

Day 1 (13/01/26)

- Mid : 15%

- End : 25%

Continuous Assessment : 60%

- In ML, experience \equiv dataset

Machine + Dataset \rightarrow Algorithm \Rightarrow Output

\rightarrow ground truth

- ML Paradigms \longrightarrow i) Supervised : Labelled Training dataset

ii) Unsupervised

iii) Reinforcement : learning from trial and errors

\rightarrow Includes association

rules mining

\rightarrow representation learning

(embedding)

\rightarrow Identifies relationship

Q. In unsupervised learning,

Since we don't label

the data, in that case

when there is no output
how is the class labelled ?

(The model transforms
organizes the input into
new representation)

Q. What includes in

unsupervised ? (algo?)

Includes algorithms where

i) no labels are provided

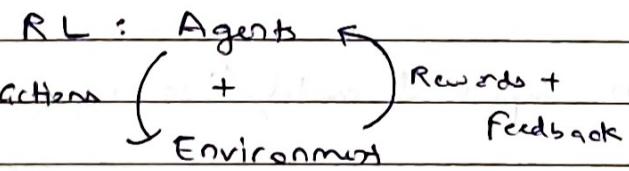
ii) the goal is to discover structure

1 \rightarrow train using supervised

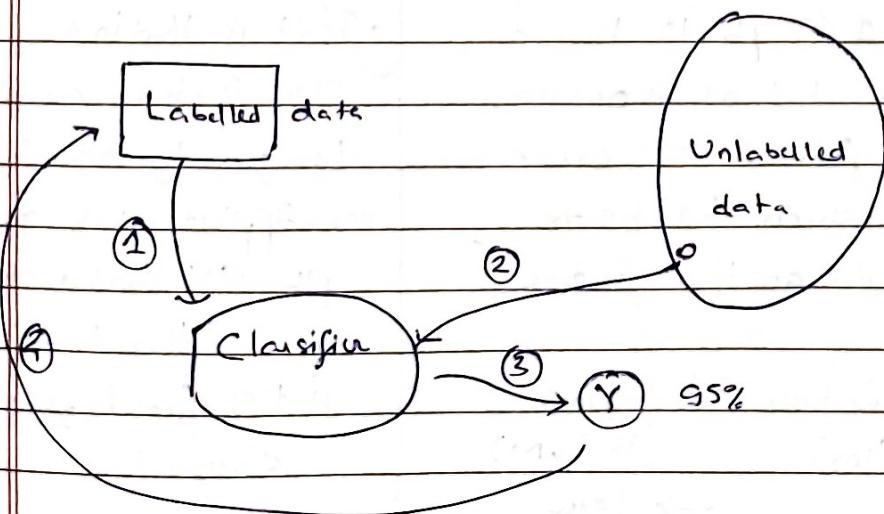
2 : predict / classify

3 : if correct

4 : add to dataset



- Semi-supervised learning



- In all ML task, a 'good representation of data' very imp? 1st
 ↳ even in DL] ↑ Step!
 as in input
- feature engineering
 (not req in DL, why?) - automatic feature learning
 one type: vector space model \Rightarrow includes feature vectors of representation

- (i) feature selection: out of all features (feature set), choosing some affects perfor- ←
 manc, computation
- (ii) dimension reduction: no removal of features, but projection in lower dimens
 $\rightarrow (3D) \rightarrow 2D$

- DL helps in feature engineering (p) Note: no elimination → done by the model itself (mathematically) of features; but automating it }

Day 2 (15/01/25)

- No of classes pre-defined in supervised model / learning
- Probabilistic model - Bayesian classifier
 - Note: when we had only one feature and we had 1/0 pb it was a good feature but as soon as, pb dropped from 1 \Rightarrow error introduced which led us to think / search for another feature
 - So it that in a classification model, we generally have softmax applied which gives pb and then we classify \rightarrow but then how is it diff from bayesian classifier?
- Increasing features \Rightarrow Kernel (upsampling) (in SVM) \rightarrow later

- Upscaling and downscaling always done on datasets

~~How do we do feature engineering?~~
 ↳ normalisation
 encoding etc

- Classical (Traditional) ML Methods for Classification

a) Bayesian classifier - (Suitable for 2 classes)

(finding plane for separation) $P(C_i | \bar{x}_j)$

Class label new sample vector likelihood prior

↑ ↑ ↳ (it is a feature)

NOTE: Bayes rule : $P(C_i | \bar{x}) = \frac{P(\bar{x} | C_i) P(C_i)}{P(\bar{x})}$

posterior Evidence

- given the sample, what's the class label? ($P(C_i | \bar{x})$)
- $P(\bar{x} | C_i)$ - prob of picking \bar{x} given the class C_i
- $P(C_i)$ - prob of picking the class C_i
- $P(\bar{x})$ - prob of picking \bar{x} in the universe

