

- ML is one of the many ways to achieve AI
  - Deep learning < ML is yet another way to achieve AI
  - GOAL OF ML
    - TO create an algo to teach machine based on prior knowledge.
  - Subset of AI
    - Teach the machine on basis of prior knowledge
    - Train a m/c to be able to make decisions
    - Train a m/c to generalise from prior data (goal is to predict for future)
    - Construct a model whose efficiency increases as we spend more time training it. ↳ incorrect term  
decision making accuracy

## ML FORMAL DEF<sup>N</sup>

Study of algos that allow computer programs to automatically improve through experience (Tom Mitchell) \*

## Role of statistics: inference from a sample

Role of CS: efficient algs to — 1) solve optimisat<sup>n</sup> problem  
2) representing & evaluating model for inference

\* optimise a performance criterion using example data or past experience.

## SEEMINGLY RELATED FIELDS TO ML

- Statistical Learning
- Pattern Recognition
- Data Mining
- Deep learning

CONTRIBUTE TO AI IN SOME WAY

COMMON OBJECTIVE: analyze data with either of 2 goals

- prediction
- description

STATISTICAL LEARNING: mostly about inferences & parameters  
(SL)  
estimations, while ML emphasises supervised, unsupervised, semi-superv., online, reinforcement learning

Statistical learn. is based on smaller dataset (few dimensions) compared to ML (high dimension)

Extent of assumptions are more stringent in statistical learn, while ML is not that assumpt<sup>n</sup> dependent

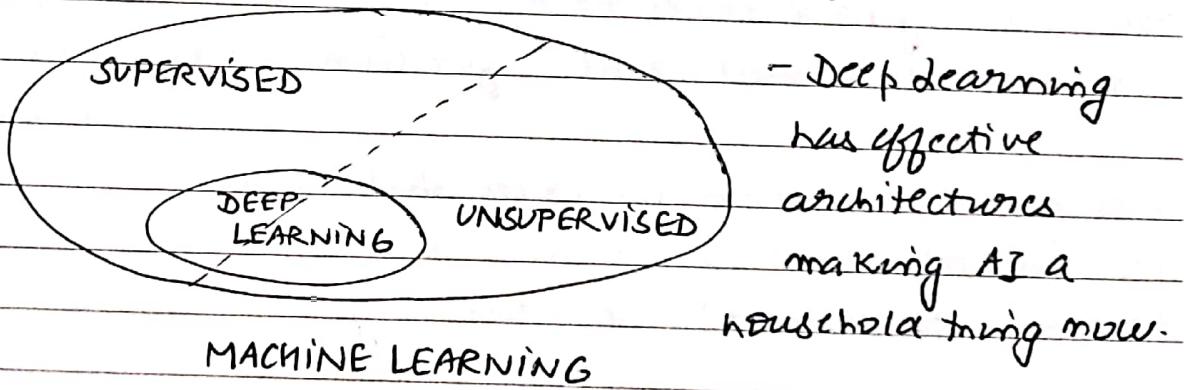
SL & ML differ in origin, history & emphasis.

PATTERN RECOGNITION (PR): typically an engg. problem to predict-based on patterns hidden in data

ML is a solution to solve problem of PR

## ML v/s DM

- DM is a process, while ML deals with science of creating algo that will make machines learn to predict
- DM is interactive & iterative process, while once a ML algo is designed, no human intervention is required.
- DM focused on scalability since its inception in 80's, ML has addressed scalability recently. (since deep learning came 5-6 yrs ago)
- Traditionally, ML didn't consider data cleaning etc. This routine of data cleaning was adopted recently.



- Deep learning has effective architectures making AI a household thing now.

There is no need to learn in case of finding prime no.

Learning is different from normal coding  
 Learning is reqd. when human expertise doesn't exist  
 (e.g. unknown territory) or when humans are unable to explain their expertise (face & speech recognition) or when env changes with time. (driving on highway vs in city)  
 or when solut<sup>n</sup> needs to be adapted on case to case basis (diagnosing disease)

Date :

Deep Neural m/w machines make very accurate decisions but almost all of them fall short in explaining them.

Machines can't explain but can figure out the patterns humans can't

HYPOTHESIS: Build a model that is a good & useful approximation to the data.

Machine learning algo generalizes an outcome based on specific instance values (eg: for a malignant cell cell wall thickness lies b/w 0.04 & 0.06).

Metrics to establish reln b/w attributes

Machines only have power to model. From the data provided to machine, it must create a function  $f(\underbrace{x_1, \dots, x_m}_{m \text{ attribute values}}) \rightarrow \{\text{dis, not dis}\}$

Finding such a function results in abstraction

ML Algo attempt to approximate this function  
A better Approximation  $\propto$  Accuracy ABSTRACT

Gini, Info gain, Gain Ratio are 3 diff. hypotheses for the same task in decision tree

Data are cheap & abundant, labelling is often expensive

## APPS OF ML

- Association
- supervised learning
  - classification
  - regression
- unsupervised learning
- Reinforcement learning ( $\alpha$ -GO)

CLASSIFICATION : diff b/w low risk & high risk customers from their income & savings

hypothesis discriminant : if  $\text{income} > \theta_1$  &  $\text{savings} > \theta_2$   
THEN Low-Risk  
ELSE High-Risk

REGRESSION : hypothesis  $\rightarrow$  eqn-of-line

Date:

## CHAPTER-2 SUPERVISED LEARNING

MACHINE LEARNING → Theoretical ML  
Applied ML

GUESS THE INTERVAL  $[a, b]$

Goal of ML is to abstract a function but since it can't be predicted accurately, we approx. it.

A hypothesis must be consistent with data if there are no conflicts.

Move the data, better the hypothesis.

ID	INSIDE	Given Data	
		Hypothesis 1	Hypothesis 2
65	Yes	$(-\infty, \infty)$	$[0, 100]$
125	Yes	$(-\infty, \infty) \checkmark$	X becomes inconsistent
-80	ND	$[0, \infty) \checkmark$	We improve it as we have more data. we make it $[0, \infty)$
-53	Yes	$[-79, \infty)$ → $[-60, 200]$	We can improve it as $[-79, \infty)$
101	Yes	$[-60, 200]$	random values to tighten the range
110	ND	Creates conflict with 125, Yes ∴ This is considered as noise as it creates inconsistencies in data	
153	ND	$[-60, 152]$	

Date:

We evaluate a hypothesis using some metric.  
Error decreases as examples increase.  
classes are called concepts.

Concept learnt → inside the interval  
complementary concept → outside the interval

we evaluate a hypothesis using some metric  
Errors decreases as examples increase.  
classes are called concepts.

→ concept learnt → inside the interval  
complementary concept → outside the interval

### GOAL OF SUPERVISED LEARNING

- Accurately predict class labels for unseen instances
- understand which features are important for class & how
- Assess quality of classifier

predictor: class in classification, md in regression

Pattern recognition came before ML

- Concept is subset of instance space  $X$ .

For predicting whether a car  $C$  is family car ( $Y$ ) or not ( $N$ )

PRICE	ENGINE POWER	FAMILY CAR	instance space, training set
-	-	$Y$	concept of family car
-	-	$N$	
-	-	$Y$	
-	-	$Y$	concept of n-family car
-	-	$N$	
-	-	$N$	

Date:

Learning of concept: Find a description that's shared by all instances of class & NOT by any example of the other class.

For training set  $X$

$$X = \{x^t, y^t\}^N_{t=1}$$

as there are  $N$  instances in training set

$x^t$  is a vector of features

$y^t$  label for  $t^{th}$  instance

$$x = \left\{ \begin{array}{l} x_1 \rightarrow \text{feature 1} \\ x_2 \rightarrow \text{feature 2} \end{array} \right\}$$

in family car eg, price is a feature vector & so is engine power

$x_1^t$ :  $x_1$  feature of  $t^{th}$  instance

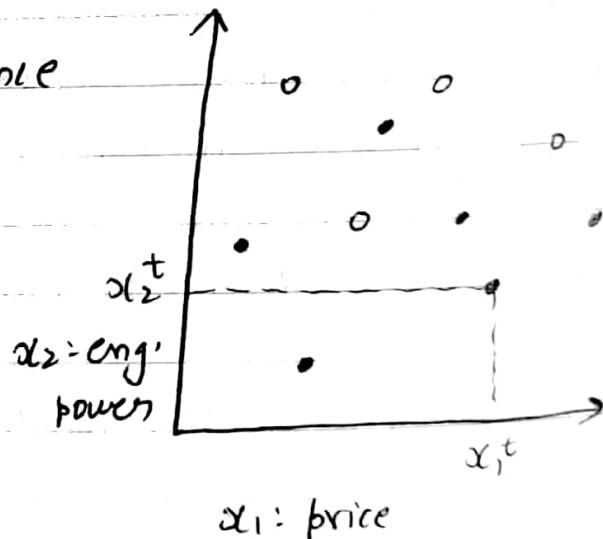
$x_2^t$ :  $x_2$  .., ..

$$y^t = \begin{cases} 1 & \text{if } x \text{ is +ve} \\ 0 & \text{if } x \text{ is -ve} \end{cases}$$

To learn a concept, of like 'yes' for family car, we must find constraints s.t

$$(p_1 \leq \text{price} \leq p_2) \& (e_1 \leq \text{eng. power} \leq e_2)$$

There can be  $\infty$  such constraint values, hypotheses



misforms  
a rectangle

All possible rectangles form hypotheses space

HYPOTHESES : educated guess

HYPOTHESES SPACE : set of possible hypotheses

GENERALIZATION : How well does hypotheses capture specific objects in training set to predict unseen objects (loss of precision & sometimes gain in understanding of description)

Machine can't be made intelligent by memorization

PAC : Probably Approximately Correct

EMPIRICAL ERROR : Model that machine constructs from data is just an approximation. ∴ there is some error which must be measurable, this measurement is empirical error

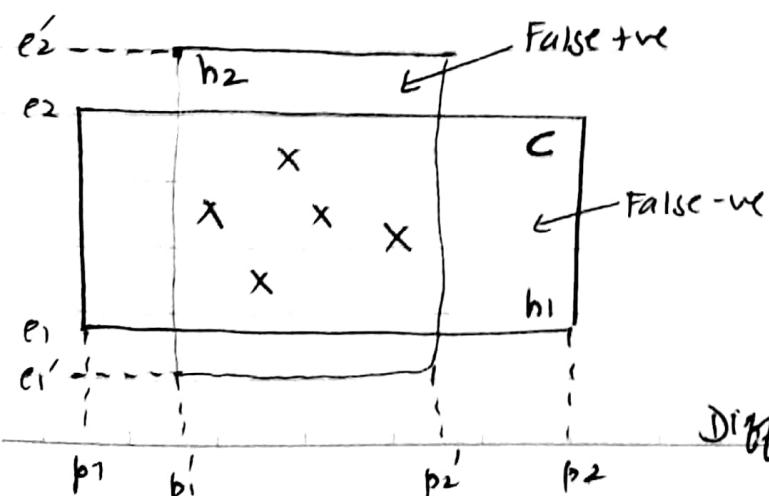
In predictive modelling, "approximate" is a theoretical concept

hypothesis :  $h$

Hypothesis space :  $H$  s.t.  $h \in H$

→ indicator funct<sup>n</sup>

→ Error of a hypothesis  $h = E(h|X) = \sum_{t=1}^N I(\text{if } \underline{h(x^t)} \neq \underline{y^t})$



denotes how many errors

choice of algs & parameters  
defines the hypothesis space

Diffr. algs have diff. hyp-space

Date : \_\_\_\_\_

task of algo is to navigate through  $H$  to obtain the best hypothesis  
This makes ML as an optimizat<sup>n</sup> problem

## CHAPTER-3

BAYESIAN DECISION THEORY

Bayesian learning is probabilistic learning.

$$\text{Estimat}^n: p_0 \text{ (for heads)} = \frac{\# \text{ Heads}}{\# \text{ Tosses}} = \frac{\sum_t x^t}{N}$$

$x^t = 1$  for heads

0 for tails

Prediction of next toss: heads if  $p_0 > 1/2$ , tails otherwise

Bernoulli Thm HW

Bayesian learning combines prior knowledge with observed data & provides practical learning algos :

1) Naive Bayes learning

2) Bayesian Belief NW learning

→ don't reject directly

Bayesian learning provides flexible approach to learning compared to algos that eliminate inconsistent hypothesis. B

Provides "gold standard" for evaluating other learning algos, even though sometimes they may be computationally intractable

Occam's razor NW

Bayesian learning provides the bound of error

Date :

How to decide class without seeing  $x$ ?

A: Acc. to Bayesian theory, use prior probabilities of classes.

- choose  $C_1$  if  $P(C_1) > P(C_2)$ , else  $C_2$

If  $\exists K$  classes, choose the class with max. prior probability.

seeing features of an instance will lead to better decision making

## COURSE PREDICTION PROBLEM

No. of Eco students = 800      No. of engg. students = 1500

observe marks of Maths in 12<sup>th</sup> (we are given marks of maths)

collect data

plot histogram of Marks in Maths v/s Count / Probability  
for eco & engg. separately & notice the density of distribution

Notice average marks in both (assuming avg in E = 78, B = 88)

Assume avg. marks in Eco in Maths = 70

$$P(M=70|E) \quad P(M=70|B)$$

since avg. of Eco is closer to 70

engg.

$$\therefore P(M=70|E) > P(M=70|B)$$

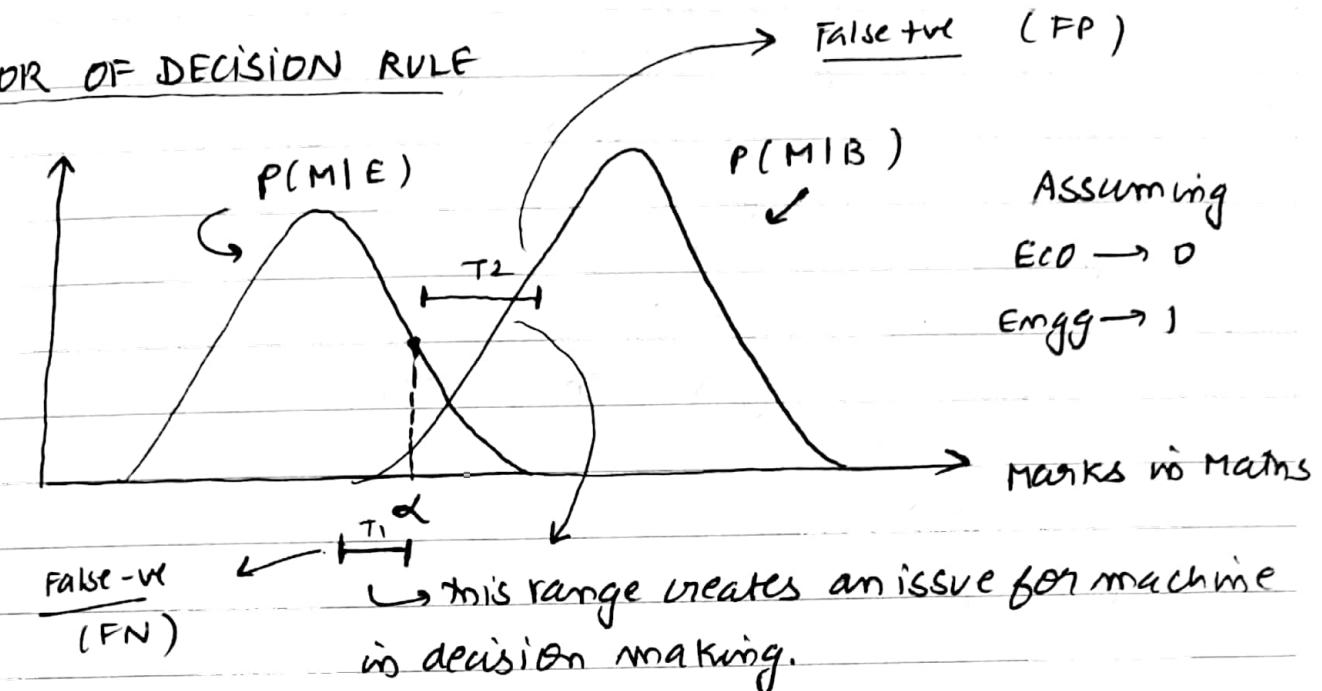
And why  $P(M=94|B) > P(M=94|E)$  this doesn't mean  
that high scorers in M always join B, they are ~~most~~ more likely  
to join B.

DECISION RULE

HYPOTHESIS : let there be some marks threshold  $\alpha$   
if marks  $\leq \alpha \rightarrow$  Eco

else  $\rightarrow$  Engg.

What should be value of  $\alpha$ ?  $\rightarrow$  Based on populat<sup>n</sup> statistics

ERROR OF DECISION RULE

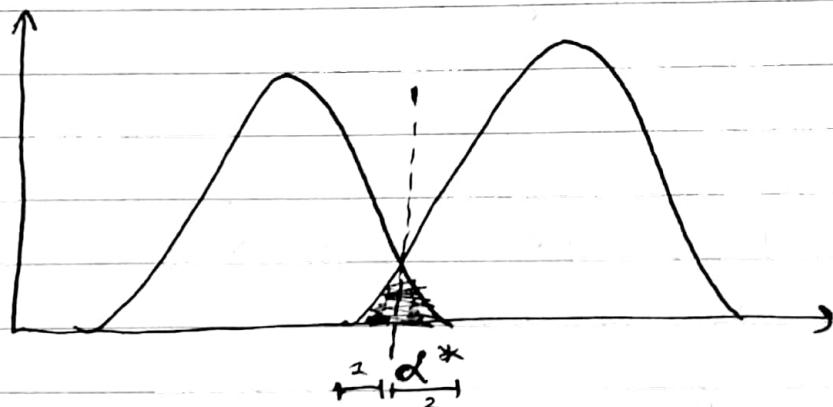
$T_1$ : TYPE 1 : Actually BTech, classified ECO

$T_2$ : TYPE 2 : " ECO , " BTech

More the area under the other curve, more is the occurrence

$\therefore$  FP occurs more than FN

$$\text{ERRORS} = T_1 + T_2$$

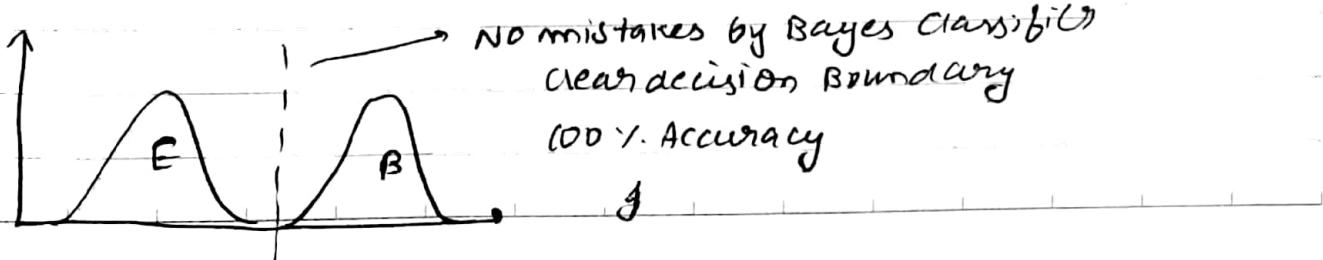


If  $T_1$  &  $T_2$  errors have diff. costs :  
optimal boundary shifts

$\alpha^*$  : optimal value of  $\alpha$  : Min-possible error

$$P(\alpha^* | E) = P(\alpha^* | M B)$$

Mistakes = 1+2 = area under overlap curve



Date :

The decision boundary in the plot = hypothesis of our model  
No classifier gives error < Bayes Error  $\rightarrow$  min. error probability  
Bayes Error can't always be computed. It depends on the complexity of data available.

The distribution changes with Location / Time of observation

- job opportunities
- cost of courses
- family background

### APRIORI PROBABILITY

- without measuring marks, what can we guess about the course student may opt?  $\rightarrow$  depends on locat' / time
- Apriori Prob.,  $P(E)$ ,  $P(B)$  : property of frequency of classes during exist. Not a property of marks of student

### PREDICTION

$$\begin{cases} C=1 \text{ if } P(C=1 | x_1, x_2) > 0.5 \\ C=0 \text{ otherwise} \end{cases}$$

OR

$$\begin{cases} C=1 \text{ if } P(C=1 | x_1, x_2) > P(C=0 | x_1, x_2) \\ C=0 \text{ otherwise} \end{cases}$$

## BAYES' RULE

$$P(C|x) = \frac{P(x|C) \cdot P(C)}{P(x)}$$

↓ likelihood      ↓ prior probability

↓ normalising term      ↓ evidence (object that we want to predict)

ULE

→ observed the data  
 → after we observe the data)

- $P(C=0) + P(C=1) = 1$
  - $P(x) = P(x|C=1) \cdot P(C=1) + P(x|C=0) \cdot P(C=0)$
  - $P(C=0|x) + P(C=1|x) = 1$

For 'd' dimensions / attributes, we call it joint probability.

$P(x_1, x_2, x_3, \dots, x_d) \rightarrow$  joint probability

a random variable  $X$  of 'd' dimensions where each dimension can take a value from a range of values. This is an instance of possible combination of values that can be assigned to the dimensions of  $X$ .

Frequency = 0  $\Rightarrow$  Prob. = 0

For plotting a histogram for 2 subjects, assuming  $X_{\text{coord}} =$   
Maths marks,  $Y_{\text{coord}} = \text{English marks}$ , Frequency of students  
with that many marks in M & E = Height of histogram

For 3 subjects, random variable  $X = [x_1 \ x_2 \ x_3]$   
 $\therefore$  we will have  $(100)^3$  possible combinations of values  
 If we have only 1000 records, then most of the instances will have prob = 0 & most non-zero prob instances would be very less in frequency.

$\therefore$  Bayes classifier is computationally intractable.

### FOR $> 2$ CLASSES (BAYES' RULE)

We must find posterior prob for each class

$$P(C_i|x) = \frac{P(x|C_i) \cdot P(C_i)}{P(x)} = \frac{P(x|C_i) \cdot P(C_i)}{\sum_{m=1}^K P(x|C_m) \cdot P(m)}$$

- $P(C_i) \geq 0$  &  $\sum_{i=1}^K P(C_i) = 1$
- DECISION RULE  $\rightarrow$  choose  $C_i$  if  $P(C_i|x) = \max_K P(C_k|x)$

### BAYES CLASSIFIER

APRIORI: Without any measurement based on just location / time - what can we guess about class membership (estimated from size of class populations)

CLASS CONDITIONAL: Given the instance belongs to a particular class, Likelihood: What is the prob that its (combination of) features exist (estimated from training set)

Eg:  $P(<80, 65 | E)$  we know the class is E, we check 60% instances having class = E with instance values =  $<80, 65$

optimal Bayes classification works on joint prob. distrib' of data  
 Bayes classification works on Bayes theorem

MAXIMUM A POSTERIORI HYPOTHESIS (MAP HYPOTHESIS / HMAR)  
 choose hypothesis with max. posterior probability

MAXIMUM LIKELIHOOD HYPOTHESIS : If all classes are equally likely  
(ML HYPOTHESIS) choose hypothesis which is most likely

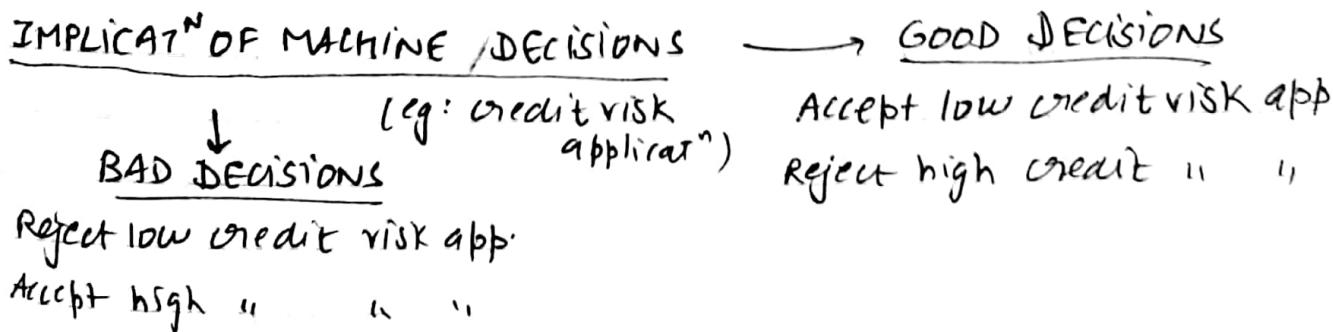
$$\text{eg: } P(C|x) = \underbrace{P(x|C) \cdot P(C)}_{P(x)} \rightarrow \text{this is same for all} \therefore \text{we classify}$$

a basis of  $P(x|C)$

$\rightarrow$  same throughout  $\therefore$  can be ignored

SALIENT FEATURES OF BAYESIAN LEARNING : This is a more flexible approach than many other algos

- PRACTICAL DIFFICULTIES :
- 1) Typically require initial knowledge of many probabilities
  - 2) Significant computational cost reqd. to determine Bayes optimal hypothesis



Reject low credit risk app  
 Accept high .. .. "

Date:

		PREDICTED			
		1	0		
ACTUAL	1	TP	FN	Rejecting low credit risk app = FN ∵ we wrongly reject a good app.	
	0	FP	TN	Accepting high credit risk app = FP	
		GOOD DECISIONS			

LOSSES & RISKS : Loss incurred in case of wrong prediction / decision for  $x$

Action  $\alpha_i$  : choose class  $c_i$

loss of  $\alpha_i$  when correct class is  $c_k$  :  $\lambda_{ik}$

Expected Risk for taking action  $\alpha_i$

$$R(\alpha_i | x) = \sum_i \lambda_{ik} P(c_k | x)$$

choose  $\alpha_i$  if  $R(\alpha_i | x)$

$$= \min_k R(\alpha_k | x)$$

LOSS MATRIX

loss related to

$\lambda_{ii}$ : correct class 1 predicted as 1

Possible to define loss funct" for predict" problem  
Sometimes it's not possible to define all possible combinations

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i=k \\ 1 & \text{if } i \neq k \end{cases}$$

$$R(\alpha_i | x) = \sum_i \lambda_{ik} P(c_k | x)$$

$$= \sum_{k \neq i} P(c_k | x)$$

$$= 1 - \underbrace{P(c_i | x)}$$

there is only 1 case when  $\lambda = 0$  when  $i = k$

For minimization, choose most probable class

REJECT OR DOUBT

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i=k \\ \lambda & \text{if } i=k+1, 0 < \lambda < 1 \\ 1 & \text{otherwise} \end{cases}$$

$$R(\alpha_{k+1}|x) = \sum_{k=1}^K \lambda P(C_k|x) = \lambda$$

$$R(\alpha_i|x) = \sum_{k \neq i} P(C_k|x) = 1 - P(C_i|x)$$

choose  $C_i$  if  $P(C_i|x) > P(C_k|x) \forall k \neq i$  &  $P(C_i|x) > 1 - \lambda$   
reject otherwise

DISCRIMINANT FUNCT<sup>N</sup> AS CLASSIFIER

- works for both binary & multi-way classification
- for every class  $i = 0, 1, \dots, K$  define a function  $g_i(x)$   
mapping  $X \rightarrow R$
- when decision of on I/P  $x$  is to be made choose class with highest value of  $g_i(x)$

