partial E}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{ji}} \cdot \delta_j$$

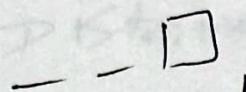
$$\sum_{k \in \text{Downstream}(j)} \frac{\partial E}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_k}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j}$$

$$o_j(1-o_j)$$

$\sum_k \alpha_{kj} \sum_k \delta_k \cdot o_j (1 - o_j) w_{kj} \rightarrow$  For hidden layers

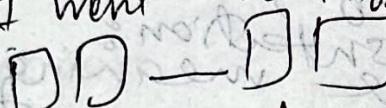
## # Word To Vec

### Fake Task



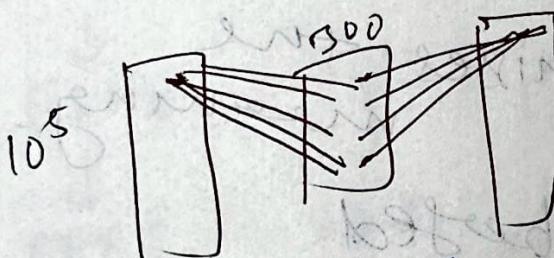
Predict the neighbour words.

I went school today



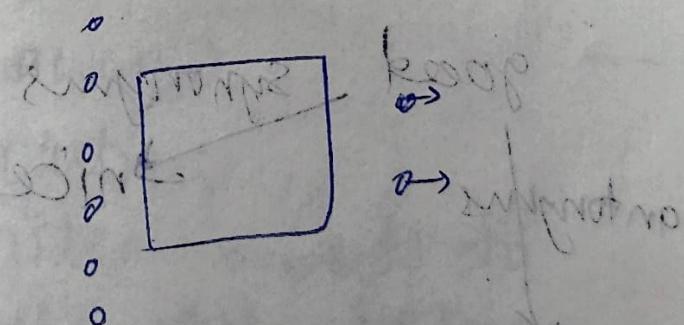
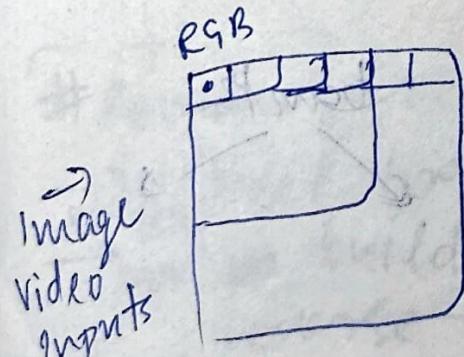
Predict the missing word.

→ static representation.



# How to give Input to ANN?

Data - Nos, text, Images, audio, ~~time series, graphs, etc.~~



Waves → waveforms time domain  
freq. domain  
→ audio inputs

Graphs can be given as input as  
Adjacency lists or non matrix.

How to give text as inputs.

earliest way  $\rightarrow$  One hot representation

$d_1$	1	apple	0
$d_2$	0		0
:	0		0
:	0		1
$d_{10^5}$	0		0

Problems with one hot representation:

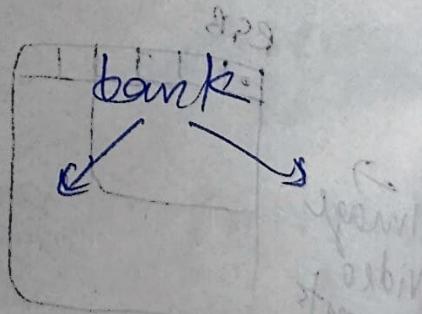
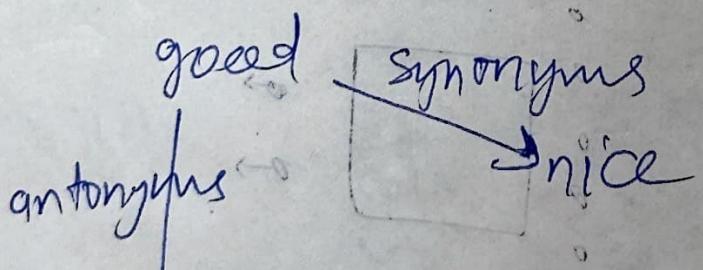
• It completely misses the meaning,  
New words comes up with time

2nd way  $\rightarrow$  ASCII representation  
represent words at character level.