NOTE: if $P(C_i | \bar{x}) > P(C_j | \bar{x})$

$\Rightarrow \bar{x}$ is C_i

$$P(\bar{x} | C_i) * P(C_i) > P(\bar{x} | C_j) * P(C_j)$$

↳ This implies no significance of denominator

NOTE: if prior probability not known, then it is considered same! (That is distribution equal)

$$\Rightarrow P(\bar{x} | C_i) = P(x_1, x_2, x_3, \dots, x_k | C_i)$$

$$(Pb1) \quad \quad \quad = P(x_1 | x_2, x_3, \dots, x_k, C_i) * P(x_2 | x_3, \dots, x_k, C_i) \dots \dots P(x_k | C_i)$$

• (expensive task) - computationally

(A.2) if dataset set sparse, then (Condition)

May turn to 0 //



What's the way out?

(i) features are conditionally independent - considered as
↳ bag of words

$$\begin{aligned} \text{(Naive Bayes)} & \quad \text{(may look seeming)} \Rightarrow P(u_1 | u_2, u_3, \dots, u_k, C_j) \\ \text{Classifier} & = P(u_1 | C_j) \end{aligned}$$

(ii) (after assumption)

Note: none of $P(u_i | C_j)$ shouldn't be zero
after applying (i)

\Rightarrow so we apply smoothing

$$P(u_i | C_j) = 0 + \text{small}$$

how to estimate $|C_j| + \text{large}$

this prob? probability

among all elements in C_j ,

how many of them having u_i

with Gaussian / feature)

↳ Multinomial distribution

here the assumption is: it is already there in the

dataset / dataset has gaussian distribution

\rightarrow (representative finding)

(b) K Nearest Neighbors

&

Centroid Based Classifier.

\Rightarrow Lazy classifier / Instance Based classifier

x_1	c_1
x_2	c_2
:	:
:	:
x_n	c_n

- find distance with all elements in dataset
 - pick K least one : closest one
- \Rightarrow This is known as inference time : where we identify the class label for a new input sample!
- \Rightarrow dependent on size of dataset for (KNN) - The inference time
- K : no. of feature $\leq O(nk)$ - for kNN
 n : no. of inputs in dataset
- $O(k)$ - for Bayesian

(b) Centroid Based Classifier

- We have representatives for each class
- take mean of all feature vectors = centroids

for a new input : find the distance from that from each class's representation

Note: Training over here = finding centroid for each class

Inference time less!

(c) Decision Tree - (Rule based approach)

(Random Forest)

- Objective : to find rules for classification (ie build the tree)
- Leaf nodes - class labels
- Internal nodes - the rules / condition
- Many trees possible = (but need to find best one)



Day 3 (19/01/26)

(C) Decision tree contd

Q if there are 2 features, which one should be root?

\Rightarrow either of them can be taken and form tree

\Rightarrow but one is ~~very~~ good: to have resultant dataset as homogenous !!

achieved / done by information theory

(i.e. the uncertainty should be low)

$$X = \{n_1, n_2, n_3, \dots, n_n\}$$

$$\begin{aligned} \text{entropy } H(X) &= \sum_{x \in X} P(n) \cdot \log\left(\frac{1}{P(n)}\right) = E\left(\log\left(\frac{1}{P(n)}\right)\right) \\ &= -\sum_{x \in X} P(n) \log P(n) \end{aligned}$$

expectation information content (or uncertainty)

• Data: 1, 2, 1, 5, 4, 2

• given an event, how much info we ~~know~~ / don't know

• Mean: $\frac{\sum \text{Data}}{|\text{Data}|}$

• Note: If $P(n)$ high, $\frac{1}{P(n)}$ low

$$= \frac{1+2+1+5+4+2}{6}$$

\Rightarrow less uncertain

$$= 1 \cdot \left(\frac{1}{6}\right) + 2 \cdot \left(\frac{1}{6}\right) + 1 \cdot \left(\frac{1}{6}\right) + \dots + 2 \cdot \left(\frac{1}{6}\right)$$

1 is value

$$= 1 \cdot \left(\frac{2}{6}\right) + 2 \cdot \left(\frac{2}{6}\right) + 5 \cdot \left(\frac{1}{6}\right) + 4 \cdot \left(\frac{1}{6}\right)$$

$\frac{1}{6}$ is its probability

$$\Downarrow = \sum n \cdot P(n) = \text{expectation}$$

This gives: $X = \{1, 2, 4, 5\}$

$$H(n) = n \cdot \frac{1}{n} \log(n) = \log(n)$$

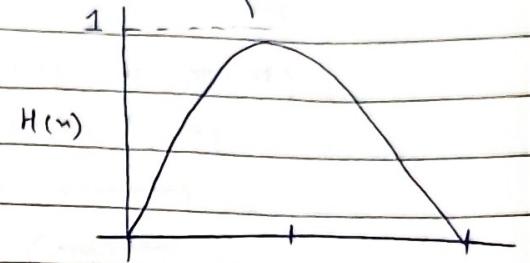
(If distribution is equally likely, highest entropy)

at middle (c)

$$X = \{H, T\}$$

p : probability of get head
 $1-p$: " " tail

$$H(n) = -p \log p - (1-p) \log(1-p)$$



- Most fundamental Decision tree is ID3!