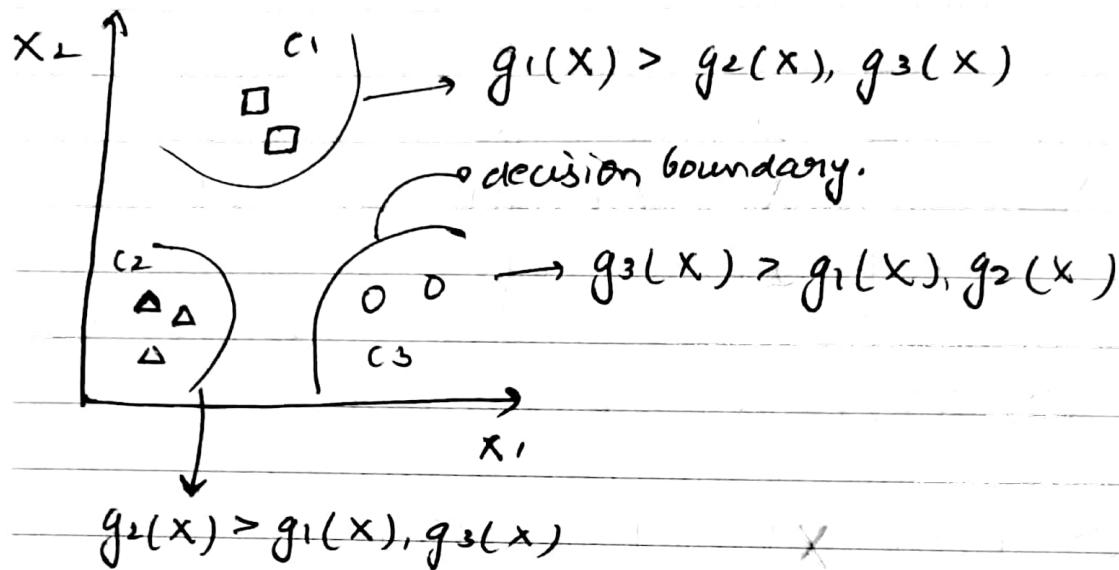
DISCRIMINANT FUNCTIONS FOR BAYES CLASSIFIER

$$g_i(x) = \begin{cases} -R(\alpha_i|x) \rightarrow \because g_i \text{ is to be maximised} & \\ P(C_i|x) & \text{risk is to be minimized} \\ P(x|C_i) \cdot P(C_i) & \end{cases}$$

Generates  $K$  decision regions  $R_1, \dots, R_K$

$$R_i = [x | g_i(x) = \max_k g_k(x)]$$

Date :



K=2 CLASSES      Dichotomizer (K=2) v/s Polychromizer (K>2)

UTILITY : Utility of  $\alpha_i$  if class is  $k = U_{ik}$

Expected Utility =  $EV(\alpha_i | x) = \sum_k U_{ik} P(s_k | x)$   
Choose  $\alpha_i$  with max utility

Value of info & increase in expected utility

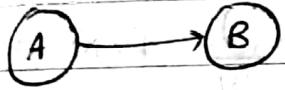
$$EV(\alpha_1 | x_1, 3) > EV(\alpha_1 | x)$$

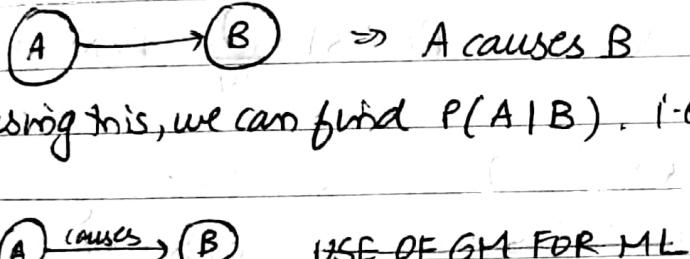
(GM) (UPTO 14.3)

CHAPTER-14 GRAPHICAL MODEL (PROBABILISTIC LEARNING)

- representat^n method for data where Nodes = variables and edge b/w nodes = relationship
- GOAL: extract a model which can be used to predict class label of an unseen object

Graphical Models are statistical models based on Bayesian learning

  $\Rightarrow$  A causes B eg: smoking causes lung cancer  
using this, we can find  $P(A|B)$ . i.e. we are trying to diagnose

  
 $\xrightarrow{\text{diagnostic}}$

- 1) Rel^n b/w features in training data are the quantities of interest for prediction

For Bayes classifier, you need joint prob. distribution to derive class conditional probabilities

For d-dimensional boolean data, we require  $O(2^d)$  instances to estimate joint prob. distribution for each class. We simplify by assuming that attribute values are conditionally independent given class  $C_j$  i.e.  $P(x=x_1, x_2, \dots, x_d | C_j) = \prod_{i=1}^d P(x_i | C_j)$

HYPOTHESIS FOR NAIVE BAYES CLASSIFIER

$$h_{NB} = \operatorname{argmax}_j P(C_j) \prod_{m=1}^d P(x_m | C_j)$$

Date : \_\_\_\_\_

Q-  $C = \text{PlayTennis} (\text{Yes} | \text{No}) \quad |C| = 2$

$$|\text{OUTLOOK}| = 3$$

$$|\text{Temp}| = 3$$

$$|\text{Humidity}| = 2$$

$$|\text{Wind}| = 2$$

<sup>same</sup>  
If 36 combo are there for Yes & for No then there  
is no distinction

$3 \times 3 \times 2 \times 2$  combinat<sup>n</sup>s are possible

= 36 combinations

$$P(C|X) = [\pi P(X|C)] \cdot P(C)$$

We need 8 prob. values given X

$$P(\text{Yes}) = \frac{9}{14} \quad P(\text{No}) = \frac{5}{14}$$

If OUTLOOK :  $P(\text{sunny} | \text{Yes}) = 2/9$

$$P(\text{sunny} | \text{No}) = 3/5$$

$$P(\text{overcast} | \text{Yes}) = 4/9$$

$$P(\text{overcast} | \text{No}) = 0$$

$$P(\text{Rain} | \text{Yes}) = 3/9$$

$$P(\text{Rain} | \text{No}) = 2/5$$

Exact running  
time is  $\Theta(n^k)$

dependent

WIND :  $P(\text{weak} | \text{Yes}) = 6/9$

$$P(\text{weak} | \text{No}) = 2/5$$

$$P(\text{strong} | \text{Yes}) = 3/9$$

$$P(\text{strong} | \text{No}) = 3/5$$

TEMPERATURE :  $P(\text{Hot} | \text{Yes}) = 2/9$

$$P(\text{Hot} | \text{No}) = 2/5$$

$$P(\text{Mild} | \text{Yes}) = 4/9$$

$$P(\text{Mild} | \text{No}) = 2/5$$

$$P(\text{Cool} | \text{Yes}) = 3/9$$

$$P(\text{Cool} | \text{No}) = 1/5$$

All these are calculated to  
create the model

we maintain  $(2 \times (\underbrace{3+3+2+2}_{\text{Yes/No}}))$   
 $= 20$  attribute value

counters for the model

HUMIDITY :  $P(\text{High} | \text{Yes}) = 3/9$

$$P(\text{High} | \text{No}) = 4/5$$

$$P(\text{Normal} | \text{Yes}) = 6/9$$

$$P(\text{Normal} | \text{No}) = 1/5$$

$O(mn)$

Time complexity  $\approx O(m)$

Space complexity =  $O(C \times (\underbrace{A_1 + A_2 + \dots + A_d}_{\text{class value}}))$   
 $\downarrow$   
 $\text{no. of A}_i \text{ values}$

Date :

not

Predict the label for  $\langle \text{sunny}, \text{high}, \text{high}, \text{weak} \rangle$

$$P(\text{No}) = \pi(P(X_1 \mid \text{No})) \\ = \frac{8}{14} \times \frac{3}{8} \times \frac{2}{5} \times \frac{4}{5} \times \frac{2}{5} = \frac{24}{35 \times 25} = \frac{24}{875} = 0.0274$$

$$P(\text{Yes}) = \pi(P(X_1 \mid \text{Yes})) \\ = \frac{8}{14} \times \frac{2}{8} \times \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} = \frac{4}{63 \times 9} = \frac{4}{567} = 0.00705$$

not

Predict the label for <sunny, high, high, weak>

$$P(\text{No}) \cdot \prod (P(x_i | \text{No}))$$

$$= \frac{8}{14} \times \frac{3}{8} \times \frac{2}{5} \times \frac{4}{5} \times \frac{2}{5} = \frac{24}{35 \times 25} = \frac{24}{875} = 0.0274$$

$$P(\text{Yes}) \cdot \prod (P(x_i | \text{Yes}))$$

$$= \frac{8}{14} \times \frac{2}{8} \times \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} = \frac{4}{63 \times 9} = \frac{4}{567} = 0.00705$$

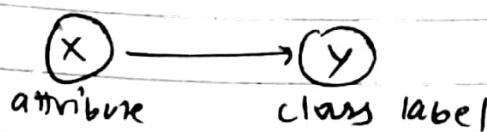
LAB WORK : Take a labelled dataset from UCI repository & find 10-fold cross validation accuracy.

- Naive Bayes is linear in time, popular because despite strong assumption, it often works surprisingly well. ] FEATURE OF NAIVE BAYES

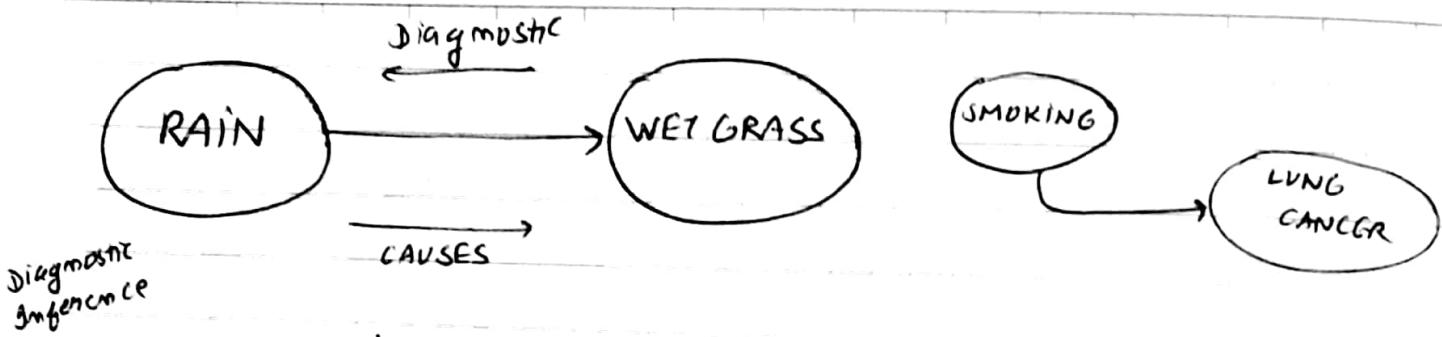
- Bayesian belief networks provide an intermediate approach which allows strong conditional independence assumptions (based on domain knowledge) that apply to subsets of variables.

$x_1, \dots, x_f$  have  $b_{C_2}$  relations  
The structure is represented as a directed acyclic graph (DAG)

Parameters are conditional prob. in arcs



Date :



$$P(R) = 0.4, P(W|R) = 0.9, P(W|\sim R) = 0.2$$

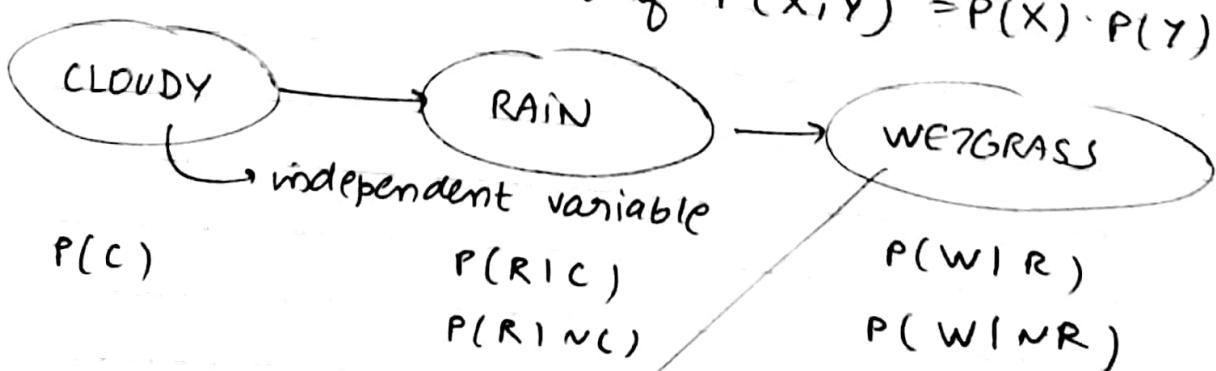
→ Knowing that grass is wet, what's the prob. that rain is cause

$$\begin{aligned} P(R|W) &= \frac{P(W|R) \cdot P(R)}{P(W)} = \frac{P(W|R) \cdot P(R)}{P(W|R) \cdot P(R) + P(W|\sim R) \cdot P(\sim R)} \\ &= \frac{0.9 \times 0.4}{0.9 \times 0.4 + 0.2 \times 0.6} = 0.75 \end{aligned}$$

If  $X$  is independent &  $X \rightarrow Y$ , then we must have  $P(X)$

### CONDITIONAL INDEPENDENCE

$X$  &  $Y$  are independent if  $P(X, Y) = P(X) \cdot P(Y)$



for finding  $P(W)$   
it doesn't matter  
if  $C$  occurs or not

$$P(W) = P(W|R) \cdot P(R) + P(W|\sim R) \cdot P(\sim R)$$

$$\begin{aligned}
 P(XYZ) &= P(X) \cdot P(Y|X) \cdot P(Z|XY) \\
 &= \cancel{P(X)} \cdot \frac{\cancel{P(Y|X)}}{P(X)} \cdot \frac{P(Z|Y)}{\cancel{P(XY)}}
 \end{aligned}$$

$$= \frac{P(X \wedge Y) \cdot P(Y \wedge Z)}{P(Y)}$$

 $\not\exists$ 

$$\begin{aligned}
 P(Z|X,Y) &= P(X,Y,Z) / P(X,Y) && P(Y|X) \\
 &= P(X) \cdot P(Y|X) \cdot P(Z|Y) / P(X) \cdot P(X+Y) \\
 &= P(Z|Y) && [\text{event } Y \text{ subsumes } X] \\
 &&& [\text{Event } Y \text{ blocks event } X \wedge Z]
 \end{aligned}$$

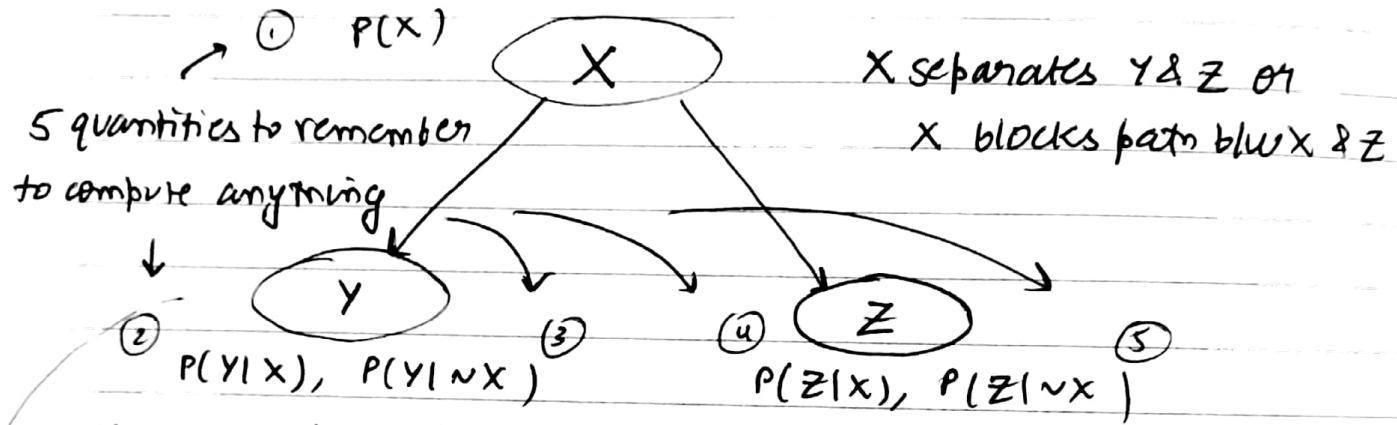
$P(Y), P(Z)$  can be computed using law of total probability  
 causal inference  $P(Z|X)$   
 diagnostic "  $P(X|Z)$

people.dur.ac.vin/nvbatmagare

$$\begin{aligned}
 P(W|C) &= P(W|R) \cdot P(R|C) + P(W|\sim R) \cdot P(\sim R|C) \\
 P(C|W) &= \frac{P(C|W) \cdot P(W)}{P(W|C) \cdot P(C)} \\
 &\quad \swarrow \text{causal inference} \qquad \swarrow \text{Diagnostic inference}
 \end{aligned}$$

Date :

If a single variable influences multiple variables, i.e. for eg:  
 $X$  influences  $Y$  &  $Z$



$P(Y|X) + P(Y|\sim X) \neq 1$  always because  $Y$  might also depend on some other factors. Same as  $Y$  for  $Z$ .

apart from  $X$

- $P(X, Y, Z) = P(X) \cdot P(Y|X) \cdot P(Z|X)$ , = joint probability of such a case

To check independent rel<sup>n</sup> b/w  $Y$  &  $Z$  given  $X$

$$P(YZ|X) = \frac{P(XYZ)}{P(X)} = \frac{P(X) \cdot P(Y|X) \cdot P(Z|X)}{P(X)} = P(Y|X) \cdot P(Z|X)$$

$$P(Y) = P(Y|X) + P(Y|\sim X) \cdot P(Y|X) \cdot P(X) + P(Y|\sim X) \cdot P(\sim X)$$

$$P(Z) = P(Z|X) \cdot P(X) + P(Z|\sim X) \cdot P(\sim X)$$

How to compute  $P(Y|Z)$  given we don't know they are independent

We know  $P(X|Y)$  &  $P(X|Z)$   
are diagnostic

$$\cancel{P(Y|Z) = P(Z|Y) \cdot P(Y)} \\ \cancel{\quad\quad\quad P(Z)}$$

Date : ↘

we need this causal for calculating diagnostic

$$P(X|Y) = \frac{P(Y|X) \cdot P(X)}{P(Y)}$$

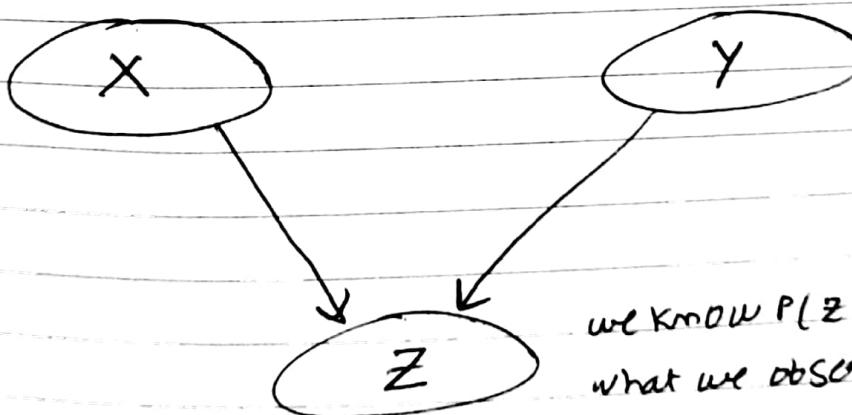
~~$$P(Y|Z) = \frac{P(YZ)}{P(Z)} = \frac{P(Y) \cdot P(Z)}{P(Z)}$$~~

↖

$$P(Y|Z) = \frac{P(YZ)}{P(Z)} = \frac{P(Y) \cdot P(Z)}{P(Z)} = P(Y) \quad \text{we don't use } P(X)$$

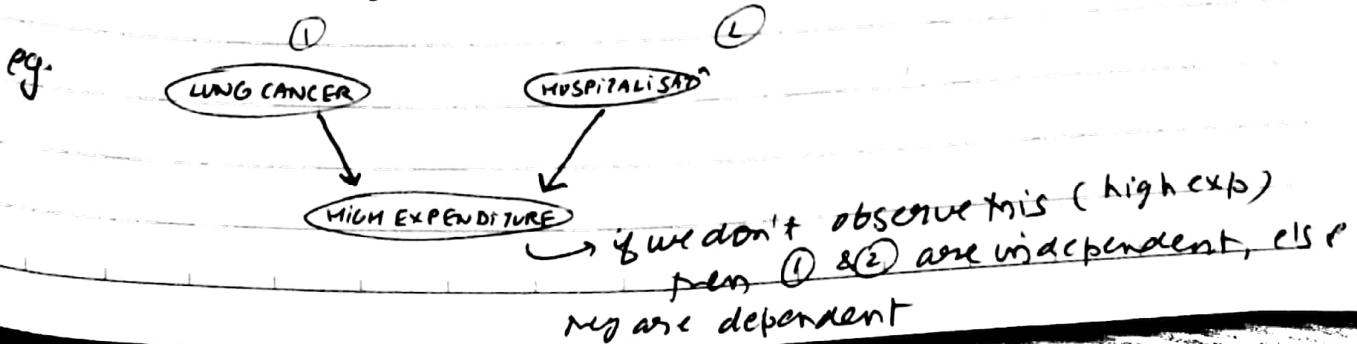
as we are assuming  
X is not known

$$\begin{aligned} &= P(Y|X) \cdot P(X) + P \\ &= P(Y|X) \cdot P(\emptyset X|Z) + P(Y|\sim X|Z) \end{aligned}$$



- $P(X, Y, Z) = P(X) \cdot P(Y) \cdot P(Z|X, Y)$

we can also obtain  $P(Z|\sim XY)$ ,  $P(Z|X\sim Y)$ ,  $P(Z|\sim X\sim Y)$  apart from  $P(Z|XY)$  from the provided data



Date :

$$P(W|R,S) + P(W|R,\sim S) + P(W|\sim R,S) + P(W|\sim R,\sim S)$$

can be  $> 1$

SPRINKLER      RAIN      Given       $P(W|R,S) = 0.95$   
↓                  ↓  
WET GRASS            $P(W|R,\sim S) = 0.9$   
 $P(W|S) = \text{causal event}$        $P(W|\sim R,\sim S) = 0.1$   
 $P(S|W) = \text{diagnostic}$        $P(S) = 0.2, P(R) = 0.4$

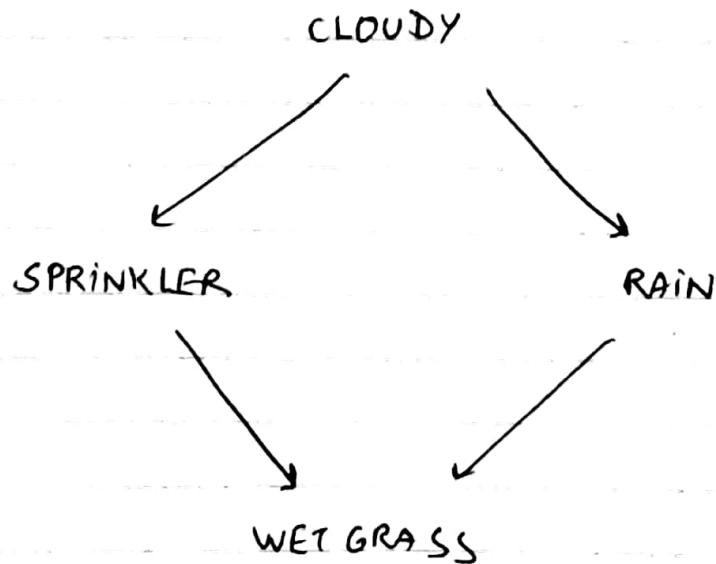
$$P(W) = P(W|RS) \cdot P(R) + P(W|\sim RS) \cdot P(\sim R)$$

using total prob. thm:

$$P(W) = P(W|RS) \cdot P(R) + P(W|\sim RS) \cdot P(\sim R)$$

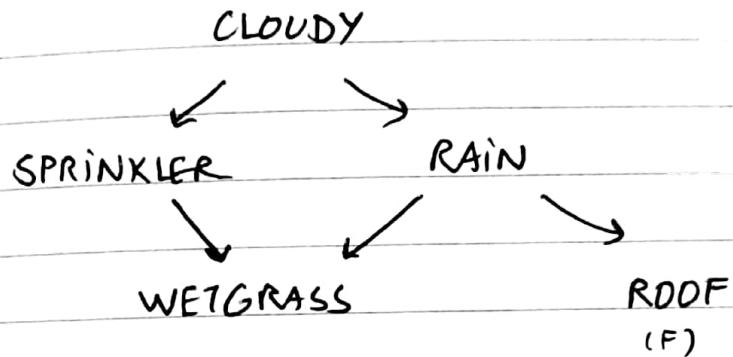
$$P(W|S) = P(W|RS) \cdot P(R|S) + P(W|\sim RS) \cdot P(\sim R|S)$$
$$= P(W|RS) \cdot P(R) + P(W|\sim RS) \cdot P(\sim R)$$

Knowing it has rained decreases the prob. that sprinkler is on



## EXPLOITING LOCAL STRUCTURE

How to extend Bayesian NW to a newly added variable?



$$\begin{aligned}
 P(CSRWF) &= P(C) \cdot P(S|C) \cdot P(R|C) \cdot P(W|SR) \cdot P(F|R) \\
 &= \prod_{i=1}^d P(X_i | \text{parents}(X_i))
 \end{aligned}$$

**RAIN** Since R is ind.  $\therefore$  we know  $P(R)$  & also  $P(W|R)$

↓ In bayes rule, we choose C & then calc. Prob.

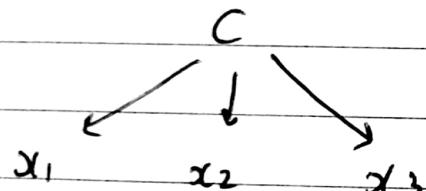
In class., we calc. prob. of a class given data

Bayes' rule inverts the arc

we have class conditional prob. from the data

Date : \_\_\_\_\_

## NAIVE BAYES CLASSIFIER



Given  $C$ ,  $x_j$  are independent

$$P(x|C) = P(x_1|C) \cdot P(x_2|C) \cdot P(x_3|C)$$

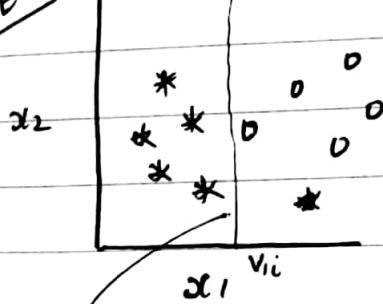
## DECISION TREES (CH-9)

→ works on Divide & Conquer

- It is a model induced from Training Set
- can be used for classification & regression
- prob. is involved in a very limited way.

CASE 1

### CLASSIFICATION TREES



→ 2 attributes + there would be 1 class

value  $\therefore = 3$  columns in training set

DISCRIMINANT FUNCTION : decision boundary b/w

( 2 classes divides data

) space into regions

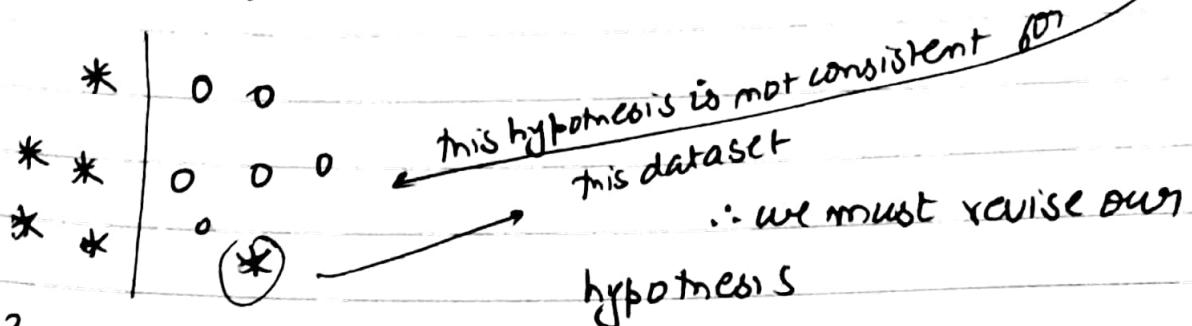
Each region contains data points of a particular class.

Can we say discriminant function is classifier?