$\rightarrow$  All words don't have same

length, it also misses the meaning.

3rd way  $\rightarrow$  Rule based



# Limitations of Symbolic representation

- Related words
- New words
- Change of meaning
- Multiple meaning.

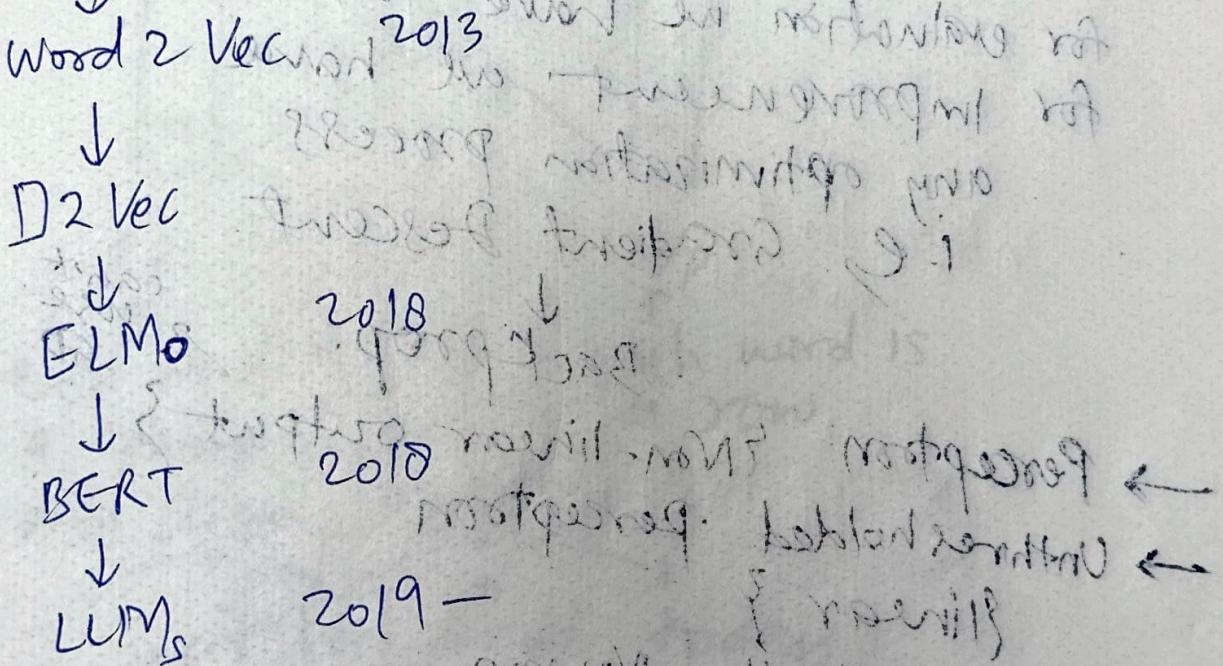
→ Distributed representation,

Distributional Semantics

Words are represented

in the form of vectors.

Back prop.

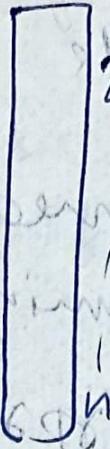
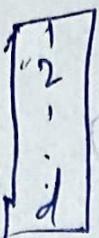
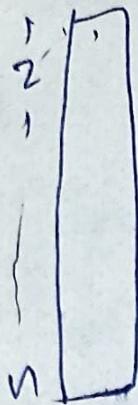


# Remember: (Takeaways)

→ Model based solutions  
we need  
→ we build statistical models to  
reverse engineer the system