\Rightarrow Information gain: difference in the entropy

(entropy is 0)!

as certain outcome

comes

\Rightarrow if uncertainty present \rightarrow keep on dividing till homogeneous

\Rightarrow if leaf nodes has only 1 element that means the tree construction is weak (i.e. rules weak)

↓

Random forest } all applied on ID3 (so extension)
Ensembling method

(d) Support Vector Machine (SVM)

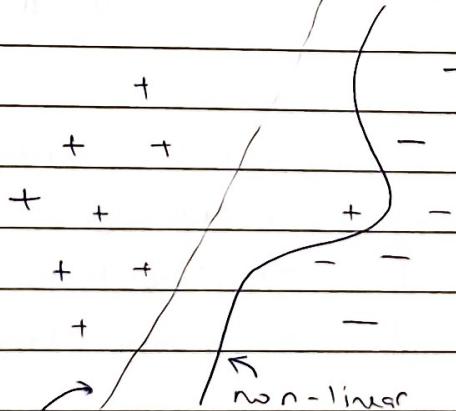
- discriminant based method
- finds a separating hyper-plane

Linear SVM

$$g(u) = w^T u + w_0 = 0$$

if $g(u) > 0$: positive class

if $g(u) < 0$: negative "



linear SVM

(LSVM)

SVM (complex)

\Rightarrow but we will

solve it by LSVM

(project in \Leftrightarrow not possible in \Leftrightarrow req data modification (or selection (or separation!))

high or

low dimension

cation (or se-

dimension!)

(in terms of Separating parable)

(Note: we find that hyperplane ($g(\mathbf{w})$) which results in Maximum Separating margin)

Q How to find \mathbf{w} of hyperplane

- take two points on $g(\mathbf{w}) = 0$

$$g(\mathbf{x}_1) = g(\mathbf{x}_2) = 0$$

$$\Rightarrow \mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0$$

$$\Rightarrow \mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$$



This indicates \mathbf{w} is orthogonal to $g(\mathbf{w})=0$!

- $g(\mathbf{w})$ better than $g(\mathbf{w})'$ parallel

• Note: $g(\mathbf{w})^+$ is the hyperplane to $g(\mathbf{w})$ passing through nearest positive class similarly $g(\mathbf{w})^-$

- Note: we have infinite ($g(\mathbf{w})^+ - g(\mathbf{w}) - g(\mathbf{w})^-$) planes

but ideally we take the middle one!

NOTE: $\gamma = \frac{|g(\mathbf{w})|}{\|\mathbf{w}\|}$: distance of point from hyperplane

Day 4 (20/01/26)

(d) SVM (Linear)

→ Video lesson will be uploaded for implementation

→ binary classifier

Q What is the magnitude of separating margin?

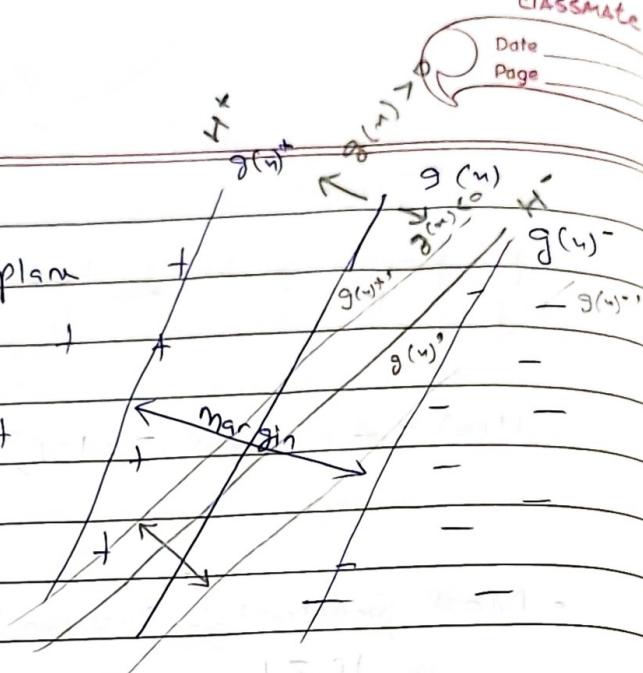
$$H^+ : g(\mathbf{x}) \geq 1$$

$$H^- : g(\mathbf{x}) \leq -1$$

$$\Rightarrow y_i (\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 \quad \forall \mathbf{x}_i$$

$$y_i = 1 \text{ if } H^+$$

$$= -1 \text{ if } H^-$$



| Margin

$$\begin{aligned} &= \text{distance of origin from } H^+ \\ &\quad - \text{distance of origin from } H^- \\ &= \frac{2}{\|w\|} \end{aligned}$$

$\left\{ \begin{array}{l} \text{What about points b/w} \\ +1 \text{ to } 0 \text{ and } 0 \text{ to } -1? \end{array} \right\}$

Considered inside margin

penalized later

NOTE: margin boundaries

$$\delta(n) = 1 \text{ or } g(n) = -1$$

• Maximizing $\frac{2}{\|w\|} \equiv \text{minimising } \frac{\|w\|}{2}$

$$\equiv \frac{1}{2} \|w^T w\| \text{ ? Objective function}$$

Q What if data is not linearly separable?