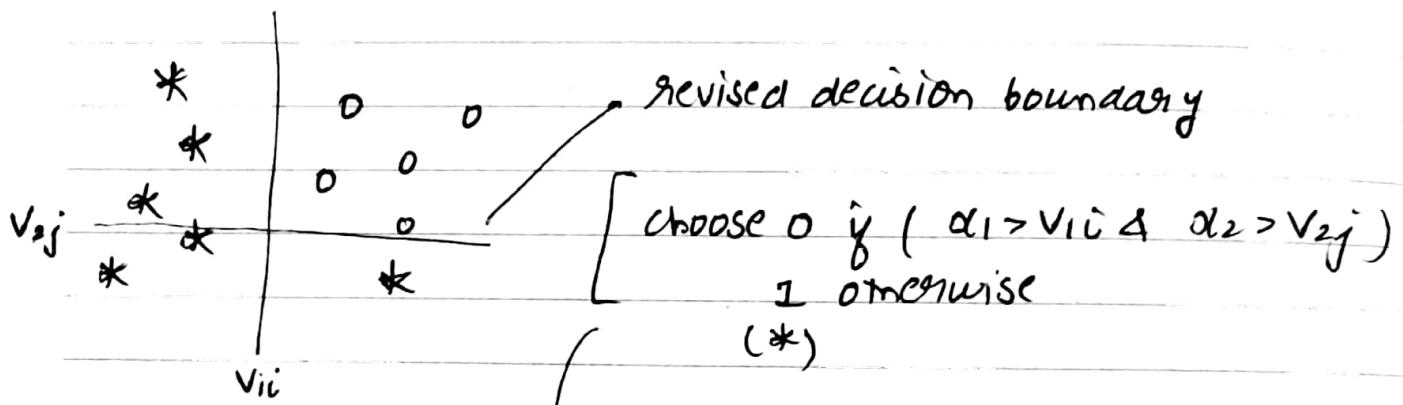
YES, we can. → choose 0 if  $x_1 > v_{1i}$  ] HYPOTHESIS  
 (\*) I otherwise

In reality, we don't get such perfect decision boundaries

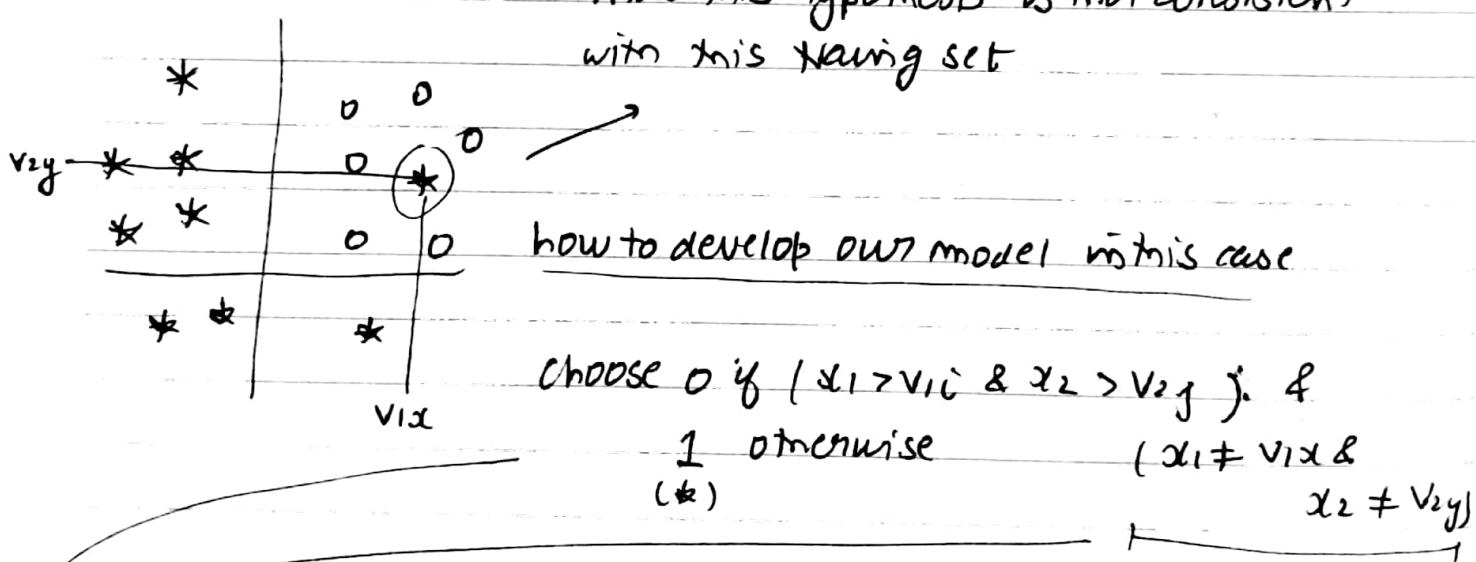
CASE 2



Date :



CASE 3



→ this hypothesis is particular for this data point, ∴ we tend to overfit our model over the provided data. We are making our hypothesis complex in order to be consistent.

→ regions that do not lie in it (or on its boundary)

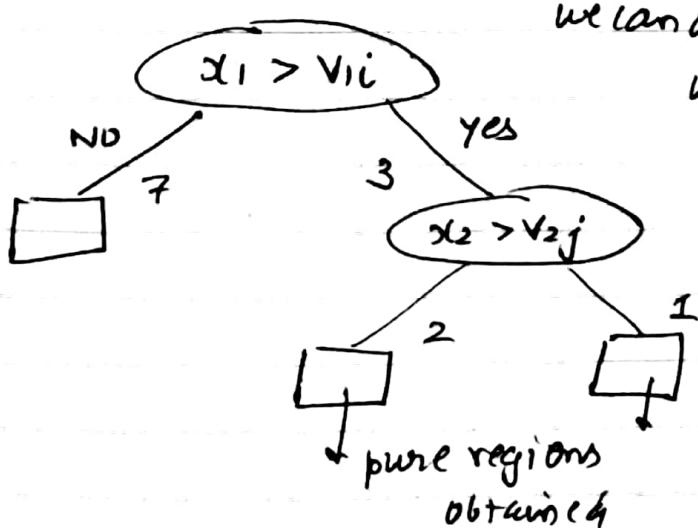
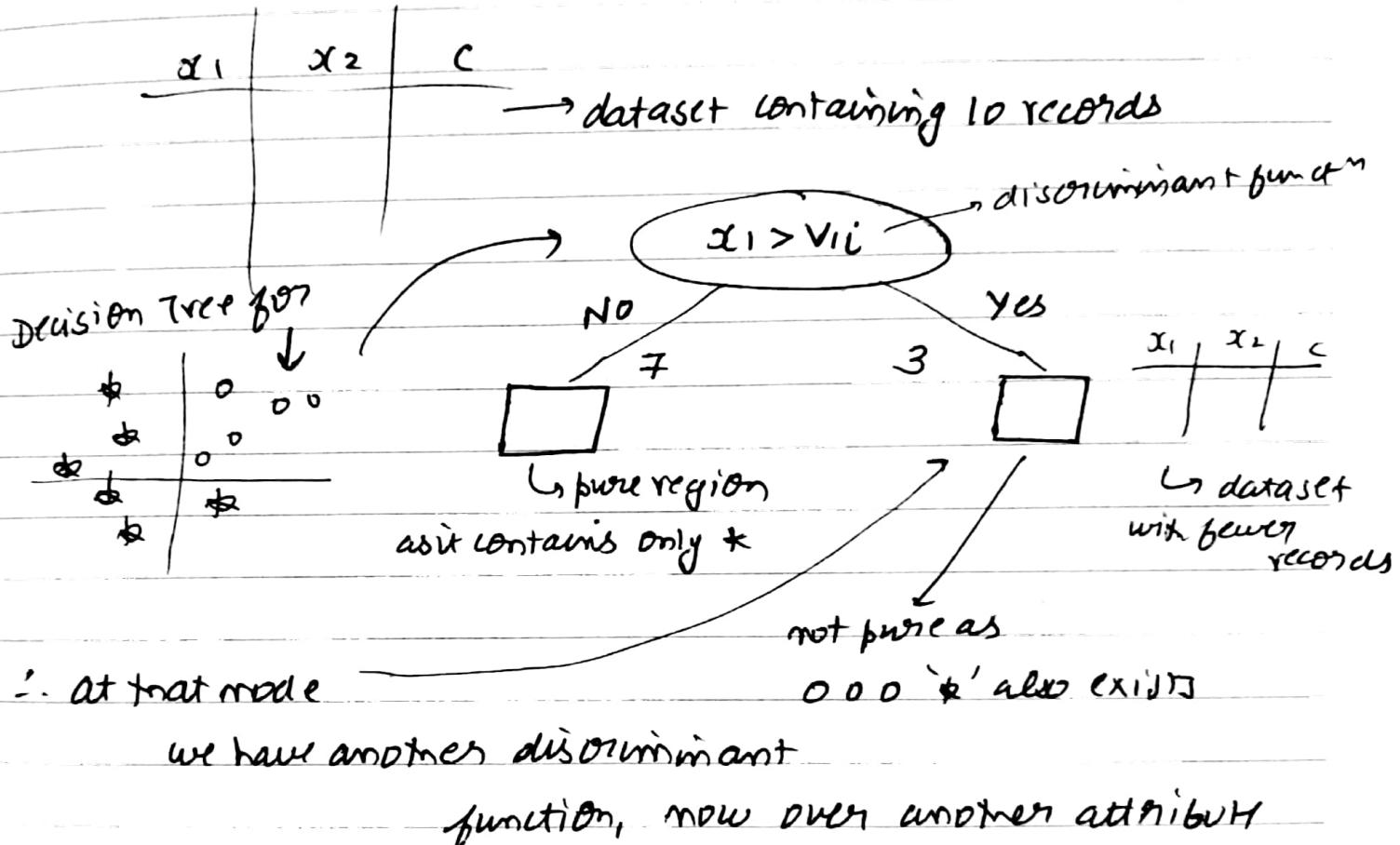
→ a better hypothesis would be to specify a non-overlapping region around  $(V_{1j}, V_{2j})$  instead of simply specifying coordinates

Decision tree provides a recursive method to divide the region until each pt is in a class.

PURE REGION : all data pt of 1 class in a region

→ aim of our hypothesis

Diagram: Root & internal nodes denoted by  leaf nodes denoted by 



Date : \_\_\_\_\_

Internal Decision Nodes : univariate  $\rightarrow$  using single attr.  $x_i$   
multivariate  $\rightarrow$  1, all attrs.,  $x$

Leaves : classification  $\rightarrow$  class labels or proportions  
regression  $\rightarrow$  Numeric ; n average or local fit

### FORM OF DECISION NODES

A	O	V
(attribute)	(operator)	(value)

Discrete



categorical



ordinal

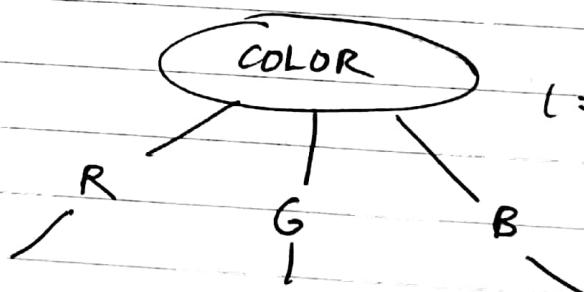
Nominal

- { Rank      { City  
  { Good, Better,  
    best }      { colors  
                  { red, green, blue }

for numeric values ( $>$ ,  $<$ ,  $\geq$ ,  $\leq$ )  
attr. binary split

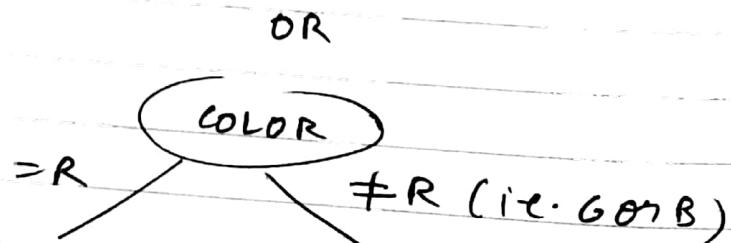
Decision on discrete values

→ m-way or  
binary split  
( $=$ ,  $\neq$ ) usually  
used



learning in decision

tree is greedy ; finds  
the best split recursively  
till some terminating  
criteria is met



Decision Tree algo is a recursive partitioning algo.  
Doesn't usually give accurate results but is popular  $\because$  it is transparent

& we can explain our decision at each step

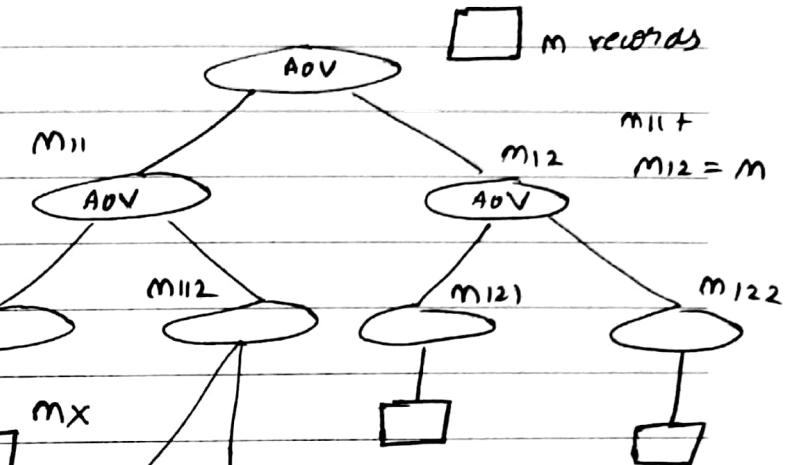
CLASSIFICATION TREES

(ID3, CART, CHS)

popular recursive partitioning algos

we define a region is pure using probability.

$$\begin{aligned} &= m_{11} \\ &= m_{111} + m_{112} + m_{121} \\ &\quad + m_{122} \quad ] = m_{12} \\ &= m \end{aligned}$$



To check purity of  $m_x$

$$P_{m_x}^i = \frac{\# \text{instances of } C_i}{|m_x|}$$

total no of records at  $m_x$

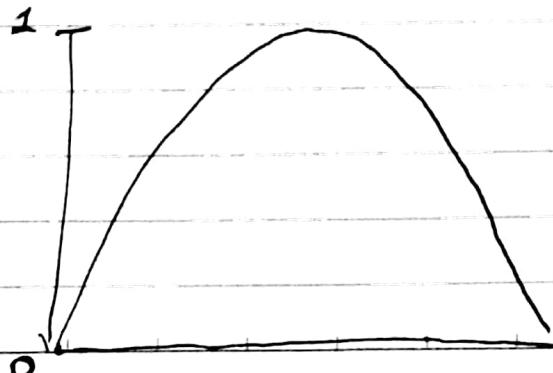
Node  $m_x$  is pure if  $P_{m_x}^i = 0$  or  $1$

~~$f(m_x) = N_{m_x}^i$~~  = # instances of class  $i$  at node  $m$

( $\hat{n}$ ) → denotes estimate

$$P(C_i | x, m_x) = P_{m_x}^i = \frac{N_{m_x}^i}{(N)} \rightarrow |m_x|$$

ENTROPY : A more rigorous way to check purity



entropy is a measure  
of disorder  
in a system

entropy  $\propto$  disorder

Date :

$$I_m = - \sum_{i=1}^K p_m^i \log_2 p_m^i$$

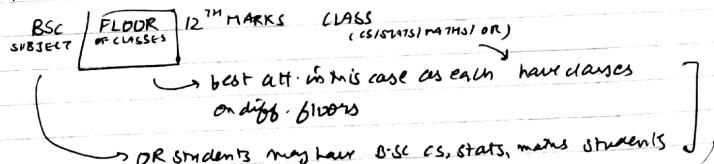
*informat<sup>n</sup>*  $\rightarrow p_m^i$  is in  $[0, 1] \therefore \log_2 p_m^i < 0$   
 in bits  $\therefore$  we have -ve sign  
 entropy measures the *ext. informat<sup>n</sup>* outside  
 in the data

If there are  $K$  classes each with  $K$  records, then entropy is max  
 $\& \text{prob. } = \frac{1}{K}$  for each class,  $0 \log_2 0 = 0$

Entropy is min i.e. = 0 if there is only 1 class

BEST SPLIT : If node  $m$  is pure, generate a leaf & stop,  
 otherwise split & continue recursively

We should be able to predict with least questions possible  
 pick attribute with highest discriminating power



how to decide this?  $\rightarrow$  using measures of goodness of split

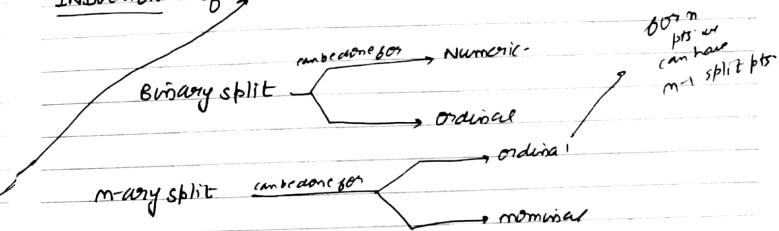
Decision Tree is a non-parametric, hierarchical model  
 importance of discriminant function decreases down the tree  
 / Power  
 $\therefore$  Root node has the discriminant funct<sup>n</sup> of highest power, importance

Date :

Each leaf defines a localised region in  $d$ -dimensional data space  
 with high density of data pts. of a class.  
 Decision Trees are one-way structures

$$\text{Entropy} \propto \frac{1}{\text{Purity}} \propto \text{Impurity}$$

INDUCTION:  $f_b \rightarrow$



How to do split for nominal?

$\rightarrow$  by grouping  
 we find all possible subsets to get discriminant function

Eg: weight is a continuous attribute,  $\therefore$  we can have  $m$  split pts even for 100 records. We select best hypothesis (split pts) from hypothesis space

$$I_m = - \sum_{j=1}^m \frac{N_{mj}}{N_m} \sum_{i=1}^{(K)} p_{mj}^i \log_2 p_{mj}^i$$

Date : \_\_\_\_\_

CHAPTER - 19

# DESIGN & ANALYSIS OF MACHINE LEARNING EXPTS.

Each decision tree is a hypothesis

How to decide which DI will perform best when deployed?

How to quantify / assess quality of prediction?

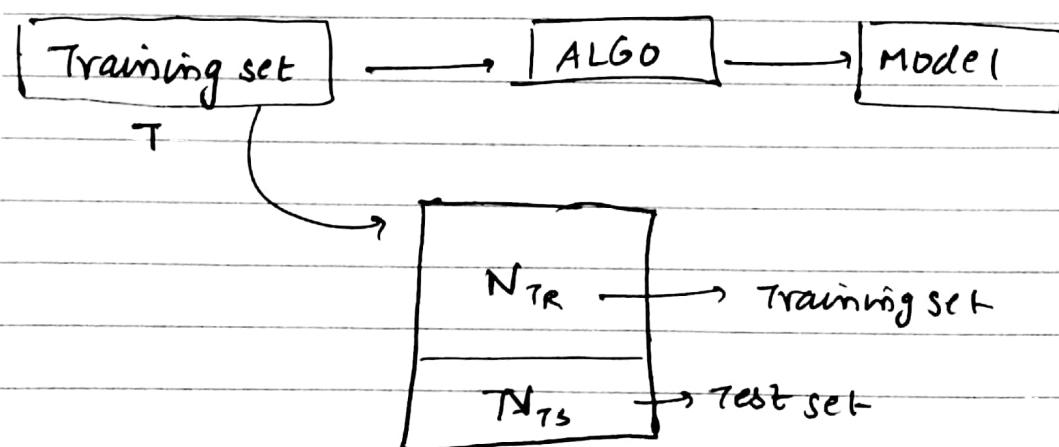
We can create diff. DT each with its own impurity threshold  
we choose the one with min-error / best accuracy.

$\Delta T_i$ with	$\rightarrow \theta_1$	$\theta_2$	$\theta_3$
threshold	hypothesis I		
impunity		hyp. 2	hypothesis 3
$\beta_i$			

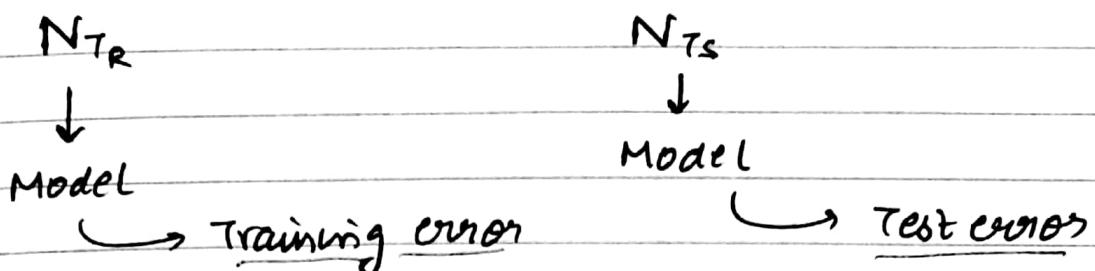
## Training error, Test error

Model is a hypothesis

"Classificat" model constructs a hypothesis

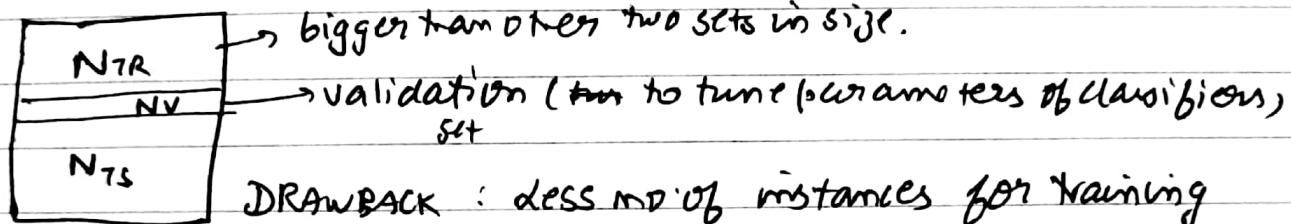


$$|T| = N_{T_R} + N_{T_S}$$



Training error < Test error  $\therefore$  Our model is based on training error  
it might be that N<sub>IS</sub> contains well behaved records which might  
degrade performance of our model when deployed in field.

To overcome this, we split our data into  $N_{tr}$  &  $N_{ts}$  multiple times randomly, so that test set can't be predicted



MOTIVATION → Model selection v/s Model estimation

↳ Assessing quality of classifier / model / algo

→ No free lunch thm (formalizes the statement that there is no universally best ML Algo)

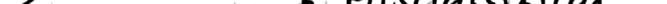
OR

On an average, all algos perform equally badly

One algo may work accurately for one data but may perform worse for other data ( even from same domain )

No algo works equally well on all datasets

→ Dividing dataset into  $N_{tr}$ ,  $N_v$ ,  $N_{ts}$  reduces size of training data & time to build the model

Algo preference on  Time complexity Misclassification error

## EVALUATING

• K-cross Validation

→ Holdout

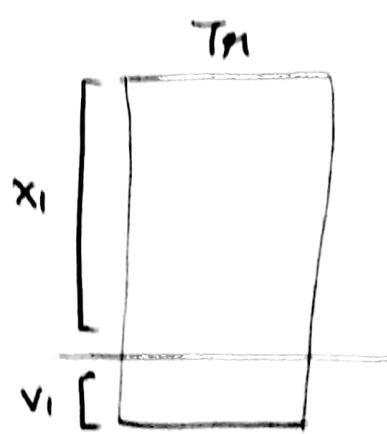
→ Bootstrap

## RESAMPLING & K FOLD CROSS VALIDATION

K partitions → (K-1) for training  
→  $K^{\text{th}}$  for testing

# iterations = K

this makes sure all distributions are represented equally in  $T_r$  &  $T_s$  set



e.g.: for 100 instances, if  $K=10 \rightarrow 10$  folds

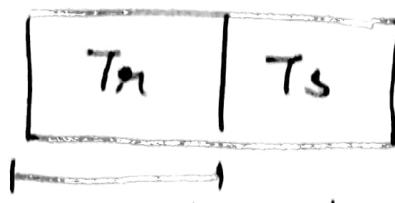
$\Rightarrow 90$  records for  $T_r$  & 10 records for  $T_s$

# of iterations  $\neq$  # of models = K

We take avg performance of classifiers over K iterations. Every instance is used  $(K-1)$  times for training & exactly once for testing

NOTATION:  $\{x_i, v_i\}$ :  $T_r$  &  $T_s$  validation sets of  $i^{\text{th}}$  fold  
For smaller dataset, 10 fold CV is not option  
 $\hookrightarrow$  do 5 times 2 fold CV

## HOLDOUT

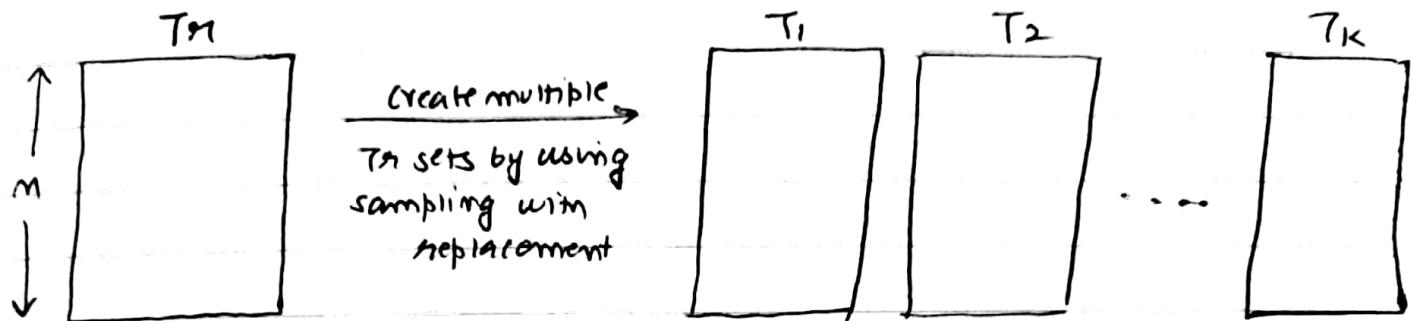


$\rightarrow$  holding this throughout for training

## BOOTSTRAPPING

- draws instances from data set with replacement

→ done for comparatively smaller datasets.