→ Linear regression  
make a guess → evaluate → improve

$V_{dxn}$



$V_C$

$U_0^T$

$$E(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_j \log(P(w_{t+j} | w_t))$$

for evaluation we have loss function

for improvement we have

any optimization process

i.e. Gradient Descent

Backprop.

can't use  
GD

→ Perception {Non-linear output}

→ Unthresholded perception

{linear}

→ Differentiable Neuron

Activation function → Sigmoid

→ How to represent input

Distributional  
semantics  
distributed representation

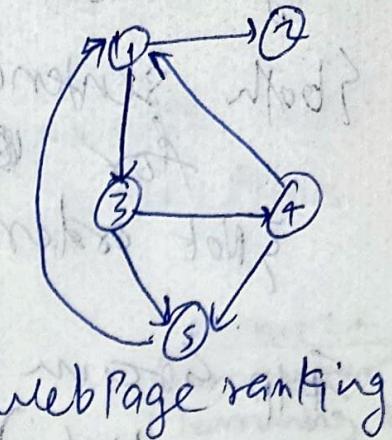
10% weightage syllabus till last Friday  
open notes

Quiz on Thursday → during class time.

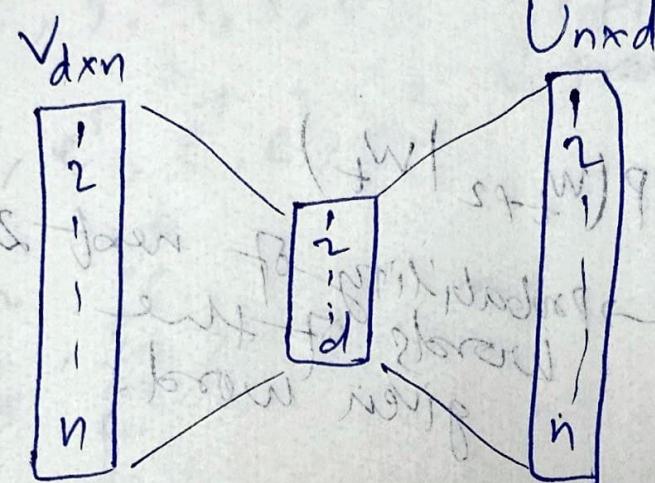
→ Distribution of probability of how often a word is used.

apple: fruit, red, sweet, etc

: company, steve jobs, etc



# Word to Vec :-



I/P  
each word is  
a column

O/P  
each word is  
a row

skip gram model

I went To school today  
predict the 2 words back & front  
of the given word?

CBOW {cont. bag. of words} model.

I went To school today  
guess the missing word?

Today school  $\boxed{\text{to}}$  went I

Both sentences are same  
for  $\{\cdot\}$  (BOW)

{Not order dependent}

Skip Gram Model:-

$$P(W_{t+1} | W_t)$$

next word      given word

Probability

$$P(W_{t+1} | W_t) \cdot P(W_{t+2} | W_t)$$

↳ probability of next 2 words of the given word.

Independent events

$$\prod_{j=-2}^2 P(W_{t+j} | W_t)$$

~~if short words~~      ↳ probability of surrounding words at distance 2

$$\sum_{j=-2}^2 \log(P(W_{t+j} | W_t))$$

↳ more value means how good the model is

↳ take -ve of this gives the cost function  
show how bad the model is

for All tokens (e.g. words)

$$\Rightarrow - \sum_{t=1}^T \sum_{j=-2}^2 \log(P(w_{t+j}|w_t)) \rightarrow \text{cost function}$$

(Window size of 2.  
{Hyper Parameters})

softmax:

$$\begin{array}{c} 2, 5, 1, 7, 0 \\ \downarrow \\ e^2, e^5, e^1, e^7, e^0 \end{array} \quad \begin{array}{l} \text{Softmax} \\ \text{Transformation} \end{array}$$
$$\frac{e^2}{\sum 5}, \frac{e^5}{\sum 5}, \dots = 1$$

Resource  
Parameter  
learning in  
Word2Vec  
christopher  
mcmurry  
stanford  
lectures