→ 2 approaches

→ Approach 1: Slack

variable

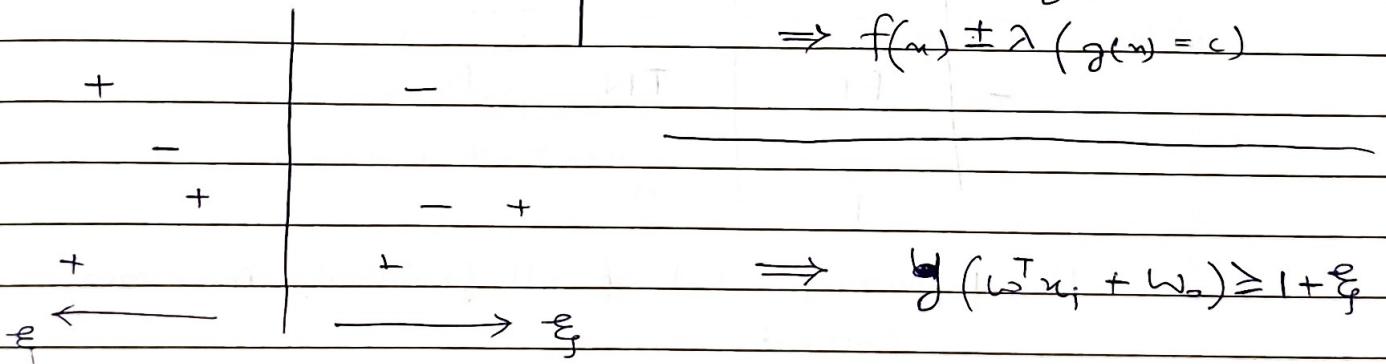
$$\text{Subject to: } y_i (w^T x_i + w_0) \geq 1 \forall x_i;$$

⇒ using Lagrange multiplier

$$f(w)$$

$$\text{s.t. } g(w) = c$$

$$\Rightarrow f(w) \pm \lambda (g(w) = c)$$



→ Not linearly separable

→ So pull both the sides by some

$$\xi$$

$$\Rightarrow y (w^T x_i + w_0) \geq 1 + \xi$$

→ Approach 2: Kernel (How to formulate: look at vid?)

• not separable in low dim → go to higher dim (but costly)

• So Kernel idea is: don't transform to high dim

but solve it like how it is solved there!

(points that touch decision and margin boundary or fall between them)
(check)

Q What are support vectors?

- Vectors lying on separating hyperplane
- $\lambda_i \neq 0$ for non-support vectors

~~Only 3 points?~~
3 types of support vectors
(on margin, inside margin, misclassified)!

Q How to evaluate our model? (# Evaluation)

(i) Dataset types:

- Training Data set \Rightarrow Building
- Testing " \Rightarrow Evaluation
- Development data : used for fine-tuning parameters
~~may or may not req~~

		Predicted Label		
		+	-	
Actual Label	+	TP	FN	P
	-	FP	TN	N
		P	n	

(iii) Accuracy : out of all predicted, how many correct

$$\frac{TP + TN}{TP + FN + FP + TN}$$

(iv) Precision : among all positive predicted, how many

- Binary \rightarrow actually correct in nature

\$ o keep one var as

\uparrow Positive } *

Similarly for Recall

$$\frac{TP}{TP + FP}$$

(iii) Recall : out of actual positive, how many predicted correct?

$$\frac{TP}{TP + FN}$$

(iv) F measure: harmonic mean of Precision and recall

Clustering

~~i~~ K-mean

~~ii~~ hierarchical clustering

(iii) density based "