There will be only 36.8% new instances

→ over multiple runs, all instances can be used in  $T_m$  as well as  $T_s$

### PERFORMANCE MEASURES

$$\text{Total} = TP + FP + FN + TN$$

#### CONFUSION MATRIX

		PREDICTED	
		1	0
ACTUAL	1	TP	FN
	0	FP	TN

$$\frac{TP + TN}{\text{Total}} = \text{Accuracy} \quad (\text{ideally } 1 \text{ or } 100\%) , \text{ error rate} = 1 - \text{Accuracy}$$

RECALL : Of all the instances, how many can model predict as true

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

$$= \frac{TN}{TN + FP} \quad (\text{recall for -ve}) = \text{specificity}$$

Precision : correct predictions of the class out of total true predictions

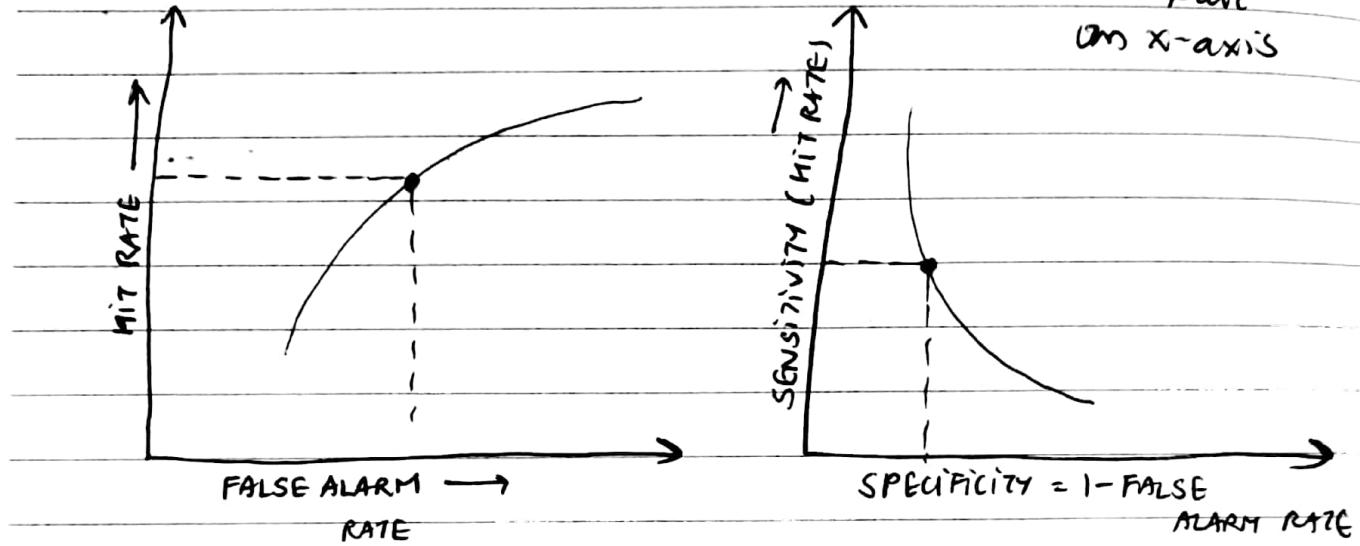
$$= \frac{TP}{TP + FP}$$

FALSE ALARM RATE :  $\frac{FP}{FP + TN}$  (causes extra expenditure & expense)

Date : \_\_\_\_\_

## ROC CURVE ( Receiver Operating characteristics )

curve reverses  $\therefore$  we took  $1 - \text{False Alarm rate}$

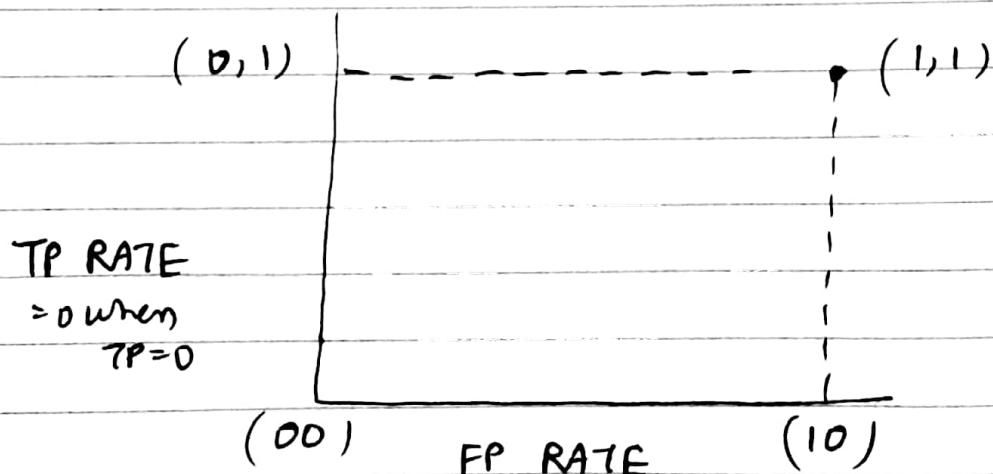


$$\text{Hit rate} = \frac{TP}{TP + FN}$$

$$\text{False Alarm Rate} = \frac{FP}{FP + TN}$$

$$\text{TP Rate} = \underbrace{\text{Recall of the class}}_{[0,1]} \rightarrow \text{sensitivity} = \frac{TP}{TP + FN}$$

$$\text{FP Rate} = 1 - \underbrace{\text{Recall of -ve class}}_{[0,1]} \rightarrow \text{specificity} = \frac{TN}{FP + TN} = \frac{FP}{FP + TN}$$



When does (00) occur

$\therefore$  model always predicts

-ve (either TN or FN)

(1D) :  $FP\ RATE = 1 \Rightarrow \frac{FP}{FP + TN} = 1$ ,  $TN = 0$

$TP\ RATE = 0 \Rightarrow \frac{TP}{TP + FN} = 0$ ,  $TP = 0$

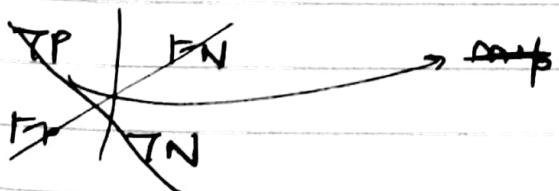
$\therefore$  no correct predictions

(1,1) :  $TP\ RATE = 1 \Rightarrow FN = 0$  & # all true sets are predicted true  
 $PP\ RATE = 1 \Rightarrow TN = 0$   $\therefore$  all are predicted true

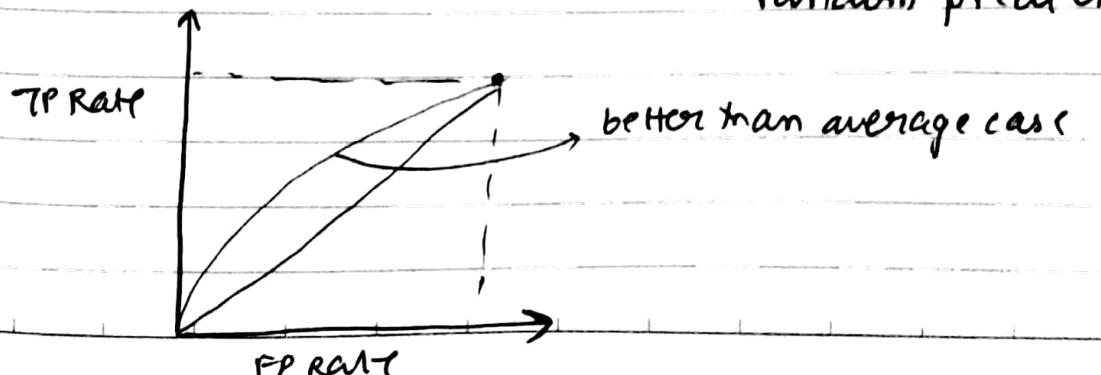
(0,1) :  $TP\ RATE = 1$   
 $FP\ RATE = 0$  ] ideal classifier

$TP\ RATE = 1 \Rightarrow FN = 0$  ] ideal case  
 $FP\ RATE = 0 \Rightarrow FP = 0$

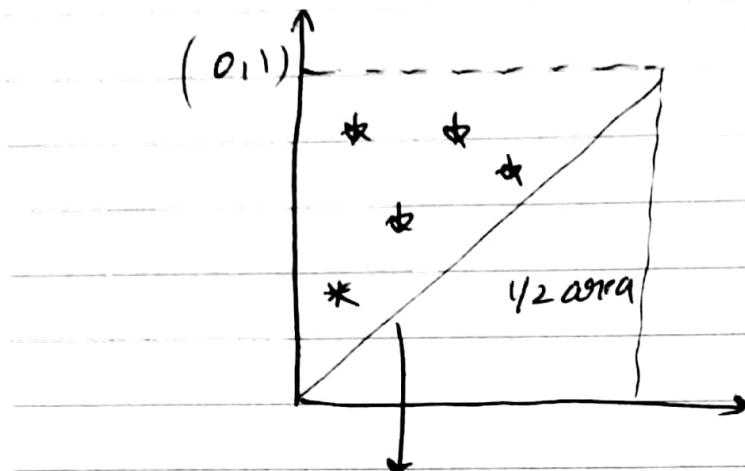
If classifier is predicted true class with 'p' prob. what is TPR? ] \*  
Total instances =  $m_+ + m_- = m$



ROC curve for TP rate = FP rate when classifier is making random predictions



Date :

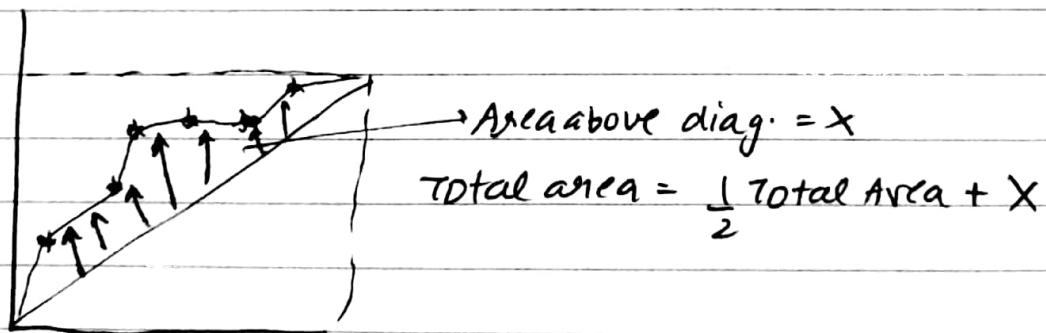


\* closest to (0,1) (ideal case) must be chosen as the most appropriate classifier.

diagonal is the classifier that predicts with fixed probability (average performance case)

AUC : Area under curve

The curve which gives most area → AUC is best classifier  
(done during model selection)



MLE : Max likelihood Estimat<sup>n</sup> (CH-4)



Regression (CH-4)



Reg. Trees

# PARAMETRIC ESTIMATION

## MAXIMUM LIKELIHOOD ESTIMATION (MLE)

↳ estimating parameters of some density function

Given  $m=7$  records following Normal Distribution

$$N(\mu, \sigma^2)$$

75
75
85
92
83
62
77

$$N_1(55, 5) \quad N_2(60, 7) \quad N_3(75, 10) \quad N_4(82, 15)$$

higher variance explains deviation ✓

values are closer around

as around 75 intuitively

Parametric Estimation is applied when we know the distribution

→ independently identically distributed variables

$X = \{x^t\}_{t=1}^N$  : training set with  $N$  instances

$x$  is a random variable

$P(X, \theta)$  : prob. density function

→  ${}^{2nd}$  training element

$$l(\theta | X) = p(x^1 | \theta) \cdot p(x^2 | \theta) \cdots p(x^m | \theta)$$

↓ likelihood  
of parameter  $\theta$  =  $\prod_{t=1}^m p(x^t | \theta)$  — (1)

→  $m^{th}$  train. element

goal : identify  $\theta$  which maximises (1)

Date :

$$\log (\ell(\theta | X)) = \log \left( \prod_{t=1}^n p(x^t | \theta) \right)$$

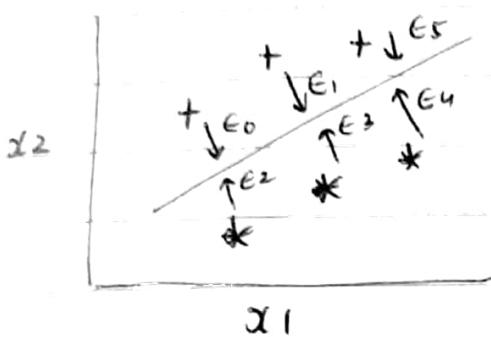
log likelihood  
(L)

$$= \sum_{t=1}^n \log p(x^t | \theta)$$

$$\log \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2} \right)$$

2 terms obtained where value is  
~~constant~~  
- ignored affected by  $\mu$  &  $\sigma$

## REGRESSION



$$\sum \epsilon_i = 0 \quad \text{--- (1)}$$

$$y^t = f(x^t) + \epsilon \quad \rightarrow \epsilon \sim N(0, \sigma^2) \quad \text{--- (2)}$$

using (1),  $\mu = 0$ ,  $\sigma = \sqrt{\frac{1}{2} \sum \epsilon^2 - (\bar{\epsilon})^2} = 0 \quad \therefore \sigma^2 \text{ is some non zero value}$

let  $g(\cdot)$  denote the approximate function

Assume we know the functional form of  $g(\cdot)$ , but still the parameters are to be estimated.

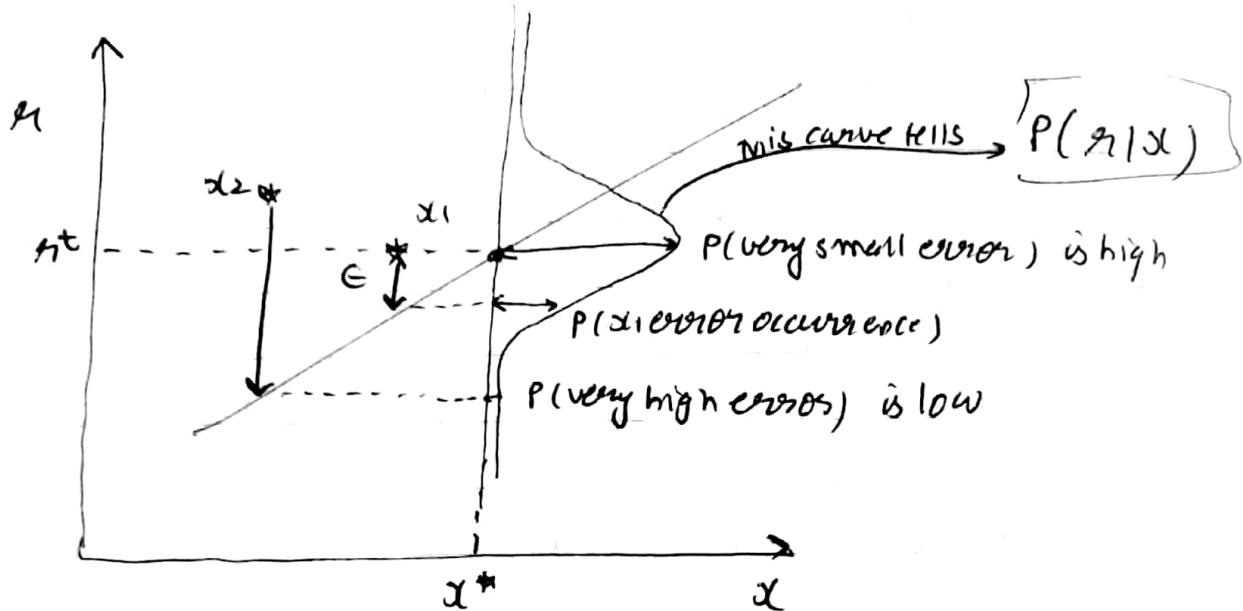
→ like we know the form of func " is normal distribution  
but we still need to know  $\mu \& \sigma^2$   
parameters

$g(x|\theta)$  approximates  $f(x^t)$

→ we don't know the parameters  $\theta$  which we need to know from (1)  
we are assuming  $g(\cdot)$  mimics  $f(\cdot)$

∴ in (2)  $\hat{y}^t = g(x|\theta) + \epsilon$

In case of a straight line :  $\hat{y}^t = w_0 + w_1 x^t$ ,  $\theta = (w_0, w_1)$   
for any new  $x^*$ , we substitute it in 1 to get  $\hat{y}^t$ , where  $(x^*, \hat{y}^t)$   
is a pt. on line



$P(r|s)$  follows  $N(g(s|\theta), \sigma^2)$

we need to find this using max. value estimation

$$L(\theta | s) = \log \left[ \prod_{t=1}^N P(s^t, r^t) \right]$$

same as  $P(r^t | s^t) \cdot P(s^t)$

$$= \sum_{t=1}^N \left[ \log (P(r^t | s^t)) + \log (P(s^t)) \right]$$

$$= \sum_{t=1}^N \underset{\text{I}}{\log (P(r^t | s^t))} + \sum_{t=1}^N \underset{\text{II}}{\log (P(s^t))}$$

since  $s^t$  is independent variable  $\therefore s^t$  is independent of  $\theta$   
 $\therefore \text{II is omitted}$

$$L(\theta | s) = \sum_{t=1}^N \log (P(r^t | s^t))$$

$$= \sum_t \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left( \frac{r^t - g(s^t|\theta)}{\sigma} \right)^2} \right]$$

$$= \sum_t \left[ \log \frac{1}{\sqrt{2\pi\sigma}} - \frac{1}{2\sigma^2} (y_i - g(x_i|\theta))^2 \right]$$

III                          IV

$$= \sum_t \log \frac{1}{\sqrt{2\pi\sigma}} - \frac{1}{2\sigma^2} \sum_t (y_i - g(x_i|\theta))^2$$

From the point of view of  $\theta$ ,  $\sigma$  is constant  $\therefore$  III is neglected

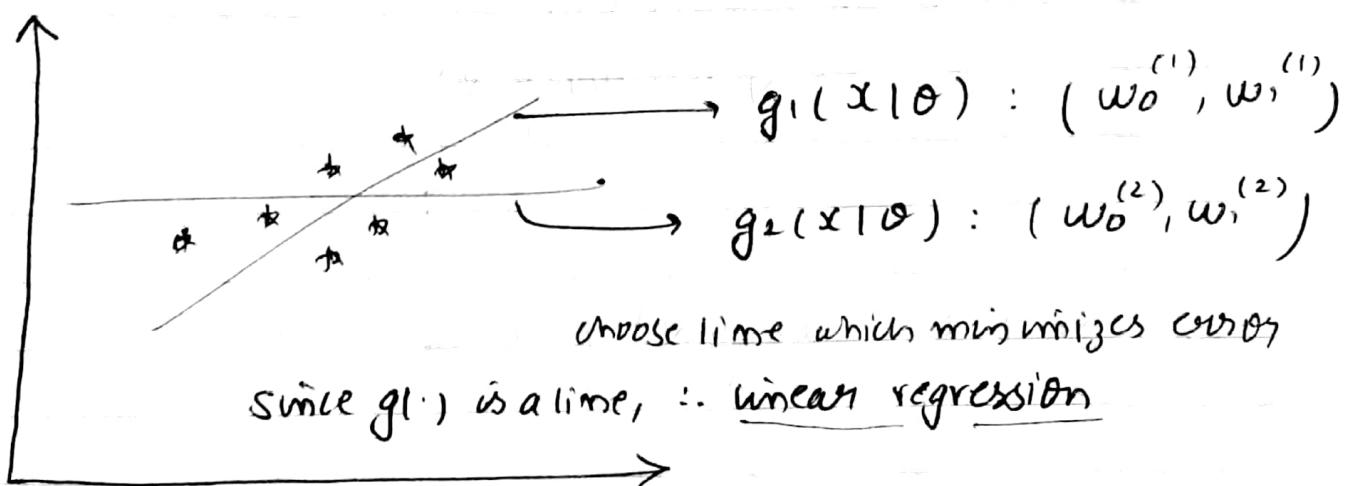
$$\therefore L(\theta|x) = \frac{-1}{2\sigma^2} \sum_t (y_i - g(x_i|\theta))^2$$

this is independent  $\therefore$  neglecting

$$L(\theta|x) = \sum_t (y_i - g(x_i|\theta))^2$$

we need to minimize this term → squared error

Likelihood  $\propto \frac{1}{\text{error}}$ , we choose the line which minimizes error



~~Sub - Sub~~

Date :

In linear regression,  $g(\cdot) = w_0 + w_1 x^t$

$$L(\theta | x) = \sum_t [r^t - g(x^t | \theta)]^2$$

$$= \sum_t (r^t - (w_0 + w_1 x^t))^2$$

x	r
:	:
:	:
:	:
:	:

We need to calculate  $\frac{\partial L(\theta | x)}{\partial w_0}$ ,  $\frac{\partial L(\theta | x)}{\partial w_1}$

$$\frac{\partial}{\partial w_0} [r^t - (w_0 + w_1 x^t)](-1)$$

→ ①

$$\frac{\partial}{\partial w_1} [r^t - (w_0 + w_1 x^t)](-x^t)$$

$$Eq. ① = 0$$

$$Eq. ② = 0$$

→ ②

$$-2 \left( \sum_t (r^t - (w_0 + w_1 x^t)) \right) = 0$$

$$-2 x^t \left( \sum_t (r^t - (w_0 + w_1 x^t)) \right) = 0$$

$$\Rightarrow \sum_t r^t = \sum_t w_0 + w_1 x^t$$

$$= N w_0 + w_1 \sum_t x^t$$

$$\sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t (x^t)^2$$

③

$$\sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t (x^t)^2 \quad → ④$$

A

X

B

$$\begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t (x^t)^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} \sum_t y^t \\ \sum_t n^t x^t \end{bmatrix}$$

we can calculate all  
these values from ..

$$\begin{array}{|c|c|} \hline x & n \\ \hline : & : \\ \hline : & : \\ \hline \end{array}$$

we find X using  $X = A^{-1}B$

### FOR POLYNOMIAL REGRESSION

$$g(x|\theta) = w_0 + w_1 x^1 + w_2 (x^1)^2 + \dots + w_k (x^1)^k$$

we need to estimate  $k+1$  parameters

$$L(\theta|x) = \sum (y_i - (w_0 + w_1 x^1 + \dots + w_k (x^1)^k))^2$$

we do  $\frac{\partial L(\theta|x)}{\partial w_0}, \frac{\partial L(\theta|x)}{\partial w_1}, \dots, \frac{\partial L(\theta|x)}{\partial w_k}$

$$A \quad X \quad = \quad B$$

$$\begin{bmatrix} (k+1) \times \\ \text{matrix} \end{bmatrix} \begin{bmatrix} x^0 \\ \vdots \\ x^k \end{bmatrix}_{(k+1) \times 1} \begin{bmatrix} \sum n^t \\ \sum n^t x^t \\ \sum n^t (x^t)^2 \\ \vdots \end{bmatrix}_{(k+1) \times 1}$$

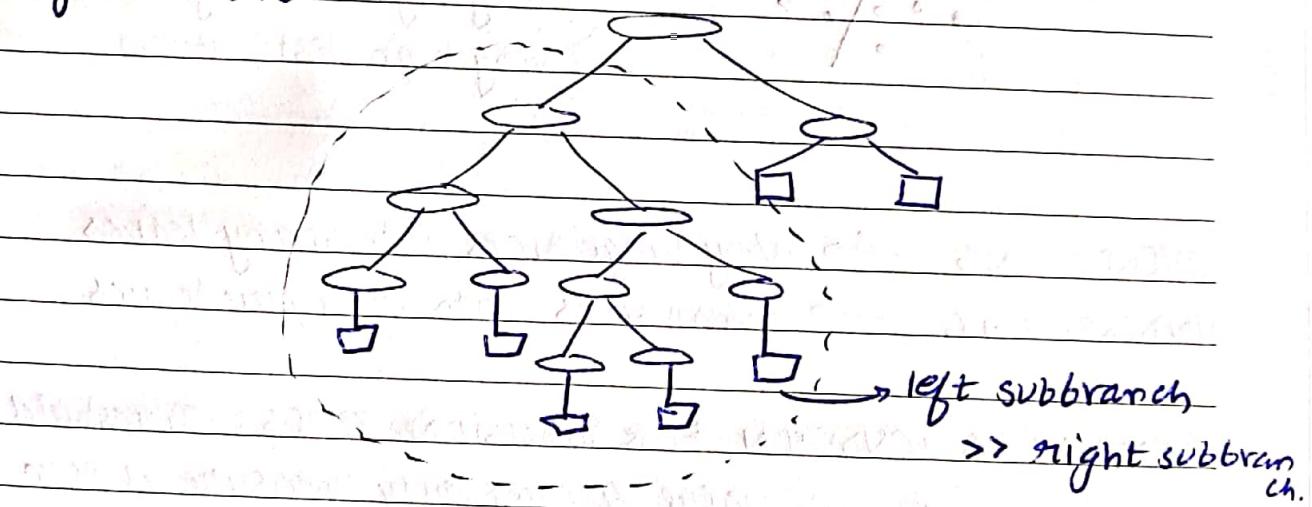
We can't find  $A^{-1}$  when  $|A|=0$  when

$$N \sum (x^t)^2 = (\sum x^t)^2$$

Regression trees are strictly binary in nature.

If threshold is small, the depth & 'bushyness' of tree increases.

The structure of tree can be highly imbalanced which depends over the nature of the data. We have no control over it apart from calculating impurity measure for selecting the splitting attribute.



### OVERFITTING v/s UNDERFITTING

we will choose  $h_1(x)$  over  $h_2(x)$

$h_1(x) \rightarrow ax^2 + bx + c \geq 0$  (eg: makes 12 errors)  
hypothesis (3 parameters) but still is easier to

$h_2(x) \rightarrow$  a piecewise curve (very difficult to describe)  
overfitted hypothesis (numerous parameters)

↳ complex of 2 reasons

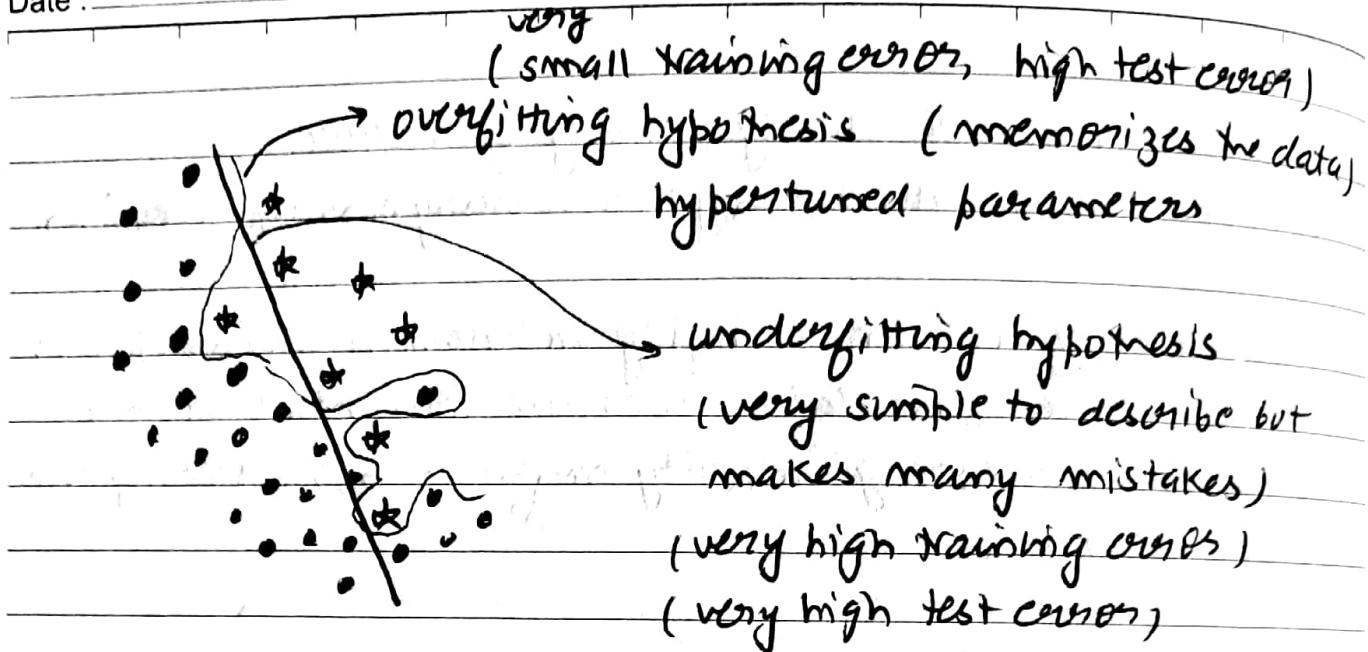
— too many ~~few~~ functions each with its parameters

simplest one is

better hypothesis

$h_3(x) \rightarrow mx + c \rightarrow$  simplest hypothesis but with way too many errors, so not always preferable to be selected. may underfit / too simple for data