Regression 2 Videos
Clustering lessons

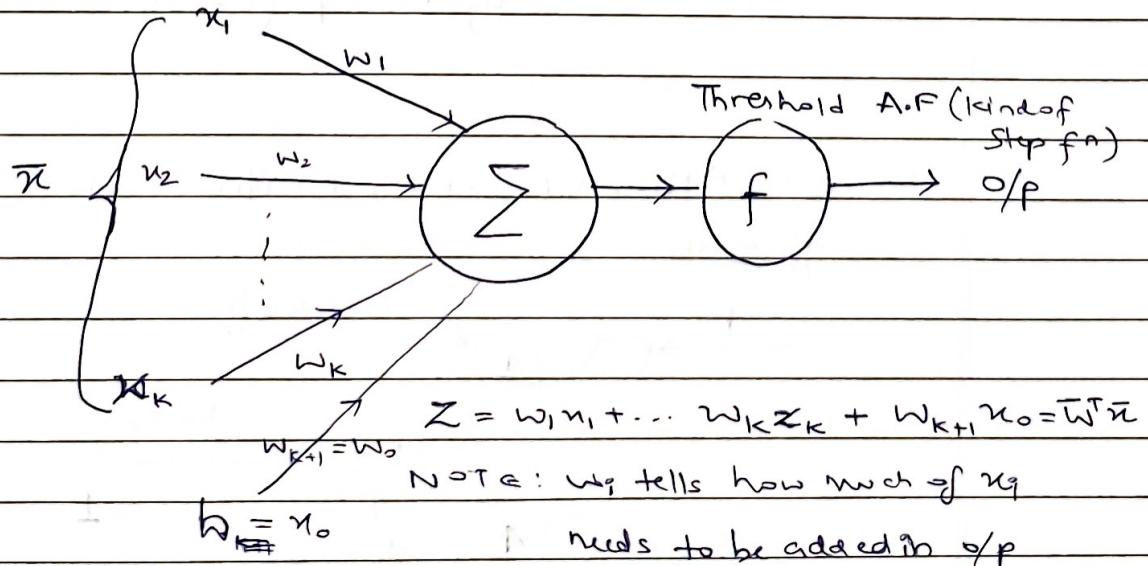
Day 5 (21/02/26)

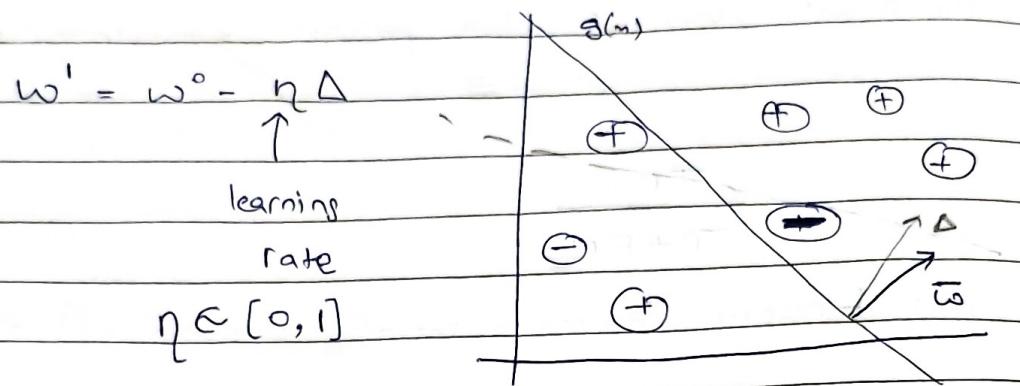
Artificial Neural Network

(i) Perceptron

→ Solves the linearly Separable problem

→ Start with some $g(\mathbf{x})$ which is not separating completely, at first but then we move it to more in certain dirⁿ
So it separates the points





$$\chi = \bar{w}^T \bar{u}$$

M: missclassified examples

NOTE: loss $l_i = \begin{cases} -y_i(w \cdot u_i) & \text{if misclassified} \\ 0 & \text{otherwise} \end{cases}$

$$E = \sum l_i = -\sum y_i(w \cdot u_i)$$

objective function: $E = -\sum_{u_i \in M} (\bar{w}^T \bar{u}_i, y_i)$ } why is the error term weird?

NOTE: if correctly classified,

$$\text{minimizing} \Rightarrow \frac{\partial E}{\partial w} = -\sum_{u_i \in M} y_i \bar{u}_i \quad \begin{array}{l} \text{for one} \\ \text{misclassified point} \end{array} \quad y_i(w \cdot u_i) > 0$$

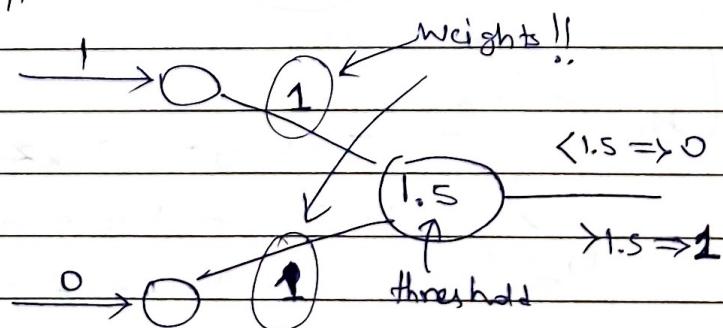
$$\Rightarrow w^{t+1} = w^t + \eta y_i \bar{u}_i \quad \begin{array}{l} \text{for one} \\ \text{misclassified point} \end{array} \quad y_i(w \cdot u_i) < 0$$

$$\Rightarrow w^{t+1} = w^t + \eta \sum_{u_i \in M} y_i \bar{u}_i \quad \begin{array}{l} \text{only for misclassified samples} \end{array}$$