Date : \_\_\_\_\_

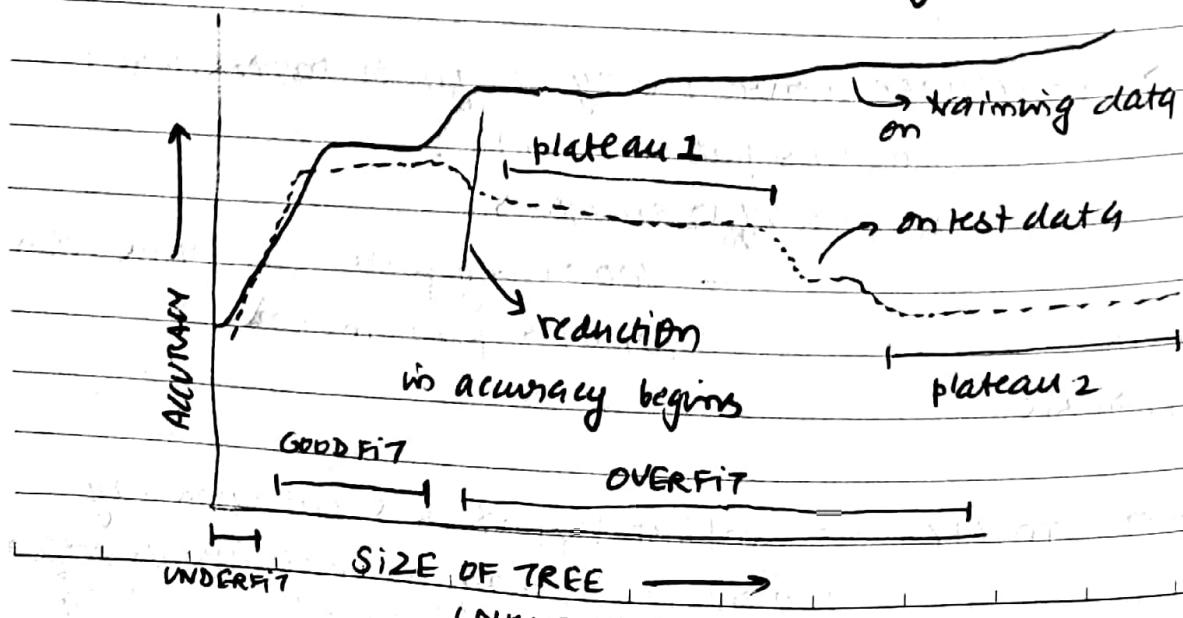


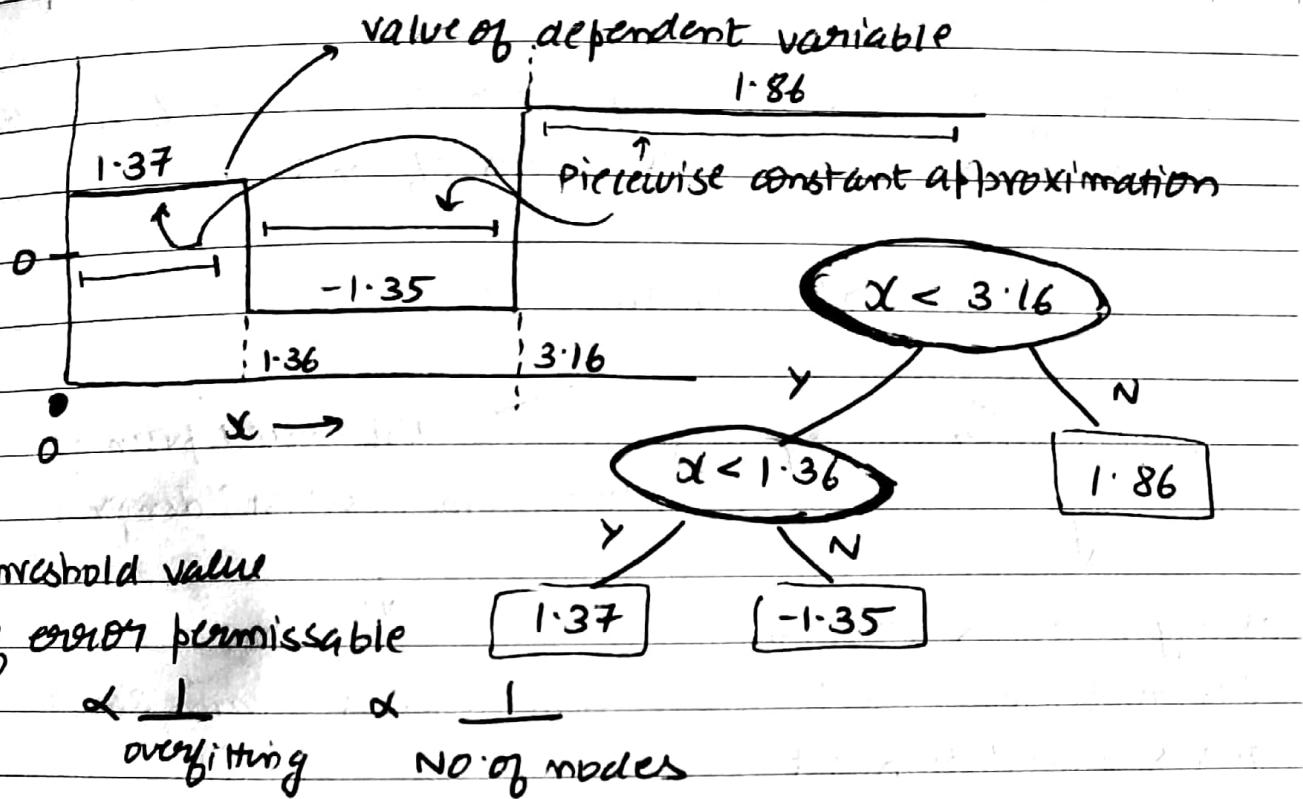
OVERRFITTING → Very large trees with many leaves

UNDERFITTING → Small trees with very few leaves.

TERMINATING CONDITION FOR DECISION TREES: Threshold value for impurity measure chosen

No. of nodes = 0 & accuracy = 65%.  $\Rightarrow$  simply guessing the class based on prior probability





### CAUSES OF OVERFITTING

→ Presence of noise

(in data / class labels)

↳ lack of representative samples

↳ excessive parameter tuning

If mean squared error (MSE)  $> \theta \rightarrow$  split else no split  
leaves with single instances  $\Rightarrow$  pure leaves

### PRUNING DECISION TREES

→ Remove subtrees for better generalisation

↳ prepruning, postpruning

→ Prepruning is faster, postpruning is more accurate  
(requires separate pruning set)

Date : \_\_\_\_\_

Post pruning involves replacing a subtree with a leaf provided the generalization error reduces because of this.

## POST PRUNING

Reduced error pruning



usually produces smallest version of tree., but it requires a validation set we are training on lesson data

MDL based pruning

(Minimum Descript<sup>n</sup> length)



formalizes Occam's Razor

tries to encode the model.

Each of K' hypotheses are encoded.

each with some length. we choose hypothesis with least length

## HOW TO ENCODE DEC. TREE?

We need to encode every discriminant function.

For 'd' attributes, we

need  $\log_2(d)$  bits

for encoding

We need to encode

hypothesis & errors

on test set

length of model

If entropy = 0, no. of bits reqd. to encode = 0

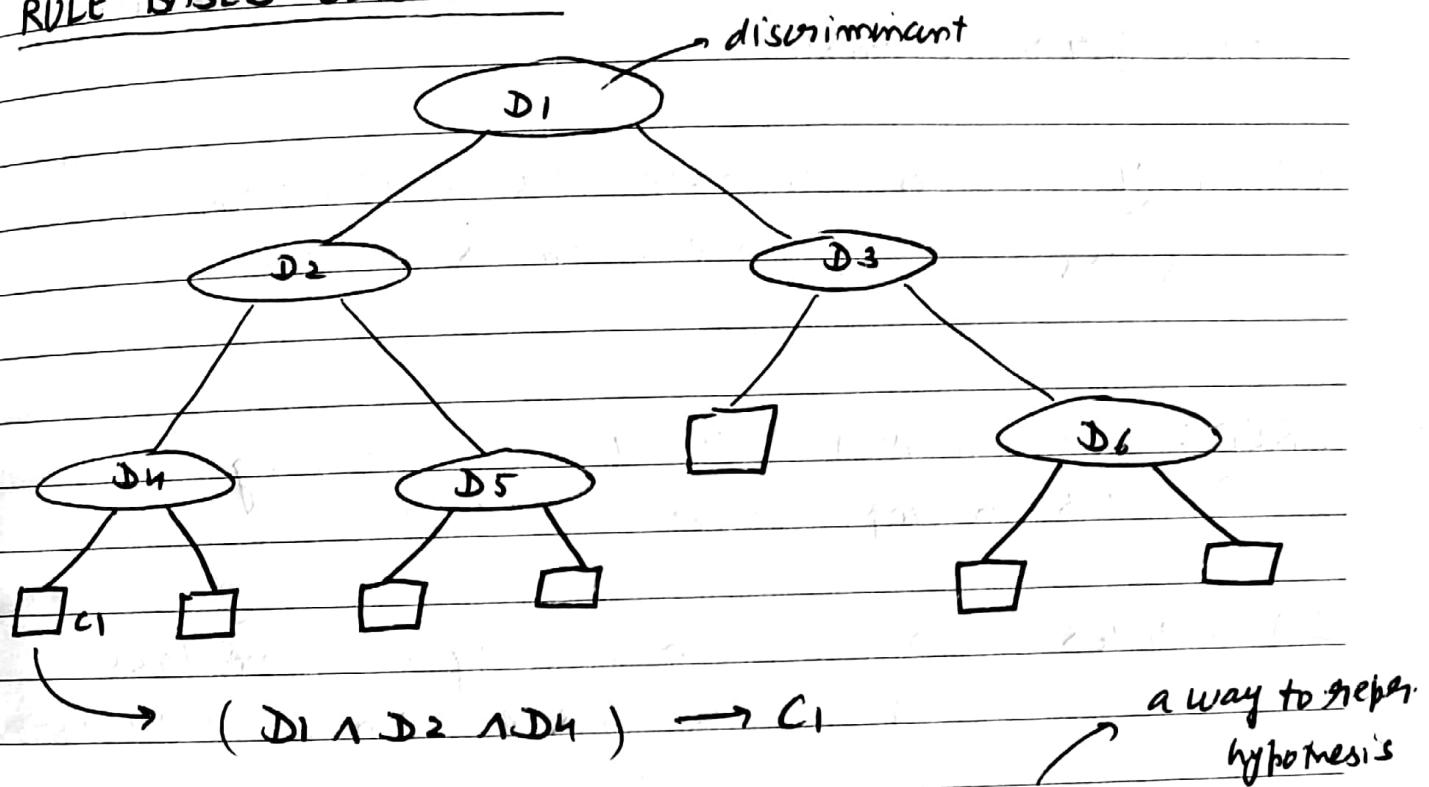
A D V

$\log_2 d$  3 X

$\log_2(d) + 3 + X$

bits.

## RULE-BASED CLASSIFIERS



If we only take all branches, then we get a RULESET

This is known as making rule based classifier from the tree

We can also make a rule based classifier from training set

Ruleset consists of rules which are conjunctions of terms.

Rule induction

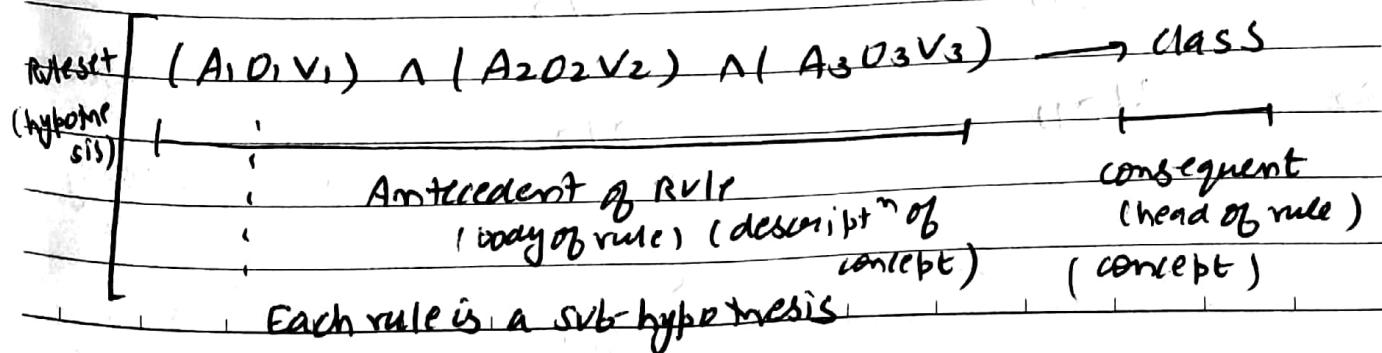
(depth-first)

one rule at a time

Tree induction

(breadth first)

adding conjunct to a rule until a quality criteria is matched



Date : \_\_\_\_\_

Rule covers an example if all terms of rule evaluate to true for the example.

SEQUENTIAL COVERING: Generate rules one at a time until all positive examples are covered.  
family of algos

IREP : Incremental Reduced Error Pruning

Ripper is an improvement over IREP

- 2 APPROACHES →
- 1) learn the body & assign class
  - 2) select the class & learn all the descriptions of that concept

There can be multiple descriptions for the same concept

Quality of classification rule :

1) COVERAGE

Fraction of records that satisfy the antecedent

2) ACCURACY

Fraction of records correctly classified

	COVERAGE	ACCURACY
R1	$3/20$	$3/3 \rightarrow$ of 3 records covered by this rule 3 are predicted correct
R2	$2/20$	$2/2$ as class predicted by rule is same as class of that record in training set

3 out of 20 records were covered by this rule

## COVERAGE

## ACCURACY Date:

R<sub>1</sub>/  
R<sub>2</sub>:

4/20  
3/20

High coverage  $\Rightarrow$  large no. of rules satisfy those combo of attribute values

We add more conditions to a rule with less accuracy ( & maybe high coverage ) , but it might lead to overfitting

We start with rules for the smallest class first

C<sub>1</sub>, C<sub>2</sub>, ... C<sub>K</sub> ( sorted in increasing order of no. of instances / priority )  
first rules for this  $\rightarrow$  most instances  
 $\therefore$  rules at last decreasing priority

$\hookrightarrow$  there are few records with such combo of attribute values

$\therefore$  fewer rules would be reqd. for such cases

for such would be very specific hypothesis

Rule 1 ( R<sub>1</sub> )

describes

this region

( x<sub>L</sub>, y<sub>L</sub> )

+	+	+	+
+	+	+	+
+	+	+	+
+	+	+	+

( x<sub>H</sub>, y<sub>H</sub> )

To develop a hypothesis for

this, we need two coordinates

to describe this region

Region covered by this hypothesis = ( x<sub>H</sub> - x<sub>L</sub> )  $\times$  ( y<sub>H</sub> - y<sub>L</sub> )

( x<sub>L</sub>  $\leq$  x  $\leq$  x<sub>H</sub> )  $\wedge$  ( y<sub>L</sub>  $\leq$  y  $\leq$  y<sub>H</sub> )  $\longrightarrow$  +

Date : \_\_\_\_\_

An empty rule  $\{\} \rightarrow C$  covers all the instances in data

$$C_1 = 20$$

$$C_2 = 40$$

$$C_3 = 60$$

$$\{\} \rightarrow C_2 \text{ coverage} = 100\%$$

$$\text{accuracy} = \frac{40}{120} = 33.33\%$$

$$\{\} \rightarrow C_1 \text{ coverage} = 100\%$$

$$\text{accuracy} = \frac{20}{120}$$

Adding more conditions to rule  $\propto$  complexity of constraints

start with empty rule &

$\propto$  1

Keep on adding conjuncts

instances covered

When to stop building a rule? based on quality

values of coverage & accuracy

Adding each new test reduces rule's coverage

stop adding if accuracy  $\rightarrow$  threshold & no significant  $\uparrow$

in accuracy

rippler alg

# classifier ensemble

Date: \_\_\_\_\_

## CHAPTER - 17

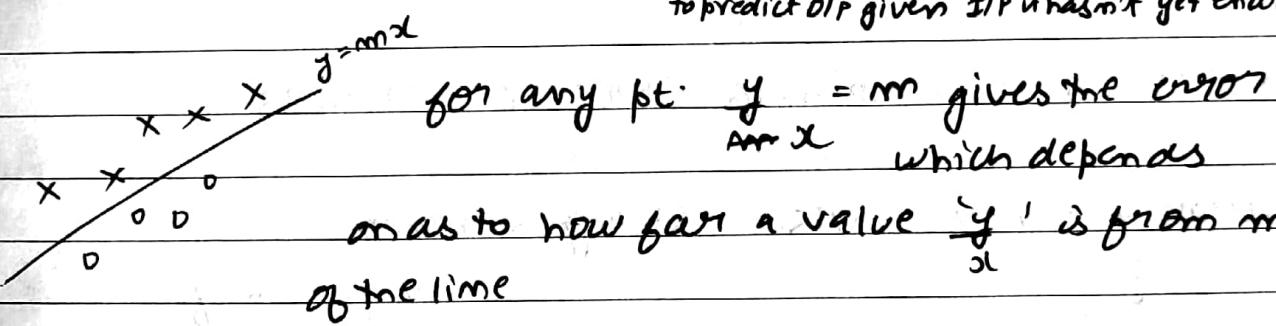
### COMBINING MULTIPLE LEARNERS

Combo of multiple classifiers

of same or different types

- hybrid model
- one model for one region R<sub>1</sub>, other for other region R<sub>2</sub>, ...
- "cross validation" accuracy is more reliable than test accuracy.
- every classifier has its own assumptions

- Assumptions of a model leads to Inductive Bias of model  
it is the set of assumption that learner uses to predict O/P given I/P it hasn't yet encountered.



$$E(\text{error}) = E(y - h(x))$$

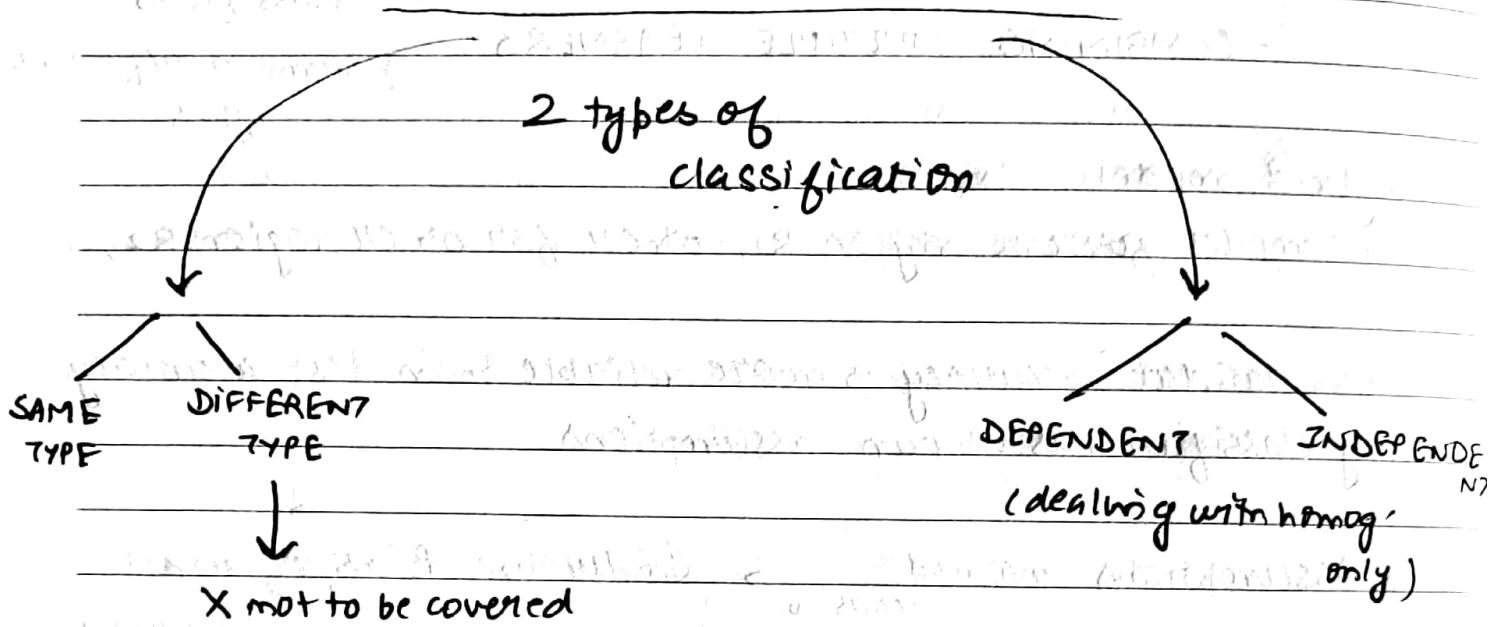
we combine multiple learners to reduce the overall errors  
error may significantly reduce if we choose some other classifiers,

- SAME TYPE
- we must —
- 1) generate multiple learners — DIFF. TYPE
  - 2) combine those learners

generating multiple learners, all of same type are called Homogeneous Ensembles. We can also have heterogeneous ensembles if all learners aren't of same type.

Date : \_\_\_\_\_

## GENERATING MULTIPLE LEARNERS



## MOTIVATION FOR COMBINING MULTIPLE LEARNERS

### i) Weak learners v/s strong learners

↓  
accuracy slightly > 50%  
which means its slightly better than random guess  
(for binary classifiers)

↳ high prob. of being correct

- Combining multiple weak learners can result into a strong learner
- Combining multiple strong learners is computationally expensive & accuracy gains are not significant & sometimes even degrades. It's generation of such strong learners is also expensive.

→ they predict correct at certain regions. So if we combine smartly we can increase accuracy

## 2) Unstable classifiers

vs

## Stable classifier

↓  
performance varies signif.

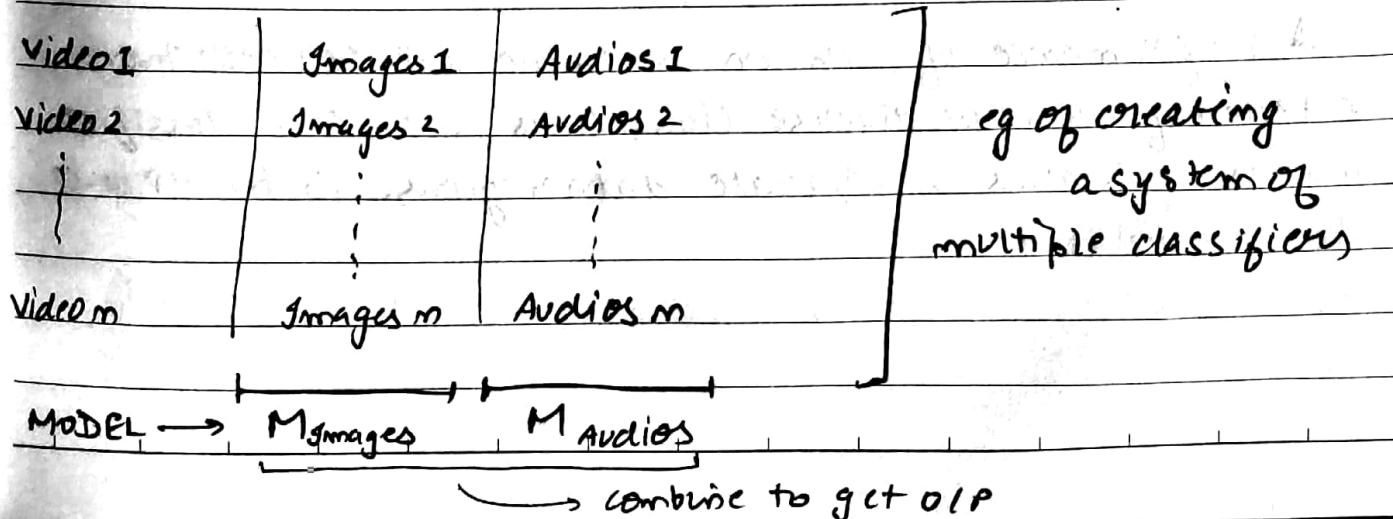
if training set is manipulated  
slightly eg: Decision Tree  
variance of model is very high  
on slight change in training set

↓  
performance won't  
affected significantly  
on slight change in  
training data

## HOW TO GENERATE MULTIPLE CLASSIFIERS (ISSUE 1)

- 1) we can generate different Algorithms (only for heterogeneous)
- 2) By changing the hyperparameters (for homogeneous)  
Re-training same NN (neural m/w) will result in a diff. NN  
if initializations are different.
- 3) Different Modalities / Representation

↳ one classifier for images ] classifying a video  
& one classifier for audio ]  
combine their results



Date : \_\_\_\_\_

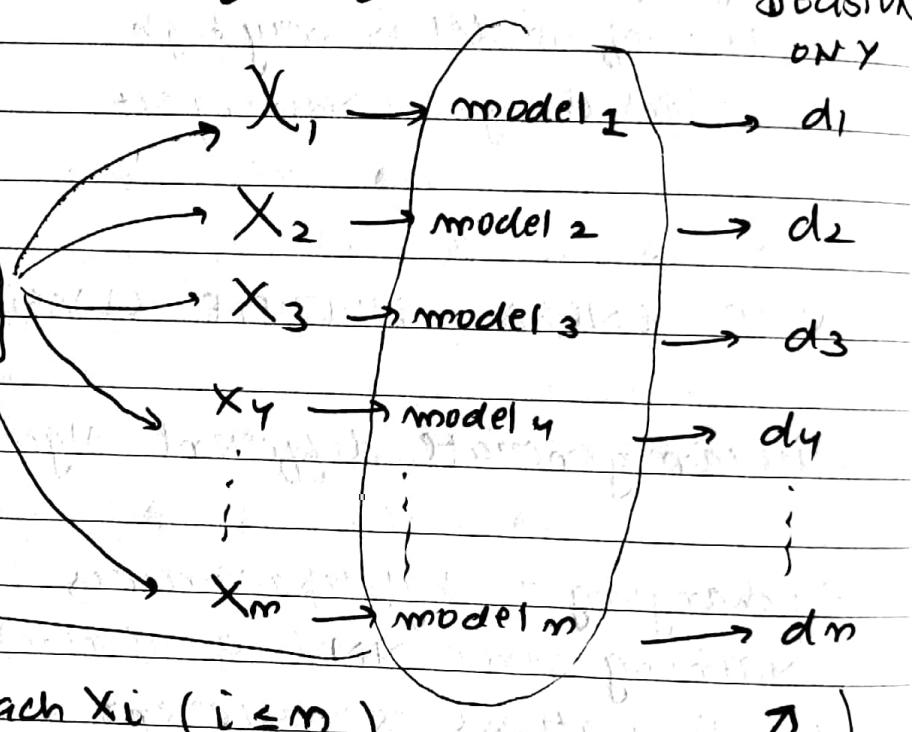
↳ Vary the training set (use different training set)

↳ how? : transform the same training set into separate versions & train the classifier for each.

X

original training set

METHOD TO TRANSFORM X



A test set goes to each  $X_i$  ( $i \leq m$ )  
(Y)

COMBINER:  $\hat{y} \leftarrow f(d_1, d_2, \dots, d_m)$

→ performance of such an ensemble will be better than expected if we have diverse classifiers i.e. each classifier makes mistakes in separate data regions  $\therefore$  the outputs will be uncorrelated.

## HOW TO COMBINE MULTIPLE LEARNERS

### Global Methods

(combine O/P of all classifiers)

eg: voting, stacking, ~~gating~~  
gating

models

$M_1 \ M_2 \ \dots \ M_L$

predictions

$\rightarrow d_1 \ d_2 \ \dots \ d_L$

voting gives weight to each classifier

$$w_1 \ w_2 \ \dots \ w_L \ \text{s.t.} \ (0 \leq w_i \leq 1 \ \& \ \sum w_i = 1)$$

combine them

$$\hat{y} = \sum w_i d_i \ (\text{weighted sum of O/P})$$

In case of classification:  $\hat{y} = \sum w_i$  (for binary classification)  
 if  $\sum w_i > 0.5$  then class 1  
 else class 0

CASE 1: odd no. of classifiers

$M_1 \ M_2 \ M_3 \ M_4 \ M_5$

$d_1$	$d_2$	$d_3$	$d_4$	$d_5$
(0)	(1)	(0)	(0)	(1)

$w_1 \ w_2 \ w_3 \ w_4 \ w_5$

$$\hat{y} = \sum w_i d_i = w_2 d_2 + w_5 d_5$$

Date : \_\_\_\_\_

### CASE 2 : even no. of classifiers

M<sub>1</sub> M<sub>2</sub> M<sub>3</sub> M<sub>4</sub>

d<sub>1</sub> d<sub>2</sub> d<sub>3</sub> d<sub>4</sub>

(0) (0) (1) (1) → 2 say '1' & 2 say '0'

w<sub>1</sub> w<sub>2</sub> w<sub>3</sub> w<sub>4</sub> in this case weights

are  
into play

$$\hat{y} = \sum w_i d_i = w_3 d_3 + w_4 d_4$$

$$= w_1 d_1 + w_2 d_2 \text{ if } w_1 = w_2 = w_3 = w_4$$

### CASE 3 : multiple classes (we must use indicator functions)

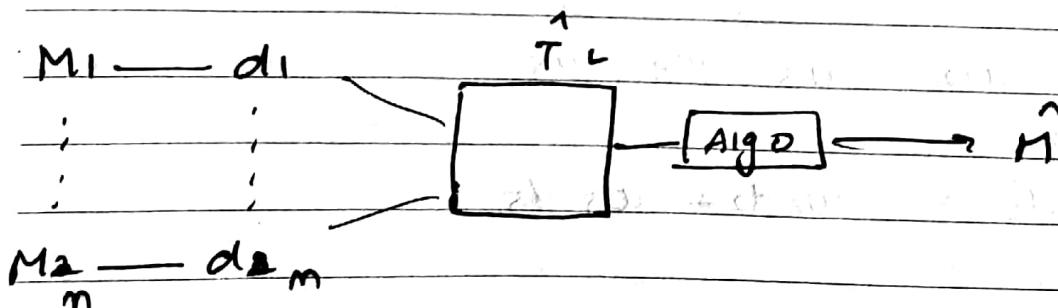
d<sub>1</sub> d<sub>2</sub> d<sub>3</sub> d<sub>4</sub> d<sub>5</sub> d<sub>6</sub> d<sub>7</sub>

(1) (2) (2) (3) (1) (3) (1)

$$y_i^j = \sum_j I(d_j = c_i) \rightarrow \text{no. of models predicting class } i$$

Final class =  $\arg\max_i y_i^j$  assuming all weights are same

### METALEARNING (Learning the learnings)

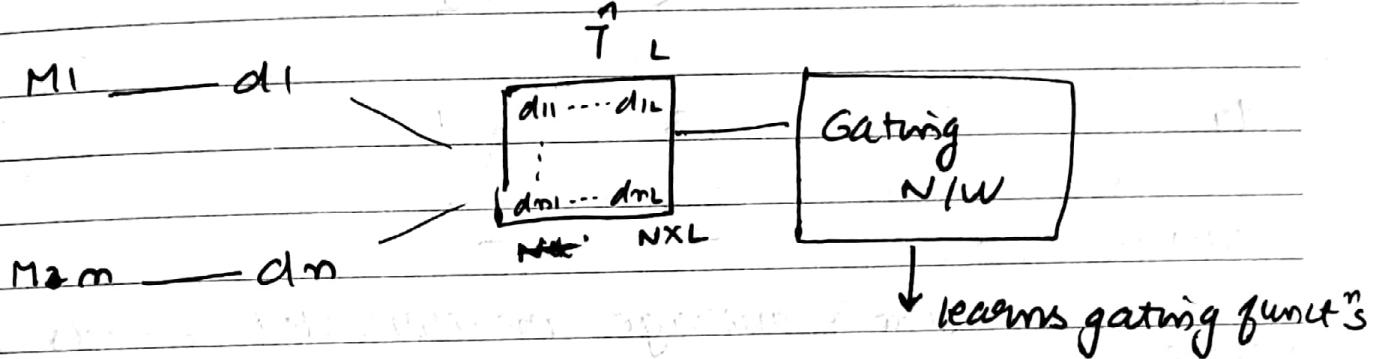


→ can maintain data privacy.

Date : \_\_\_\_\_

Stacking & Gating are also Global Methods

a metalearning method



Voting & Gating functions

are diff. as weights

in Gating functions are learnt

from the network training set

$\sum_j w_{ij} d_j$  ← each assigns probabilities to models  
↓ diff. from weights in voting (just like weights)

CASCADING : every model is assigned some confidence level (threshold) & when instance is to be predicted representation of instance is checked with each model. If confidence of a model for that instance is low, we proceed to next model. It's a local prediction as only one model is going to predict the class.

Combining Methods for Global Methods : Average, Min, Median,

↓  
Product  
in Regression

depending  
on  
applicat")

If algo is predicting Prob. of class, then average can be used

Date : \_\_\_\_\_

	$C_1$	$C_2$	$C_3$
$d_1$	0.2	0.5	0.3
$d_2$	0.7	0.2	0.1
$d_3$	0.3	0.5	0.2
$d_4$	0.6	0.3	0.1
SUM	-	-	-

sum & average are correlated methods  
as sum & average.

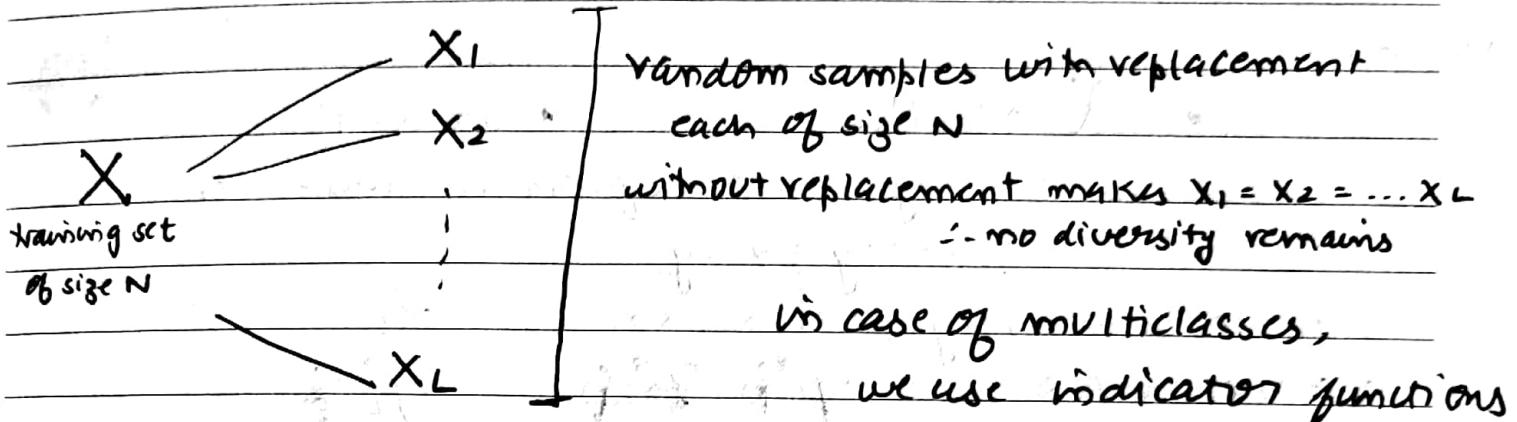
Diff algo have diff inductive bias  $\therefore$  of which they might focus on diff data regions  $\therefore$  combining such algs can result in better decision boundaries.

Gain ~~gain~~ & Diversity of classifiers

NFL (NO FREE LUNCH) THEOREM : There is no algo that is always the most accurate.

## BAGGING (ENSEMBLE METHOD) → Bootstrap Aggregation

random sampling with replacement



$$X_1 - M_1$$

$$X_2 - M_2$$

⋮

$$X_L - M_L$$

$$\sum w_j d_j \left( \frac{1}{L} \forall j \right)$$

we take average as  $w_j = \frac{1}{L} \forall j$

How much repetition can we expect?

$$\left(1 - \frac{1}{N}\right)^N : \text{prob. of an instance not coming in } X$$

$$\approx \frac{1}{e} \approx 0.632 \quad 0.37$$

There is 37% chance that an instance not coming in  $X$

There are 63% unique elements in each training set

Date : \_\_\_\_\_

$X = \langle x, y \rangle$ ,  $x \in \mathbb{R}^d$ ,  $y \in \mathbb{R}^c$  → class set

training set  
↓  
numeric classes  
↓  
in case of discrete class values

### ASSUMPTIONS

- $X \sim P_{xy}$
  - $x_0 \sim P_{xy}$
- Training set comes from the distribution defined by  $x, y$

1) Create L training sets  $x_1, x_2, \dots, x_L$

2) Create models  $\rightarrow h_1(\cdot), h_2(\cdot), \dots, h_L(\cdot)$

assumptions:  
1) weak learners  
2) countable learners

high variance

$h_1(\cdot), h_2(\cdot), \dots, h_L(\cdot)$

for  $x_0 \rightarrow$  if we use hypothesis  $h_1 \rightarrow \langle x_0, h_1(x_0) \rangle$

" " ..  $h_2 \rightarrow \langle x_0, h_2(x_0) \rangle$

In Bagging, we take average of  $h_i(x_0)$ , i.e. expected value

of  $h$  over L classifiers  $\rightarrow E(h(x_0)) = h_B(x_0)$

We want to find out how far (on average) is each prediction from  $h_B(x_0)$  or  $E(h(x_0))$

bagged estimate

$$= E((h(x_0) - E(h(x_0)))^2)$$

= variance for  $x_0$

Date : \_\_\_\_\_

e.g. of unstable classifiers: decision trees, KNN, ANN  
A linear classifier is a stable classifier

Stable classifiers have high bias -

→ the error in prediction, due to the assumptions taken.

$$x, y \in \mathbb{R}$$

$$\text{NN}(0, \sigma^2) \rightarrow \text{all terms cancel out} \therefore \mu = 0$$

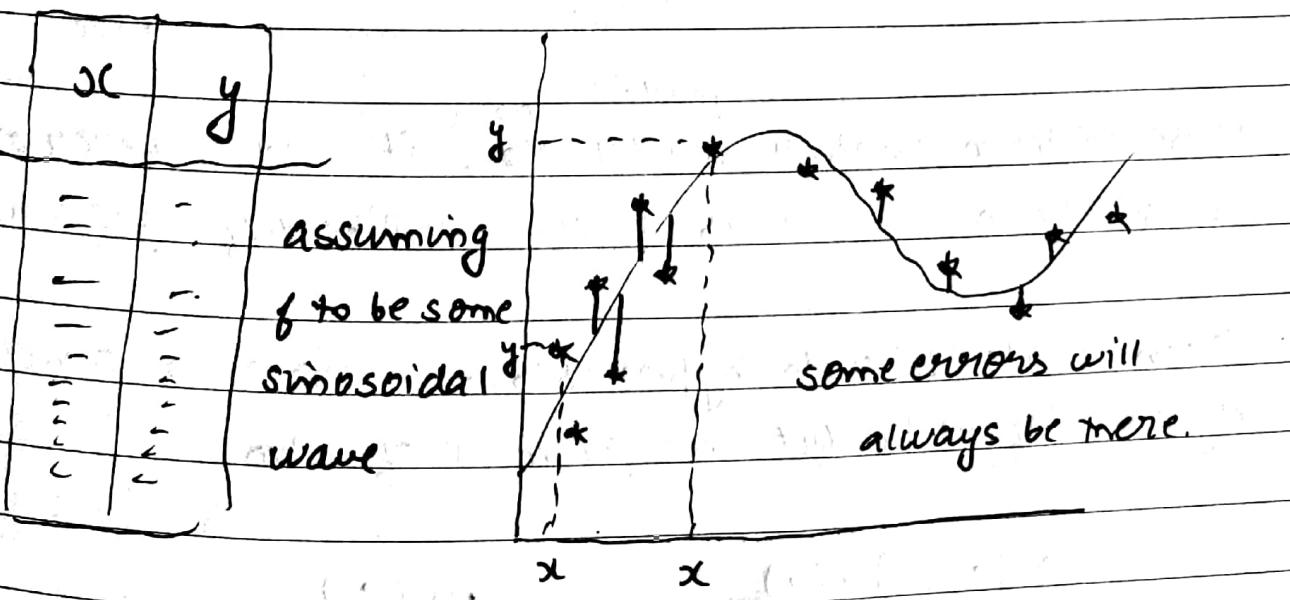
$$y = f(x) + \epsilon$$

dependent variable / predicted

independent variable / predictor

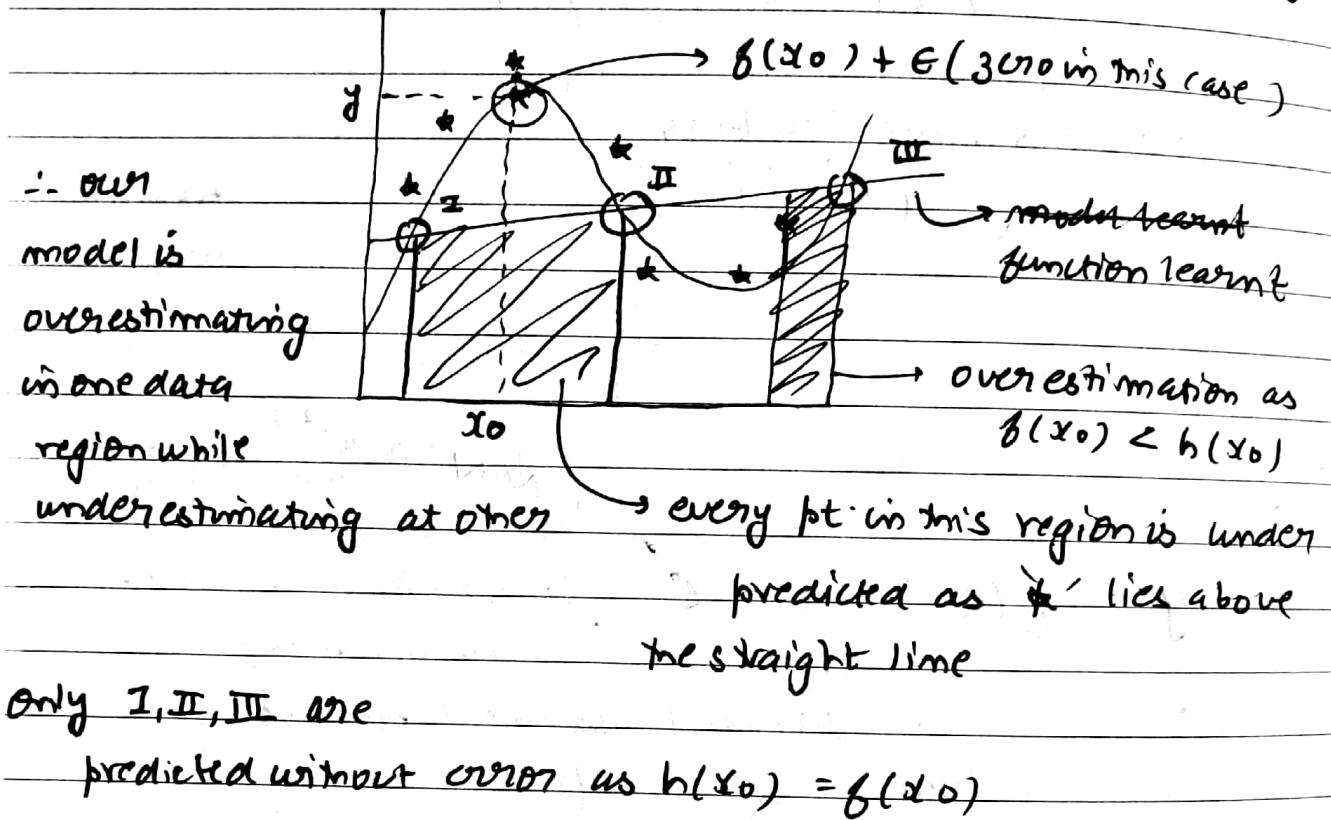
error / noise associated with each record

For simulated data, we know  $f(x)$ , but it isn't so in real life



Date : \_\_\_\_\_

Suppose we learn our model with these 2 parameters  $\alpha, \gamma$



$$\text{Bias for } x_0 \text{ (of the model)} = h(x_0) - f(x_0)$$

predicted                      true

Suppose we train another model  $h_2$ , it could be better or worse than  $h_1$ .  $\therefore$  Bias  $h_1(x_0)$  will be diff. from  $h_2(x_0)$

We can create 'L' such hypotheses each with their own error regions. If we predict  $x_0$  from all of them:

$$h_1(\cdot), h_2(\cdot), \dots, h_L(\cdot)$$



$$h_1(x_0), h_2(x_0), \dots, h_L(x_0)$$



$$E_4(h(x_0)) = \frac{1}{L} \sum h_i(x_0)$$

The variance would lie somewhere in mid of all predicted values for  $x_0$ .

$$E[(h(x_0) - E(h(x_0)))^2]$$

Bias in Decision trees comes

because of assumptions

(attribute selection measures,  
depth of tree).

Bias, Variance, Noise

(B) (V) (N)

$$\text{Error} = B^2 + V + N$$

$$\text{Squared error} = (y_0 - h(x_0))^2$$

Mean squared error for  $x_0$

$$= E[(y_0 - h(x_0))^2]$$

$$= E(y_0^2 - 2h(x_0)y_0 + h(x_0)^2)$$

$$= E(y_0^2) - 2E(h(x_0))E(y_0)$$

$$+ E(h(x_0)^2)$$

— (2)

We know

$$E((x - \bar{x})^2)$$

$$= E(x^2) - [E(x)]^2$$

$$\Rightarrow E(x^2) = E((x - \bar{x})^2)$$

$$+ [E(x)]^2$$

— (1)

Date : \_\_\_\_\_

using ① w.r.t ② to compute  $E(y_0^2)$

$$A \rightarrow E(y_0^2) = E((y_0 - E(y_0))^2) + [E(y_0)]^2$$

$$C \rightarrow E(h(x_0)) = E((h(x_0) - E(h(x_0)))^2) + [E(h(x_0))]^2$$

Putting these A & C in ②

$$= E((y_0 - E(y_0))^2) + [E(y_0)]^2 - 2E(h(x_0))E(y_0)$$

$$+ E((h(x_0) - E(h(x_0)))^2) + [E(h(x_0))]^2$$

↳ predicted      ↳ average prediction  
of weak learners

$$= E((y_0 - f(x_0))^2)$$

error

$$+ E((h(x_0) - E(h(x_0)))^2) + E(f(x_0))^2 - 2E(h(x_0))E(f(x_0))$$

variance      ↳ varying

$$= h(x_0)^2 - 2h(x_0)f(x_0) + f(x_0)^2$$

$$= E((f(x_0) - h(x_0))^2)$$

Bias  
(squared of bias)

Date :

- 1) Find out of bag accuracy
- 2) select unique out of bag instances, & for each find bias & variance

1000 pts. take randomly. learn a linear regression model from this dataset.

Create L sets :  $X_1 \dots X_L$  L Models.  
 $M_1 \dots M_L$

Bagging only reduces the variance only thereby reducing overall error  $E = B^2 + V + N$

Variance of hypothesis for  $x_0 = E((h(x_0) - E_L(h(x_0)))^2)$   
Overall variance =  $E_x [E((h(x_0) - E_L(h(x_0)))^2)]$

Date: for classification  
 Adaptive Boosting Algorithm (staged classifier)  
 works stage wise  
 incremental approach to making ensemble

ADABOOST (proposed in 1990)

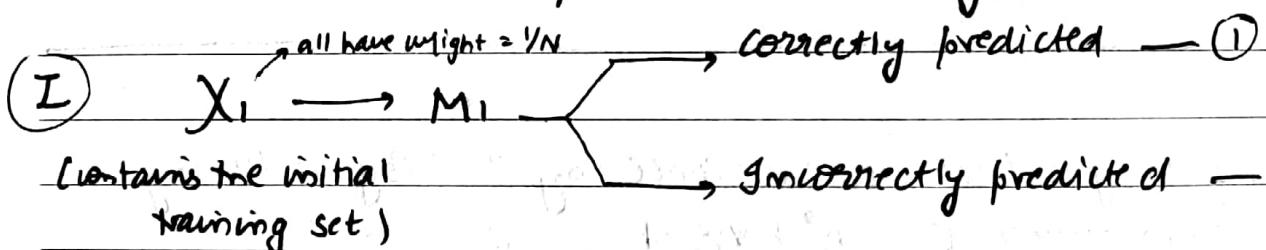
↳ all its claims are provable mathematically

↳ much better ensemble method than bagging

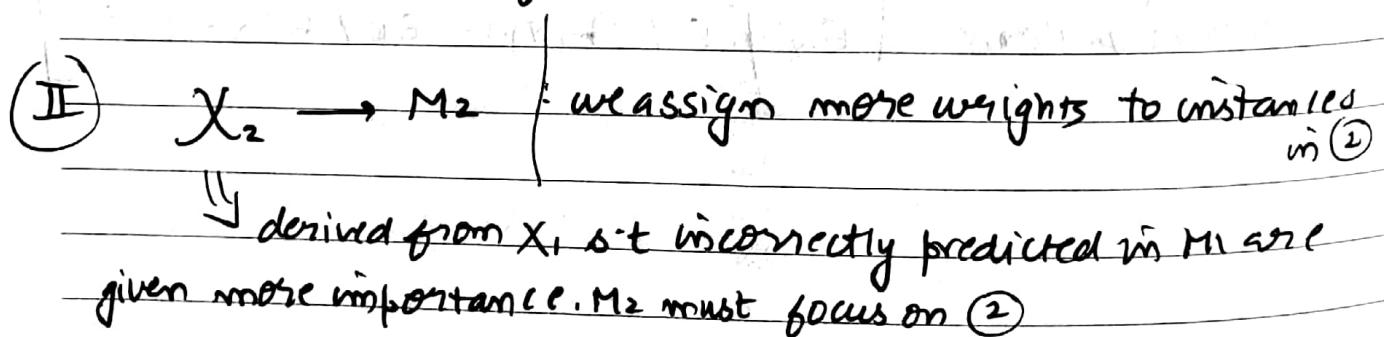
↳ boosts performance of weak classifiers

Boosting: general philosophy based on idea that

if you create a model  $M_1$  & predict the instances in training set, then  $M_1$  will predict some correctly while some will be predicted incorrectly.



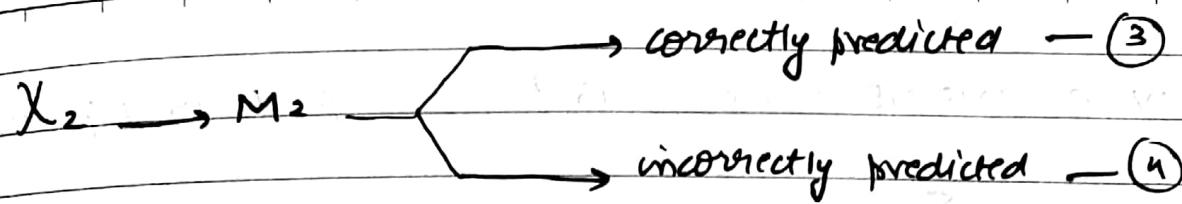
Now, we will create another  $M_2$  which focuses on incorrectly predicted instances of  $M_1$



$X_1$  = training set  $X$  from which sampling is done without replacement with  $P$  (an instance is selected) =  $\frac{1}{N}$

$$\frac{1}{N} = \text{Weight of } i^{\text{th}} \text{ instance in } 0^{\text{th}} \text{ iteration} = w_i, \sum w_i = 1$$

incorrectly classified in ② will have weight  $> \frac{1}{N}$  & correctly classified in ② will have weight  $< \frac{1}{N}$



$$w_i^{(1)} = \frac{1}{N}, \quad w_i^{(2)}$$

$\begin{cases} > \frac{1}{N} & \text{if } x_i \text{ is wrongly pred.} \\ < \frac{1}{N} & \text{if } x_i \text{ is correctly pred.} \end{cases}$

$w_i^{(3)} \begin{cases} > w_i^{(2)} & \text{if wrongly} \\ < w_i^{(2)} & \text{if correctly.} \end{cases}$

We want to stop when:

- 1) ' $L$ ' classifiers (iterations) have been completed
- 2) If accuracy  $< 50\%$ . i.e. random guess then stop (in case of binary classifier)

Working on the correctly predicted instances from next iterations so that size of set remains appreciable

$$w_i^l = f(w_i^{l-1}, \text{performance measure})$$

After (2) all incorrectly predicted will have same weight

Bagging had its classifiers working parallelly. AdaBoost works sequentially; as  $x_i$  can't be made without knowing  $x_{i-1}$ .

All classifiers are independent in bagging. A. Each training set is independent on performance of model created by prev training set in case of AdaBoost.

Date : \_\_\_\_\_

HOW TO COMBINE PREDICTIONS?  $\rightarrow$  weighted sum of predictions

$$h_{\text{ADABOOST}}(x) = \sum \alpha_i h_i(x) \quad (5)$$

no. of classifiers  $\leftarrow$   $\downarrow$  weight of  $i^{\text{th}}$  classifiers  $\rightarrow$  prediction of  $i^{\text{th}}$  classifier

Adaboost was initially used for binary classification.  $\checkmark$   
with prediction classes = {1, -1}

-1 taken instead of 0  
to denote wrong prediction for mathematical convenience  
Bagging was initially for regression purposes.

$$(5) \Rightarrow h_A(x) = \operatorname{sign} \left( \sum \alpha_i h_i(x) \right)$$

HOW TO ASSIGN WEIGHTS  $\alpha_i$ ?

$$\alpha_i = g(\text{performance})$$

DECISION CRITERIAS

1) Adjust weights (generating training set / Models)

2) Determine weights of classifiers

LOSS FUNCTION

$$\text{ERROR} = \sum_{x \in X} I(y \neq h(x)) \quad \text{--- (6)}$$

(0/1 Loss)

Loss Function followed by AdaBoost is Exponential Loss

$$\text{loss}_{\text{exp}}(x) = e^{-y h(x)} \text{ where } y = +1, -1$$

For correct prediction:

$$\begin{array}{c|cc|c} y & h(x) & y h(x) \\ \hline 1 & 1 & = 1 \\ -1 & -1 & = -1 \end{array}$$

$$\Rightarrow \text{loss} = e^{-1} = \frac{1}{e}$$

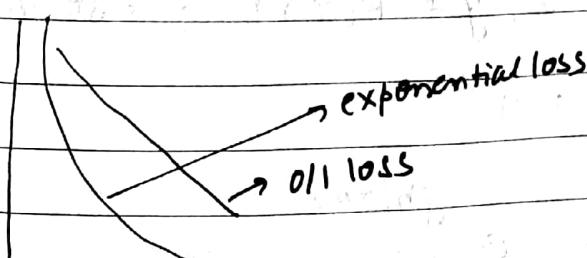
For wrong prediction

$$\begin{array}{c|cc|c} y & h(x) & y h(x) \\ \hline 1 & -1 & -1 \\ -1 & 1 & -1 \end{array}$$

$$\Rightarrow \text{loss} = e$$

$$\text{loss for classifier} = \text{loss}_{\text{exp}}(h(x)) = \sum_x e^{-y h(x)}$$

Exponential loss function will reduce faster than (6)



Date : \_\_\_\_\_

→ uses exponential loss

ADABOOST → staged, incremental, dependent classifiers, combines using weighted additive function

Sigm

$$y_i = \left( \sum_l \alpha_l h(x_i) \right) \rightarrow \because \text{class can be } +1 \text{ or } -1$$

$y = h(x)$  → correct prediction loss = 0

$y \neq h(x)$  → incorrect prediction loss = 1

$$\cancel{\text{loss}} = \sum_{x \in X} I(y \neq h(x))$$

$\sum_{x \in X} I(y = h(x))$  : correct prediction

$\left[ \sum_{x \in X} I(y \neq h(x)) \right]$  : training errors

$\sum_{x \in T} I(y \neq h(x))$  : test errors

$|X| - \sum I(y \neq h(x))$  : training accuracy.

$$|X|$$

$|T| - \sum_{x \in T} I(y \neq h(x))$  : test accuracy.

$$|T|$$

• Loss<sub>exp</sub> ( $x_0$ ) =  $e^{-y h(x_0)}$

$$\bullet \text{LOSS}_{\text{exp}}(x) = \sum_{x \in X} e^{-y h(x)}$$

$$\text{Risk or expected loss} = E(\text{LOSS}_{\text{exp}}) = E(e^{-y h(x)})$$

$$\text{if } y = h(x), y \cdot h(x) = 1, \text{ LOSS} = e^{-1} \quad \textcircled{1}$$

$$\text{if error, i.e. } y \neq h(x), y \cdot h(x) = -1, \text{ LOSS} = e^{-1} \quad \textcircled{2}$$

∴ using  $\textcircled{1}$  &  $\textcircled{2}$ ,

$$E(e^{-y h(x)})$$

$\textcircled{I}$

$\textcircled{II}$

$$\sum_{x \in X} e^{-y h(x)} = \sum_{\text{correct}} e^{-y h(x)} + \sum_{\text{wrong}} e^{-y h(x)}$$

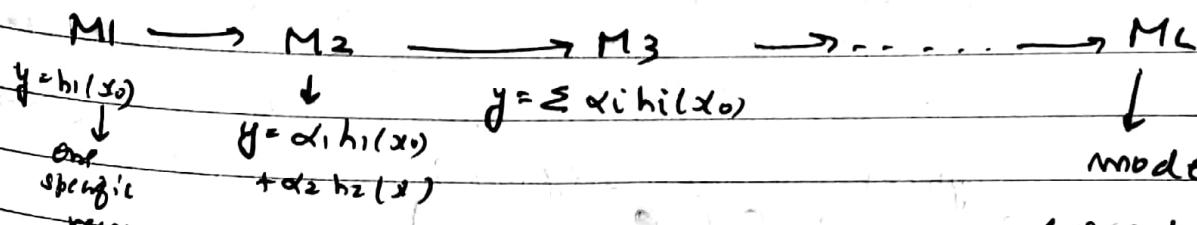
we must reduce  $\textcircled{II}$  as it contributes more

Ideally,  $\textcircled{II} = 0$

least value =  $m e^{-1}$

improving

Adaboost was originally designed for classification trees  
Bagging was in "regression trees"



Date : \_\_\_\_\_

## HOW TO ASSIGN WEIGHTS TO HYPOTHESIS?

- upto some ( $L-1$ ), we have the hypothesis  $H(x) = \sum_{i=1}^{L-1} \alpha_i h_i(x)$   
→ suppose  $h(x_0)$  is hypothesis at stage  $L$  &  
 $\alpha$  is the weight of  $h(x_0)$

we have to find  $\alpha$ .