Paper: Neural Networks for pattern recognition } Single layer
by Bishop (Ch 2) } N.N

NOTE: AND, OR, NOR, NAND : solved by perceptron

AND

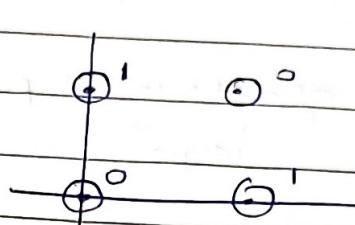


NAND $\Rightarrow [-1, -1] = w$, threshold = -1.5

OR : $[1, 1] = w$, threshold = 0.5

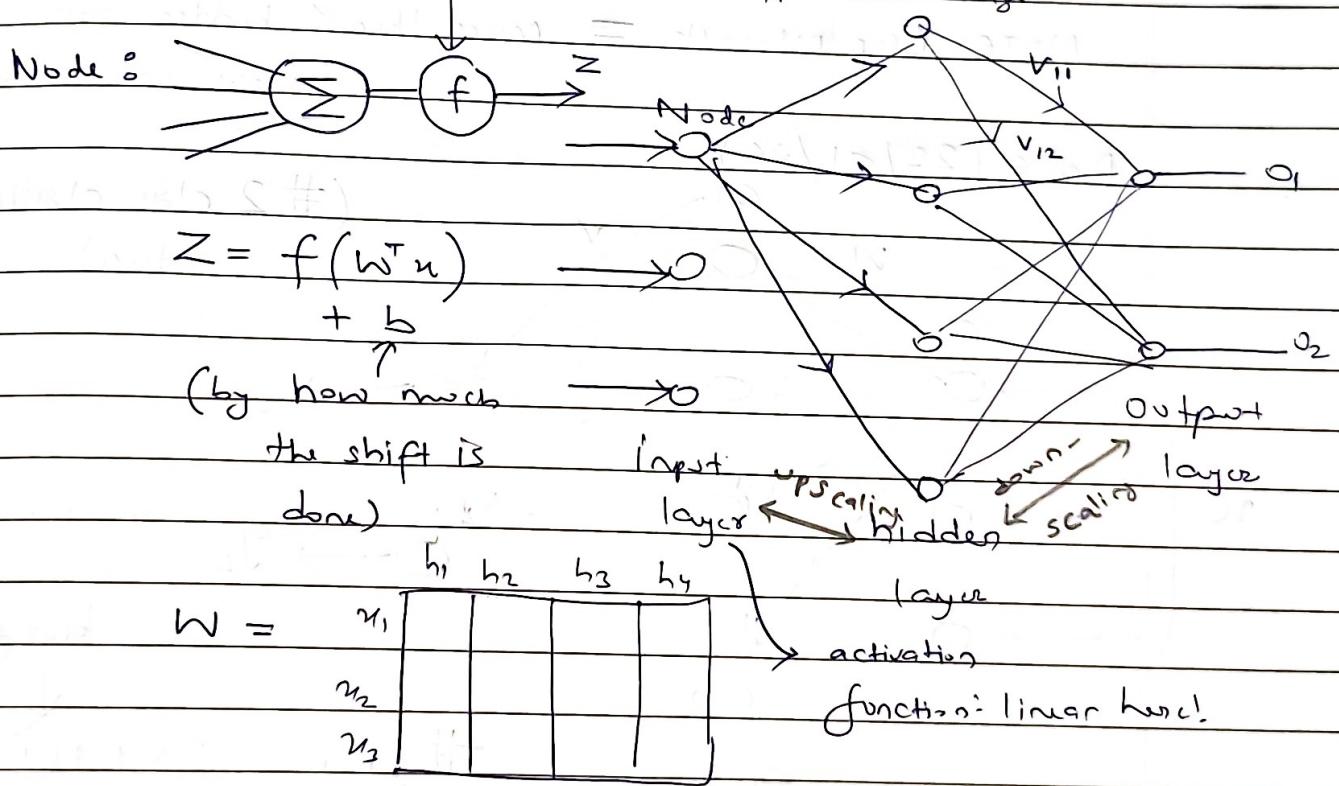
XOR : not linearly separable $\approx \begin{matrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{matrix}$

\Rightarrow add more layers



(ii) Multilayer Perceptron / Feedforward NN / ANN

\rightarrow constraints the o/p to some range



(Deep network

= multiple

hidden layers)

$$= f(\omega_y^T \bar{u})$$

$\downarrow y \in \text{no. of } \bar{u}$

feedforward

$$\Rightarrow \bar{u} = f(\bar{\omega}^T \bar{u})$$

(Output of hidden layer!)

$$\text{Similarly. } O = f_o(\bar{v}^T \bar{h}_o)$$

$$O_1 = f(v_{11}h_1 + v_{21}h_2 + v_{31}h_3 + v_{41}h_4)$$

$$\Rightarrow O = f_o(\bar{v}^T f_o(\bar{\omega}^T \bar{u}))$$

for classification pb,

- no. of nodes in input layer = no. of features in a feature vector

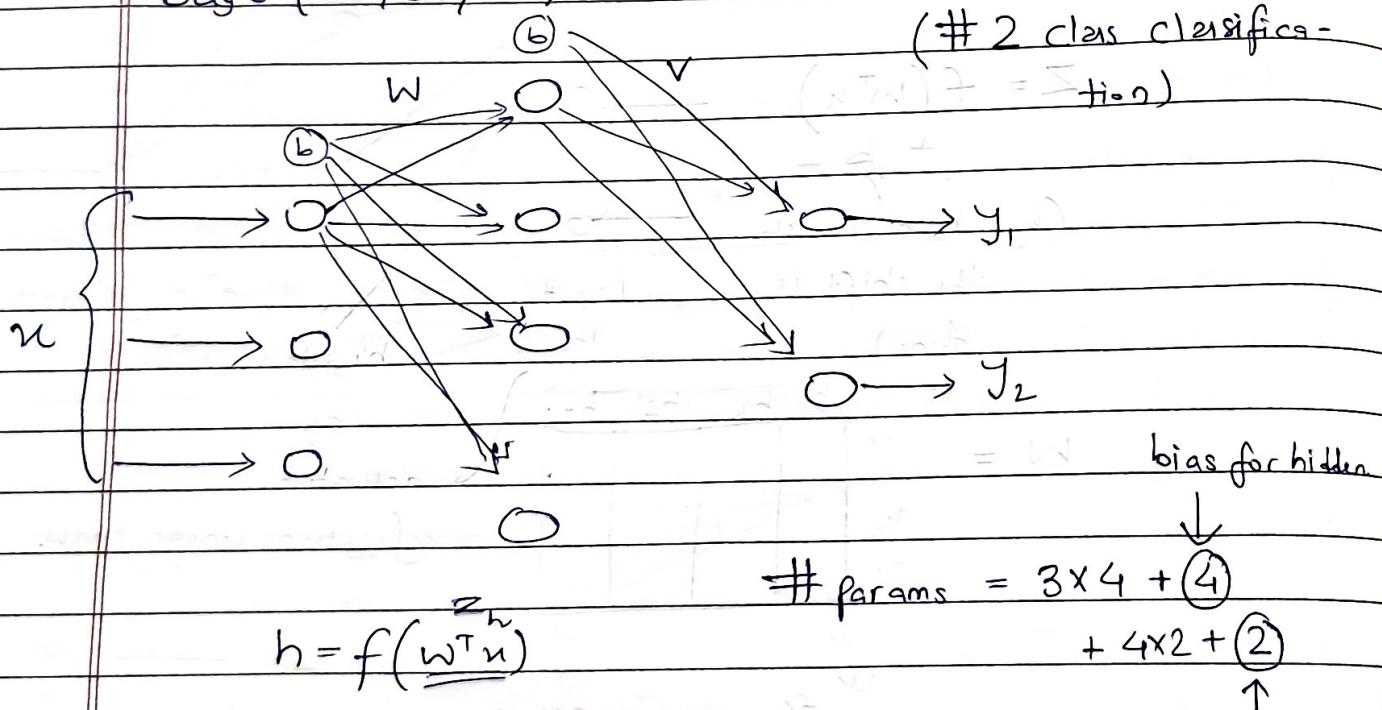
- no. of nodes in output layer = no. of classes

\downarrow
pb assigned to each class

for regression \rightarrow only 1 node in o/p layer or

NOTE: Deep network = more than 1 hidden layer

Day 6 (22/01/26)



Q What is w ? \rightarrow unknown

\Rightarrow Start with unknown values

$$\bar{y} = f(\underline{V^T b})$$

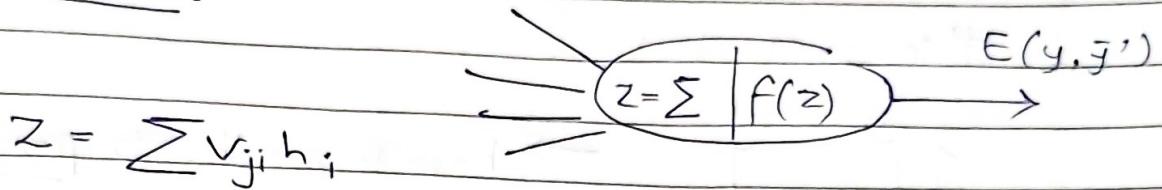
NOTE: Output of nn can be non-discrete as well