Now,  $y = H(x_0) + \alpha h(x_0) \rightarrow$  prediction for  $y_0$  at stage  $L$

$\downarrow$   
hypothesis of  $x_0$   
at  $L^{\text{th}}$  stage

loss ( $x_0$ ) at  $L^{\text{th}}$  stage =  $e^{-y(H(x_0) + \alpha h(x_0))}$

loss ( $x$ ) =  ~~$\frac{1}{2} e^{-y(H(x) + \alpha h(x))^2}$~~  =  $E_{xy}(e^{-y(H(x) + \alpha h(x))})$

at  $L^{\text{th}}$  stage

$$= E_{xy}(e^{-yH(x)} \cdot e^{-\alpha y h(x)})$$

loss at  $L-1$  stages  
we already have it  
so let it be constant

this is expected loss  
of hypothesis

$$= e^{-yH(x)} \cdot E_{xy}(e^{-\alpha y h(x)})$$

minimum loss

loss ( $x$ ) will be minimum if for every instance  
loss is minimum, in this case if  $e^{-\alpha y h(x)}$  is min.

## IMPORTANT NOTATIONS

$E_{xy} \left[ (y - h(x))^2 \right]$  : Mean squared training errors

$e^{-yh(x)}$  : exponential loss, pointwise loss

$\mathbb{E}_{xy} (e^{-yh(x)})$  : expected loss over entire training set.

$$h(x) = \sum_{i=1}^L \alpha_i h_i(x)$$

$$e^{-yH(x)} \cdot E_{xy} (e^{-\alpha y h(x)})$$

This is constant & we are assuming to have minimised it at  $(L-1)^{\text{th}}$  stage

→ minimize  $E_{xy} (e^{-\alpha y h(x)})$  (1)  
we already have  $y$  from  
training set,  $h(x_0)$  is  
random

$$\text{we know that } E(x) = \sum x_i p(x_i)$$

(1) is proceeded to have 2 cases, when  $y = h(x)$  &  $y \neq h(x)$

$$\therefore (1) = \left[ e^{-\alpha} * P(y = h(x_0)) + e^{\alpha} * P(y \neq h(x_0)) \right]$$

$\xrightarrow{\text{correct prediction}}$        $\xrightarrow{\text{incorrect prediction}}$

$y \cdot h(x) = 1+1=1$       as  $y \cdot h(x) = -1$

$y \cdot h(x)$  loss

0	0	$e^0$
0	1	$e^0$
1	0	$e^0$
1	1	$e^1$

→ in case of correct prediction we take  $\{-1, 1\}$  to reduce loss quickly which shouldn't happen  
as what is aim of predicting correctly? → ANOMALY

Date : \_\_\_\_\_

Let incorrect prediction be denoted by  $E$

∴ correct prediction be denoted by  $1-E$

$$\therefore \textcircled{1} = e^{-\alpha}(1-E) + e^{\alpha}E \quad \text{--- } \textcircled{2}$$

In order to minimize it, we must partial differentiate it w.r.t  $\alpha$ , i.e. we have to find the weight that minimizes loss

$$\frac{\partial E(\text{Loss}_{\text{exponential}})}{\partial \alpha}$$

$$= -e^{-\alpha}(1-E) + e^{\alpha}E \quad \text{--- } \textcircled{3}$$

$$\text{Put } \textcircled{3} = 0 \Rightarrow -e^{-\alpha}(1-E) + e^{\alpha}E = 0$$

$$\Rightarrow e^{-\alpha}(E-1) + e^{\alpha}E = 0$$

$$\Rightarrow \frac{E}{e^{\alpha}} - \frac{1}{e^{\alpha}} + e^{\alpha}E = 0$$

$$\Rightarrow E - 1 + e^{2\alpha}E = 0$$

$$\Rightarrow E(E(e^{2\alpha}+1)-1) = 0$$

$$\Rightarrow E = \frac{1}{e^{2\alpha}+1}$$

$$\Rightarrow \alpha = \frac{1}{2} \log\left(\frac{1-E}{E}\right)$$

$$\Rightarrow e^{2\alpha} + 1 = \frac{1}{E}$$

$$\Rightarrow e^{2\alpha} = \frac{1}{E} - 1$$

$$= \frac{1-E}{E}$$

$$\Rightarrow 2\alpha = \log\left(\frac{1-E}{E}\right)$$

Reduction  
as very fast

Date:

$$e^{-y h(x)}$$

$$\{0, 1\}$$

$$\{-1, 1\}$$

$y \quad h(x) \quad \text{LOSS}$

$y \quad h(x) \quad \text{LOSS}$

		$e^0(1)$
CORRECT PRED	0	$e^0(1)$
	1	$e^{-1}(-0.37)$

		$e^{-1}(-0.37)$
CORRECT PRED	-1	$e^{-1}(-0.37)$
	1	$e^{-1}(-0.37)$

		$e^0(1)$
WRONG PRED	0	$e^0(1)$
	1	$e^0(1)$

		$e^{(2.7)}$
WRONG PRED	-1	$e^{(2.7)}$
	1	$e^{(2.7)}$

Base-learners

penalty

smaller w.r.t  
 $\{-1, 1\}$

This has bigger penalty

decrease in error

$$= V - (2.7 - 0.37)$$

↓  
measure the  
classifiers used

decrease in error

$$= V - (1 - 0.37)$$

Decision Trees don't deal with probabilities

weights of sample records = prob. distribution of sample

For a s record, let  $p_1$  be its prob. of selection in first iteration

Let  $p_2$  be .. .. .. " second "

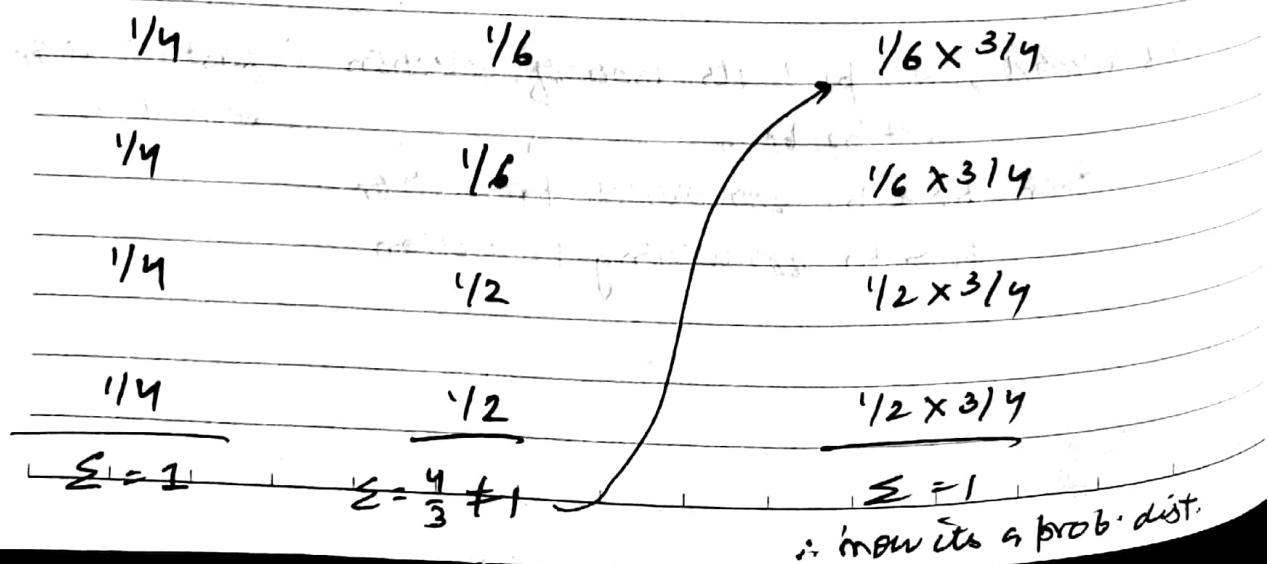
then  $p_2 < p_1$  for correct prediction

$p_2 > p_1$  for wrong prediction

ITERATION Date: 1	ITERATION 2	Normalising
$w_1^{(1)}$	$w_2^{(1)} / \sum w_2^{(1)}$	If $\epsilon = 0.5$ $\alpha = \frac{1}{2} \log \left( \frac{1-\epsilon}{\epsilon} \right) = 0$
$w_1^{(2)}$	$w_2^{(2)} / \sum w_2^{(2)}$	$\therefore$ we ignore such classifier
prob. distribution in first round	we divide by $\sum w_2^{(1)}$ to get prob distribution otherwise	If $\epsilon < 0.5$ $\alpha = +ve$ as $\log \left( \frac{1-\epsilon}{\epsilon} \right) > 1$ this classifier is better than random guess
$\text{as } \sum w_1 = 1$ as each $w_i = \frac{1}{N}$	$\sum w_2 > 1$	If $\epsilon > 0.5$ this classifier is worse than random guess & we stop reject this classifier

$$\beta = \frac{\epsilon}{1-\epsilon} \text{ we can } \frac{1-\epsilon}{\epsilon} = \beta \text{ but calculations are to be varied}$$

1<sup>st</sup> iteration      2<sup>nd</sup> iteration



Date: \_\_\_\_\_  
 one at a time ↴  
 learns description  
 of class each  
 class  
 as conjunction of  
 conditions

## DISCRIMINANT FUNCTION

Bayes classifier

Probabilistic



Rule-based classifiers

Non-Probabilistic

Decision Tree

Non-Probab.

Non-parametric

here, we assume

the form of  
probability  
distribution func.

↳ parametric

doesn't use probability for classification  
it's used only in impurity measures

learns the hierarchy, not always  
one at a time for class

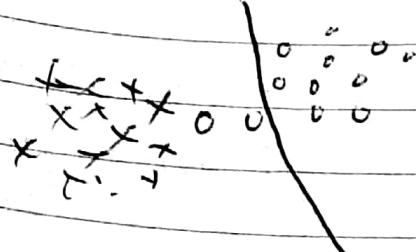
we calculate prior probabilities from the data, also the likelihood  
of class conditional probability.

Decision Trees looks at data as a whole while Rule based  
looks one class at a time.

→ importance of conditions at global level

All classifiers learn decision regions & decision boundaries

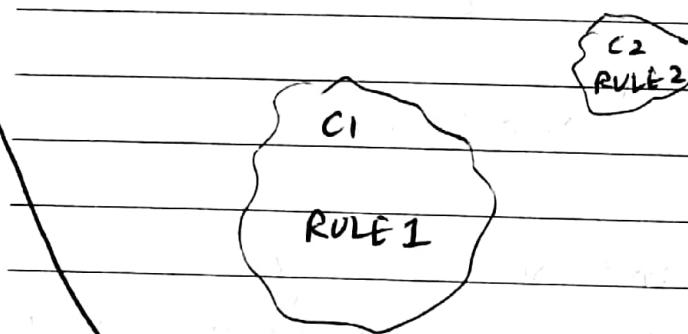
Bayes classifier learns a region using class conditional prob.  
 decision prob. dist. function



Date : \_\_\_\_\_

Decision tree learns decision boundaries based on importance of conditions in hierarchical manner.

Rule based classifier learns decision boundary (arbitrary) using class description. One class can be described by a no. of regions.



NB: Non linear classifiers :: boundaries are non linear  
Bayes classifier is typically non linear but it can be shown linear under certain constraints

DISCRIMINANT FUNCT<sup>N</sup>: A function which discriminates b/w classes

For each class, we have a dist. function  $g_i(x)$  for given  $x$

$g_1(x)$   
;  
 $g_k(x)$

→ assign class with max value.  
( used by Naive Bayes classifier )

Let  $g_0(x)$  &  $g_1(x)$  be 2 ~~dist~~ function values

$$g(x) = [g_1(x) - g_0(x)] > 0 \Rightarrow g_1(x) > g_0(x)$$

if  $g_1 = 0$  then either

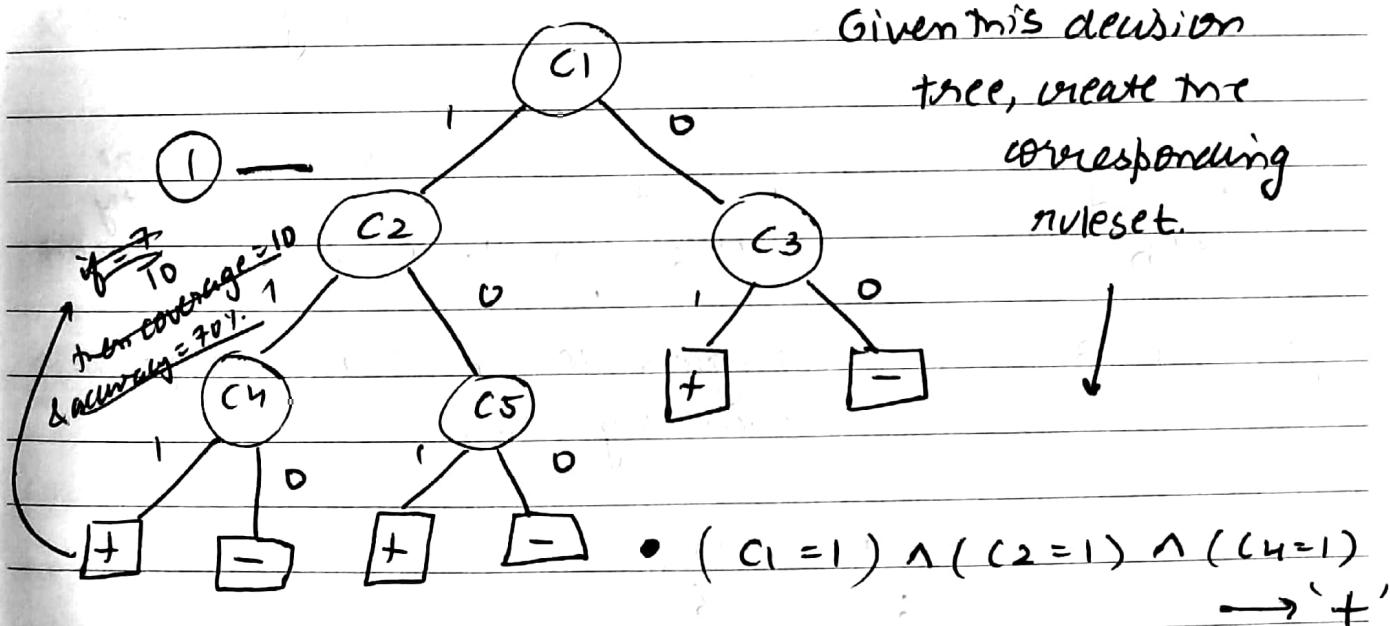
Linear Discriminant Function (in syllabus)

↳ class boundary is linear → divides data into 2 regions

1-D data : a point is a boundary

2-D data : a line (i.e.  $y = mx + c$  is a boundary)

3-D data : a plane



'+' gives:  $(C_1=1) \wedge (C_2=1) \wedge (C_4=1)$

$(C_1=1) \wedge (C_2=0) \wedge (C_5=1)$

$(C_1=0) \wedge (C_3=1) \rightarrow$

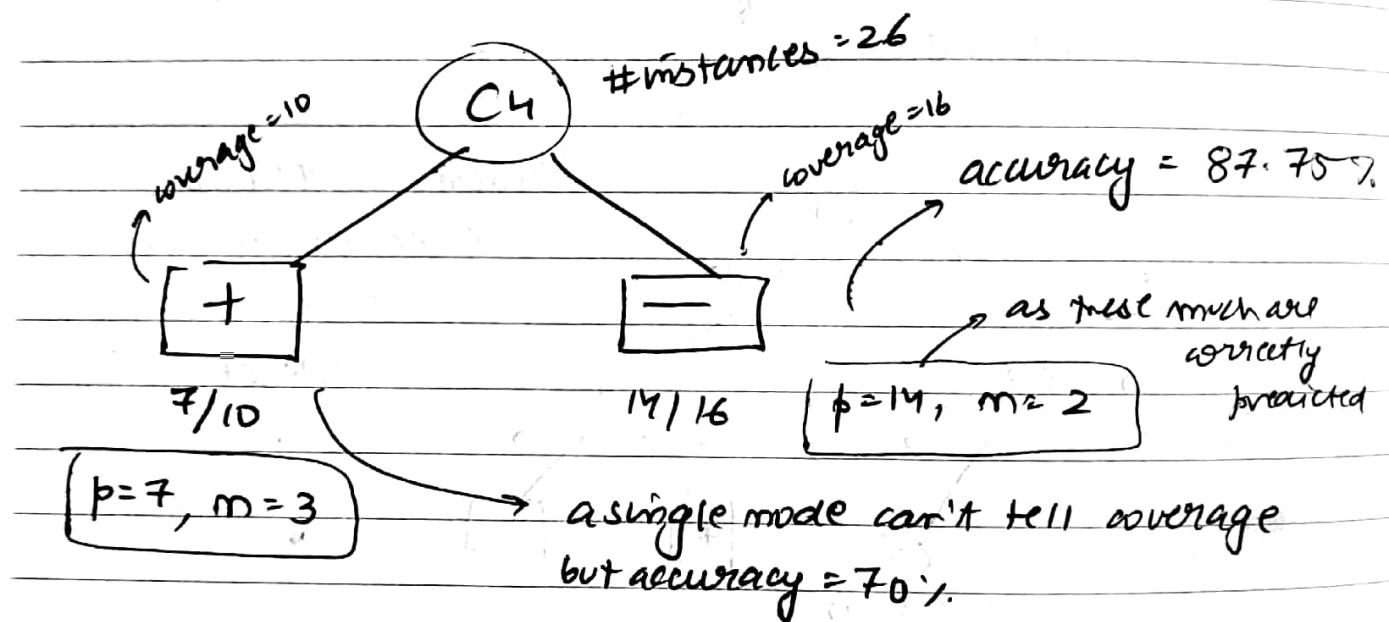
Arrange acc. to no. of rules:  $(C_1=0) \wedge (C_3=1)$

$(C_1=1) \wedge (C_2=1) \wedge (C_4=1)$

$(C_1=1) \wedge (C_2=0) \wedge (C_5=1)$

Date :

- ① Since is a binary tree,  $\therefore$  The no. of rules for each class were same. However, this might not always be true in case of multiclass rules. & to create a decision tree from multiclass, no. of rules from each class must be same.



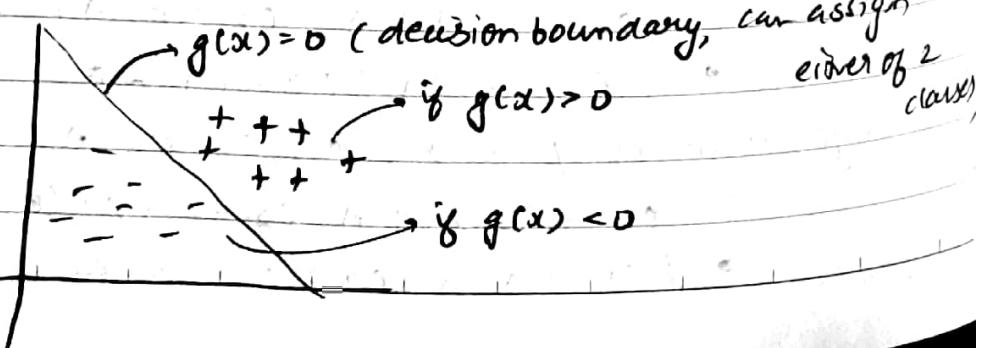
$$\frac{7+14}{26} = \frac{21}{26}$$

## DISCRIMINANT FUNCTION (CONT'D.)

### LINEAR DISCRIMINANT FUNCTION FOR 2D DATA

$$g(x) = mx + c,$$

$g(x) = 0$  defines the linear boundary



IN 3D DATA

linear deci. function takes form of a plane with one class above the plane & one class below the plane

e.g. each instance  $x = (x_1, x_2, x_3)$

$\hat{y} g(x) = 0 \Rightarrow$  pt. lies on decision boundary plane of the 2 classes

IN d-D DATA

$$x = \{x_1, x_2, \dots, x_d\}$$

called  
hyperplane.

$$= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}_{d \times 1}$$

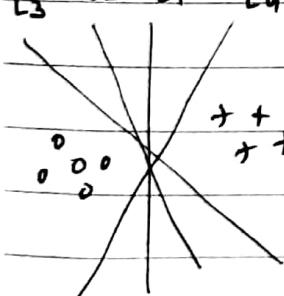
$$w = \begin{bmatrix} w_{11} \\ \vdots \\ w_{id} \end{bmatrix}_{d \times 1}$$

$$g(x) = w_0 + w_i^T x$$

$$g_i(x) = \underbrace{w_{i0}}_{\text{scalar}} + \underbrace{w_i^T x}_{\text{scalar}}$$

tells locat<sup>n</sup> of plane in dD space wrt origin

, locat<sup>n</sup> of hyperplane  
'm' slope defines orientat<sup>n</sup> of line

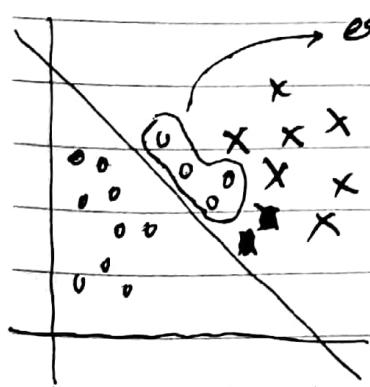


choose the line which maximizes the margin  
this can be done by SVM (Support Vector Machines)

• I)

How to predict a line? : start with a random value of

m & c



error, we need to shift the boundary  $\rightarrow$  how much to shift? ( $\eta$ )

if we shift above a lot,  $g(x)$

will drastically change & then might a lot of instances incorrectly.

Date : \_\_\_\_\_

$\eta$  should not be very big (as before)

$\eta$  should not be very small as convergence to sol<sup>n</sup> will be very late  
Each time, we should change the plane by a different amount

→ depends on the amt of error, Shift & error

If data is not linearly separable (completely)

then loss function will never be 0.

Date : \_\_\_\_\_

$\eta$  should not be very big (as before)

$\eta$  should not be very small as convergence to  $\text{GD}^m$  will be very late  
Each time, we should change the plane by a different amount

→ depends on the amt. of error, Shift  $\alpha$  error

If data is not linearly separable (completely)  
Then loss function will never be 0.

For the Bayes Classifier

Takes 'd' dimensional data  $\rightarrow$  D.P is prob.  $[0, 1]$

$B(x) : \mathbb{R}^d \rightarrow \mathbb{R}$

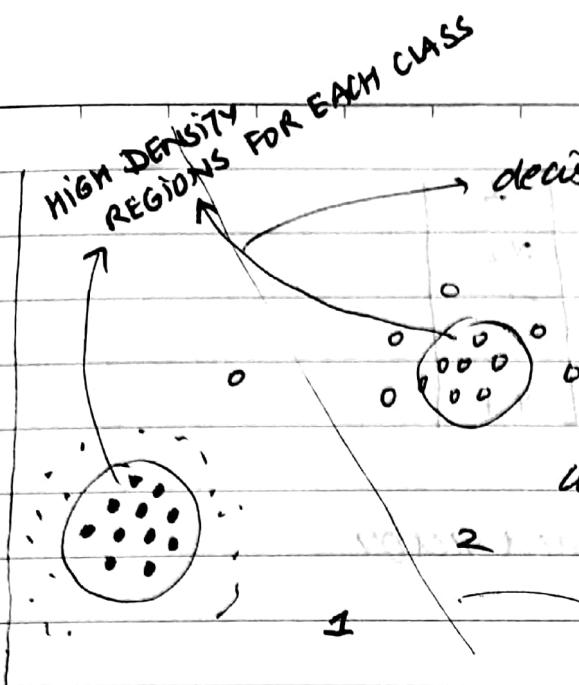
$B(x) > 0.5 \rightarrow \text{class 1}$   
else  $\rightarrow \text{class 0}$  ] only if 2 class problem

$g_i(x | \phi_i)$  : discriminant function for  $i^{\text{th}}$  class

In discriminant function we assume the model & not density  
unlike in Bayes classifier where a prob-density funct<sup>n</sup> is assumed

Discriminant functions are being assumed linear

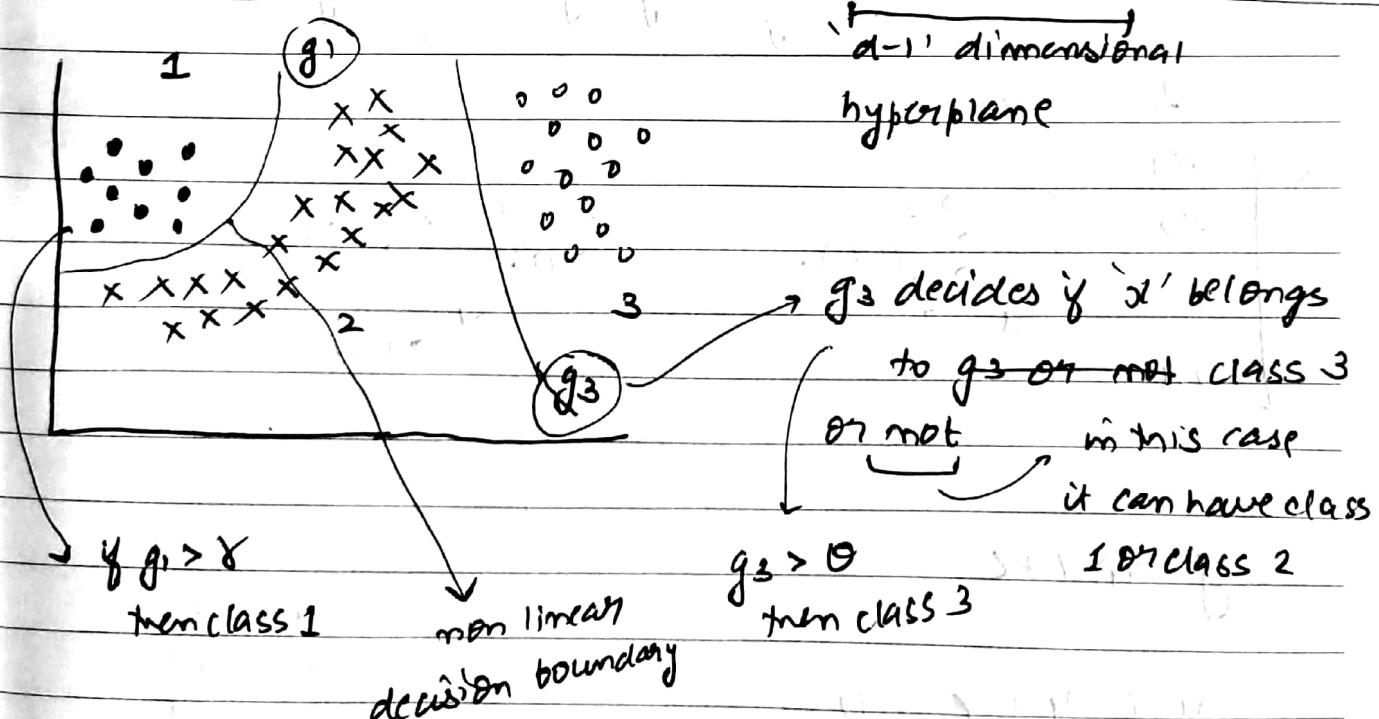
Discriminant based classification assumes a model for  
 $g_i(x | \phi_i)$ , no density estimation



How to learn the boundary & s.t loss function is minimized?

what model of  $g_1(x)$  &  $g_2(x)$  minimize loss function.