expected is ground truth (\bar{y})

observed is o/p of nn (\bar{y}')

$$E = e(y, \bar{y}')$$

Backpropagation



$$\Rightarrow z_1 = v_{11} \times h_1 + v_{21} h_2 + v_{31} h_3 + v_{41} h_4$$

$$y = f(z)$$

E = error function

$$\Rightarrow \frac{\partial E}{\partial v} = 0 = \Delta v$$

$$\left. \begin{aligned} \frac{\partial E}{\partial y} &= 0 \\ \frac{\partial E}{\partial z} &= 0 \\ \frac{\partial E}{\partial v} &= 0 \end{aligned} \right\} \text{Do this for all parameters } (v_{ji})$$

$$\frac{\partial z_i}{\partial v_{ji}} \cdot \frac{\partial y_i}{\partial z_i} \cdot \frac{\partial E}{\partial y_i} = 0$$

$$\Delta v_i = \frac{\partial z_i}{\partial v_{ji}} \cdot \frac{\partial y_i}{\partial z_i} \cdot \frac{\partial E}{\partial y_i} = 0$$

$$v^i = v^o + n \Delta v$$

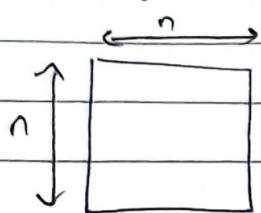
$$\Rightarrow \frac{\partial E}{\partial w} = 0$$

$$\frac{\partial z_h}{\partial w} \cdot \frac{\partial h}{\partial z_h} \cdot \frac{\partial z_o}{\partial h} \cdot \frac{\partial y}{\partial z_o} \cdot \frac{\partial E}{\partial y} = 0$$

↑ output layer's output

NOTE: if # classes = 2 we needn't req 2 nodes in o/p layer \rightarrow one is enough with threshold

Image Classification



\Rightarrow how to give this to prev ANN
Architecture?

- Convert 1024×1024 into a vector

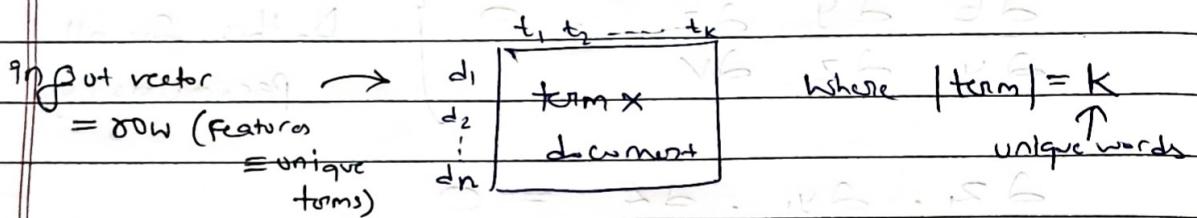
- Concatenate each rows back to back:

back: $r_1 \ r_2 \ r_3 \dots$
constitute an input vector
 $\# \text{features} = n^2$

Text Processing / Classification

What features to consider?

- (i) Unique words in corpus - so build a matrix like



A cell tells if t_j term present in d_i document

Instead we can have, \leftarrow a single word

- (ii) We can have sort of frequency stored also instead of just binary counting

bigram: 2 words at a time

trigram: 3 words at a time

NOTE: Representation imp for a good model

? prediction! \Rightarrow to be able to use in ANN
 as in trigram?
 bigram } etc.

Error functions

(i) Supervised Problem

- MSE - req continuous, differentiable, monotonic function

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\frac{\partial E}{\partial y} = \frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$

Joint probability

(ii) Cross entropy

→ uncertainty in a r.v.

$$H(x) = \sum_{x \in X} p(x) \cdot \log\left(\frac{1}{p(x)}\right)$$

to expectation

$$\text{of information content} = - \sum_{x \in X} p(x) \cdot \boxed{\log(p(x))}$$

$$P(x,y) = P(x|y) P(y) *$$

$$= P(y|x) P(x)$$

if independent: $P(x) P(y)$

Information content

↓ Amt of uncertainty

→ uncertainty of two variables together of the event

Joint Probability: Probability of both events happening

$$H(x,y) = - \sum_{x \in X} \sum_{y \in Y} P(x,y) \cdot \log(P(x,y))$$

if x, y independent $\Rightarrow H(x,y) \uparrow$
else $\Rightarrow H(x,y) \downarrow$

Relative entropy: measures how much one is diff from other

$$H(x|y) = \sum_{y \in Y} P(y) \cdot H(x|y)$$

for unsupervised learning

(KL divergence)

Marginalisation
on condition y

$$= - \sum_{y \in Y} P(y) \cdot \sum_{x \in X} P(x|y) \log(P(x|y))$$

$$= - \sum_{y \in Y} \sum_{x \in X} P(x,y) \log(P(x|y))$$