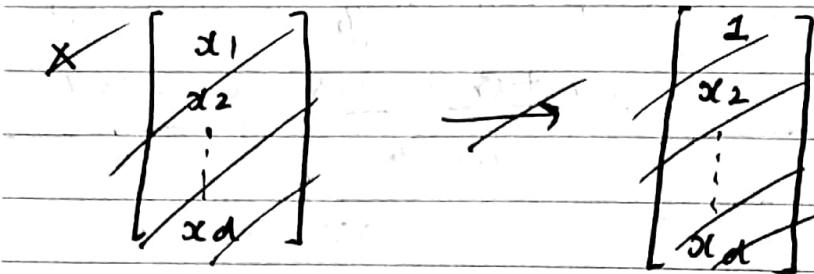
$$g_i(x|w_i, w_{i0}) = w_i^T x + w_{i0} = \sum_{j=1}^d w_{ij} x_j + w_{i0} \quad (1)$$



'm' / slope : determines orientation of the line.

in (1), the weights contribute to class -ve weights inhibit that class # hours jogging  $\propto$  diabetes

Date : \_\_\_\_\_

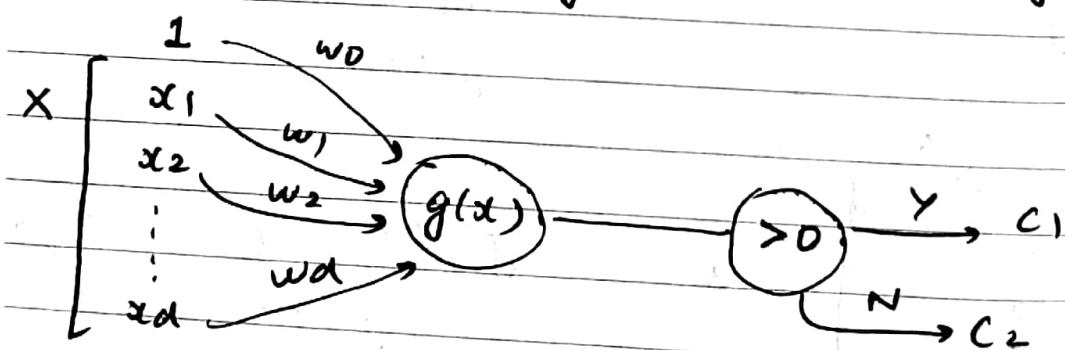


$$g(x) = w^T x \rightarrow \text{augmented vector}$$

$$x \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \longrightarrow \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

before

$$\begin{aligned} & w_0 + \sum_{j=1}^d w_j x_j \\ &= \sum_{j=0}^d x_j \cdot w_j \quad \text{with } x_0 = 1 \end{aligned}$$



$$y = mx + c$$

$$\Rightarrow mx - y + c = 0$$

substituting variables, m, y

$$\Rightarrow w_0 + w_1 x + w_2 y = 0$$

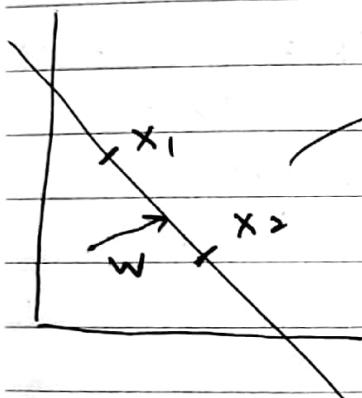
$$g(x) = w_0 + w_1 x + w_2 y$$

$$\begin{aligned} w_0 &= c \\ w_1 &= m \\ w_2 &= -1 \end{aligned}$$

$$g(x) = \underline{w^T x} \rightarrow \text{orientation of the plane}$$

$$= \begin{bmatrix} w_0 & w_1 & \dots & w_d \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

↓  
tells location of plane



$$\Rightarrow g(x_1) = g(x_2)$$

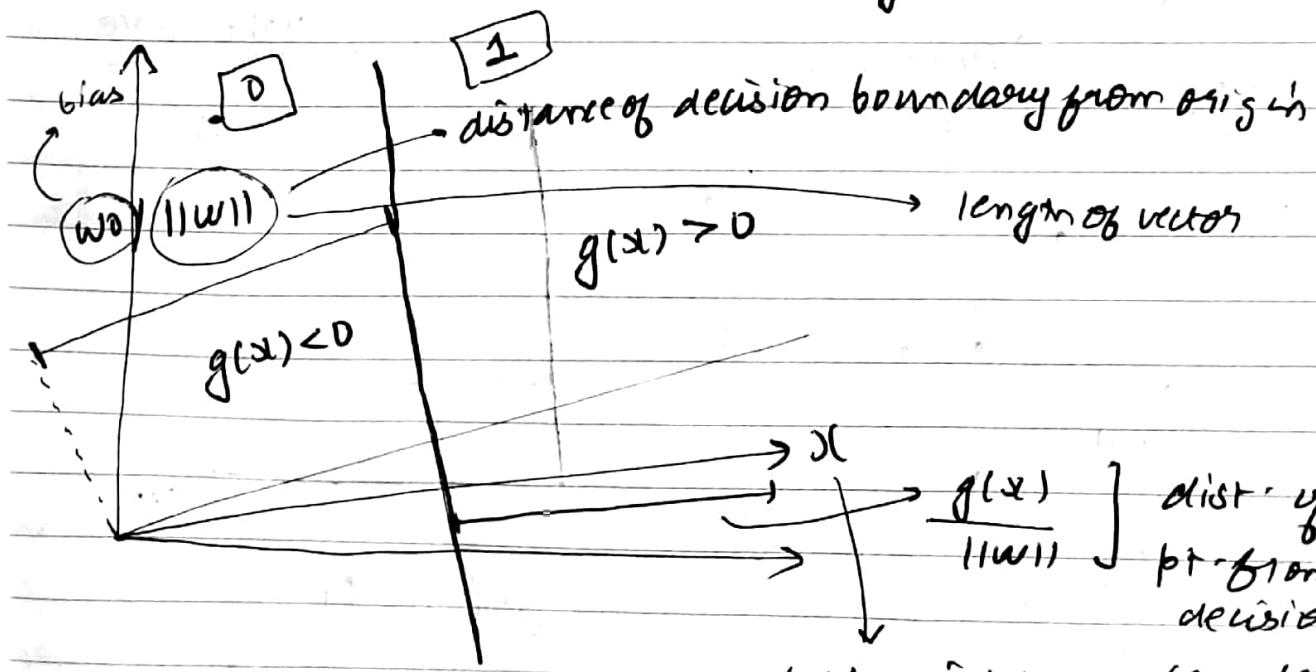
$$\Rightarrow w^T x_1 = w^T x_2$$

$$\Rightarrow w^T x_1 - w^T x_2 = 0$$

$$\Rightarrow \underline{w^T (x_1 - x_2)} = 0$$

vector

this vector is normal to any pt. that lies on plane

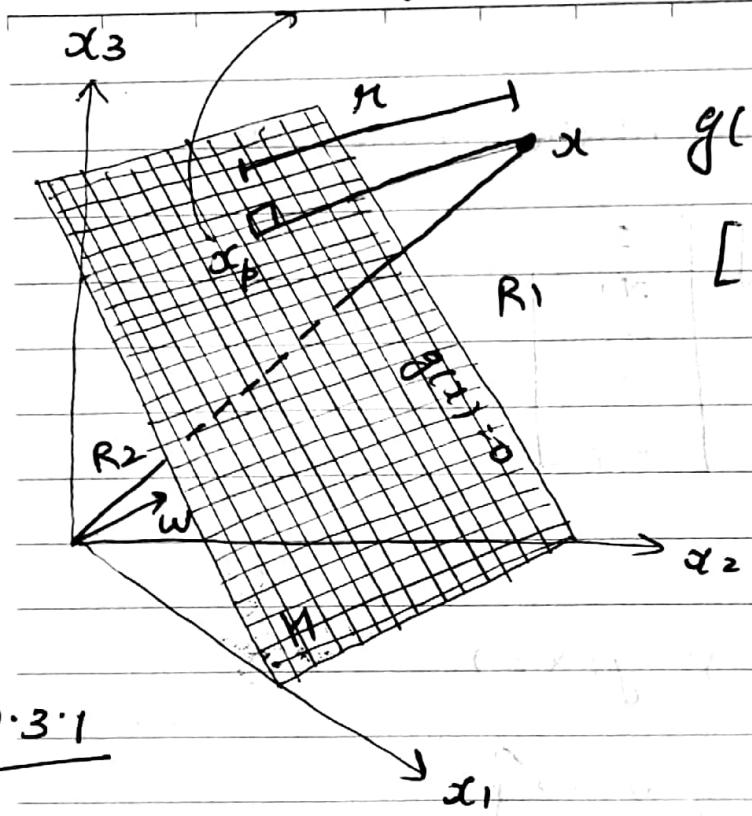


$$g(x) = 0$$

is a hyperplane's

Date :

$\perp$  of  $x$  on plane,  $r$  is the distance



$$g(x) = w^T x$$

$$[w_0 \ w_1 \ w_2 \ w_3]$$

$$\begin{bmatrix} 1 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

Show that

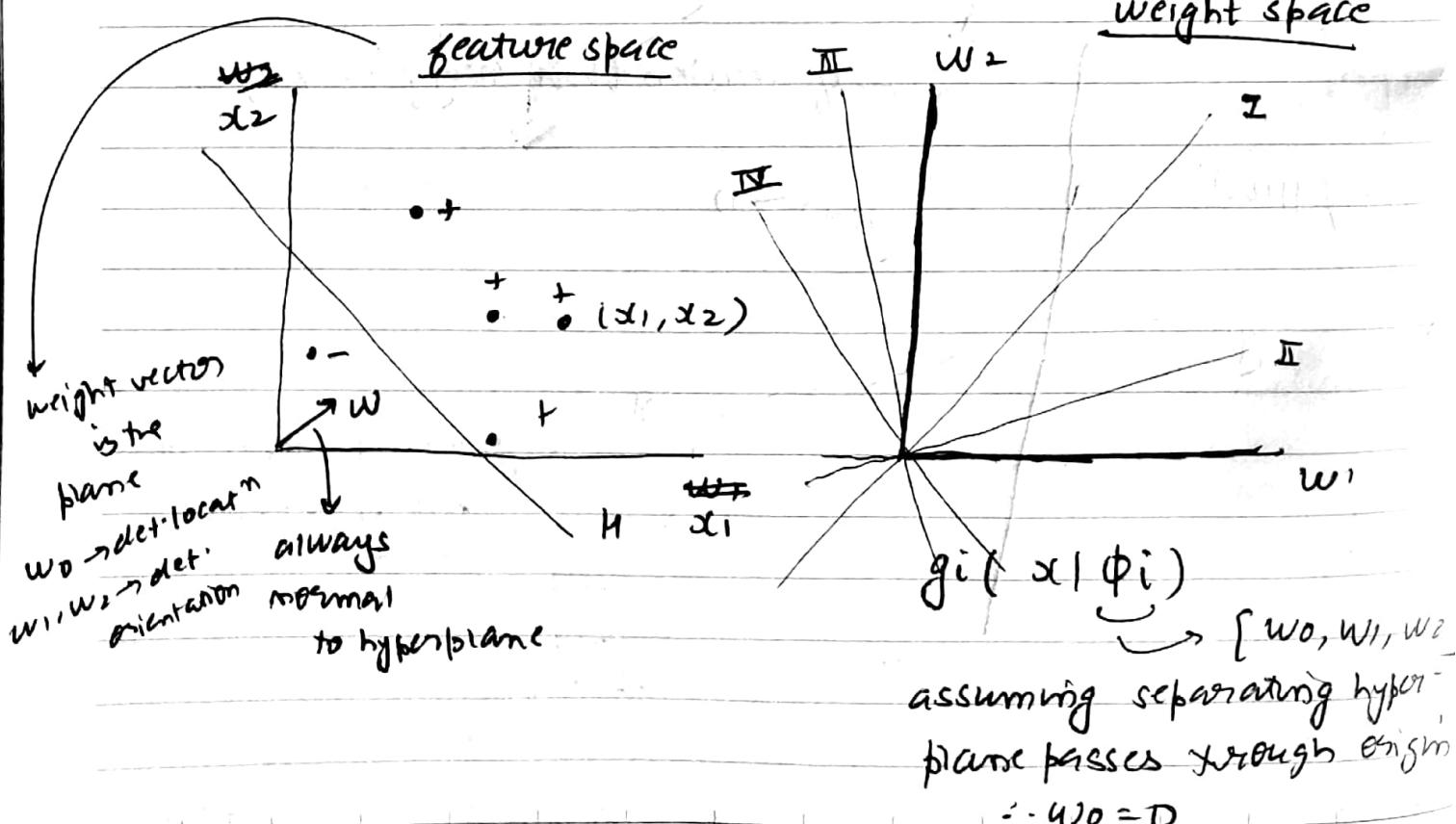
$$x = x_p + \frac{r \cdot w}{\|w\|}$$

10.3.1

$$\text{we know } r = \frac{|g(x)|}{\|w\|}$$

$$w^T x = \begin{cases} > 0 \rightarrow R_1 \\ < 0 \rightarrow R_2 \\ = 0 \rightarrow \text{Hypothesis } H \end{cases}$$

(parametric space)  
weight space

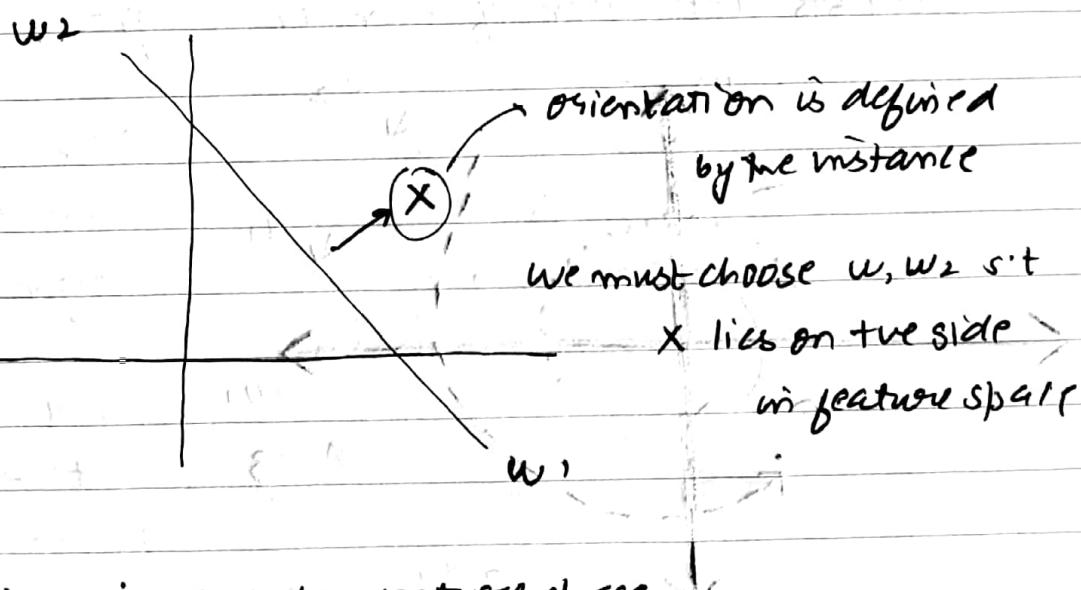


I, II, III, IV are the possible separating hyperplanes that predict '+' in R1 & '-' in R2

We must choose the hyperplane that minimizes loss function

Hyperplane in feature space must be complementary of some hyperplane in weight space.

Multiple separating hyperplanes exist in  $360^\circ$  space weight space



Weight space is comp. to feature space

$$\underbrace{w^T x}_{\text{orientation in weight feature space}} = \underbrace{x^T w}_{\text{orientation in weight space}}$$

w vector is normal to plane

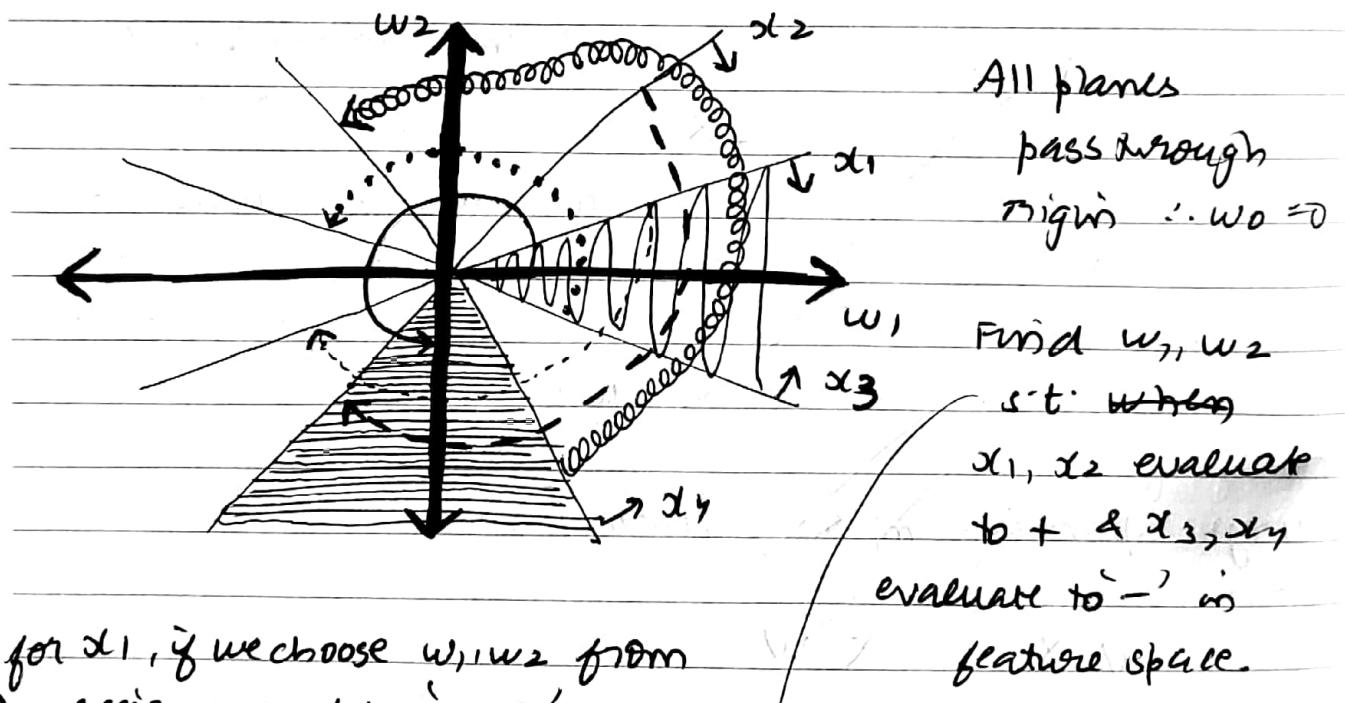
$x$  are normal to plane e.g.  $[x_1, x_2, \dots, x_n]$

Date : \_\_\_\_\_

	$x_1$	$x_2$	$c$
$x_1 \rightarrow$	$x_{11}$	$x_{12}$	$1$
$x_2 \rightarrow$	$x_{21}$	$x_{22}$	$1$
$x_3 \rightarrow$	$x_{31}$	$x_{32}$	$0$
$x_4 \rightarrow$	$x_{41}$	$x_{42}$	$0$

$x_1, x_2, \dots x_4$  are the 4 records

data pts. det. origin of plane in weight space



(I) for  $x_1$ , if we choose  $w_1, w_2$  from region covered by ' $\rightarrow$ ', then  $x_1$  evaluates to true.  $w_1, w_2$  to a plane in feature space s.t.

(II) for  $x_2$ , if we choose  $w_1, w_2$  from region covered by ' $\rightarrow$ ', then  $x_2$  evaluates to true

for  $x_1, x_2$  both to evaluate true, we must take intersection of regions in (I) & (II) to select  $w_1, w_2$ . This region is a semi- $\infty$  region

(III)

for  $x_3$ , if we choose  $w_1, w_2$  from region covered by ' $\dots \rightarrow$ ',  $x_3$  evaluates to +ve  
but  $x_3$  belongs to -ve class

for  $x_1, x_2, x_3$  to evaluate to +ve we take intersection  
of regions (I), (II), (III)

(IV) same for  $x_4$ 

(I)  $\cap$  (II)  $\cap$  (III)  $\cap$  (IV) leads to the prediction for  
all instances

But goal is  $\rightarrow$  det. +ve for  $x_1, x_2$

det. -ve for  $x_3, x_4$

$\therefore$  region denoted by MMM results to 2 false tues

A solution space is region denoted by , as it predicts  
 $x_1, x_2$  as +ve &  $x_3, x_4$  as -ve

Date :

GOAL : to determine weight vector that min. loss function

- multivariate regression
- trial error method by adjusting hyperplane

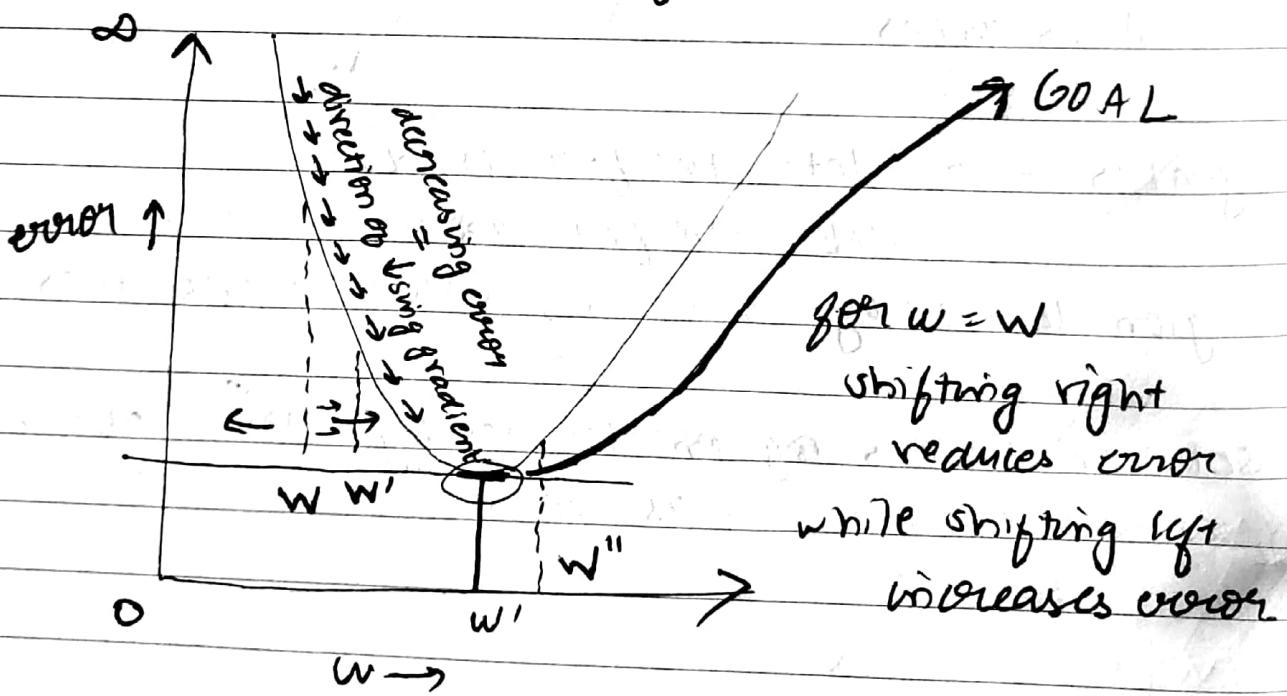
Squared error is the typical loss function

$$L(x_i) = [y - h(x)]^2$$

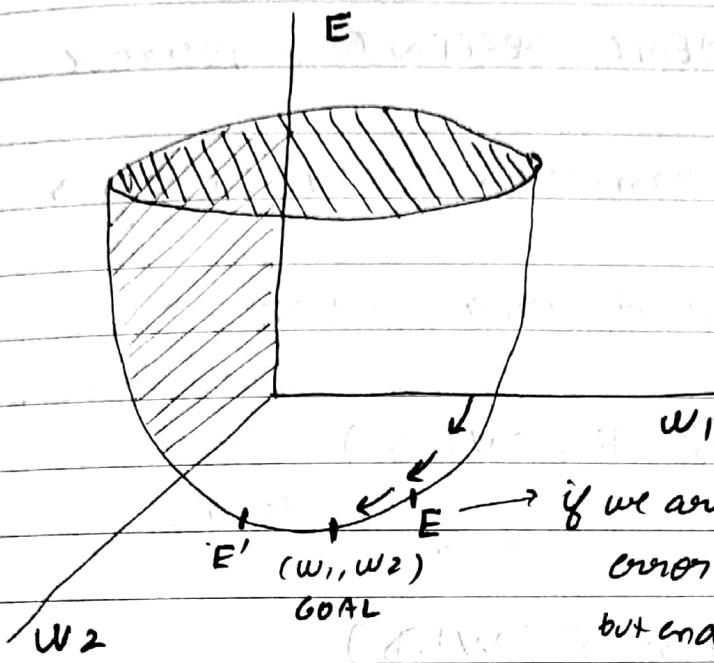
Training error =

$$\sum L(x_i)$$

convex function



Suppose we take a large step & end up overshooting to right of goal i.e.  $w''$ , error increases



$E'$  ( $w_1, w_2$ )

if we are here & to reduce error, we shift it to left but end up overshooting to  $E'$

& if  $E' = E$ , then we keep on oscillating until we don't change the step size

$\therefore$  as error becomes smaller we reduce step size

Wide data : columns  $\gg$  rows or some rows or some columns are dependent - of which  $X^T X$  becomes non invertible.

Inverse time complexity =  $O(m^3)$

Analytical method is good for small data else gradient descent method

Date :

## (GD) GRADIENT DESCENT → iterative

GOAL: Optimize parameter vector  $w$  s.t loss is minimized

LOSS FUNCTION: squared error function

We need to minimize  $E(w|D)$

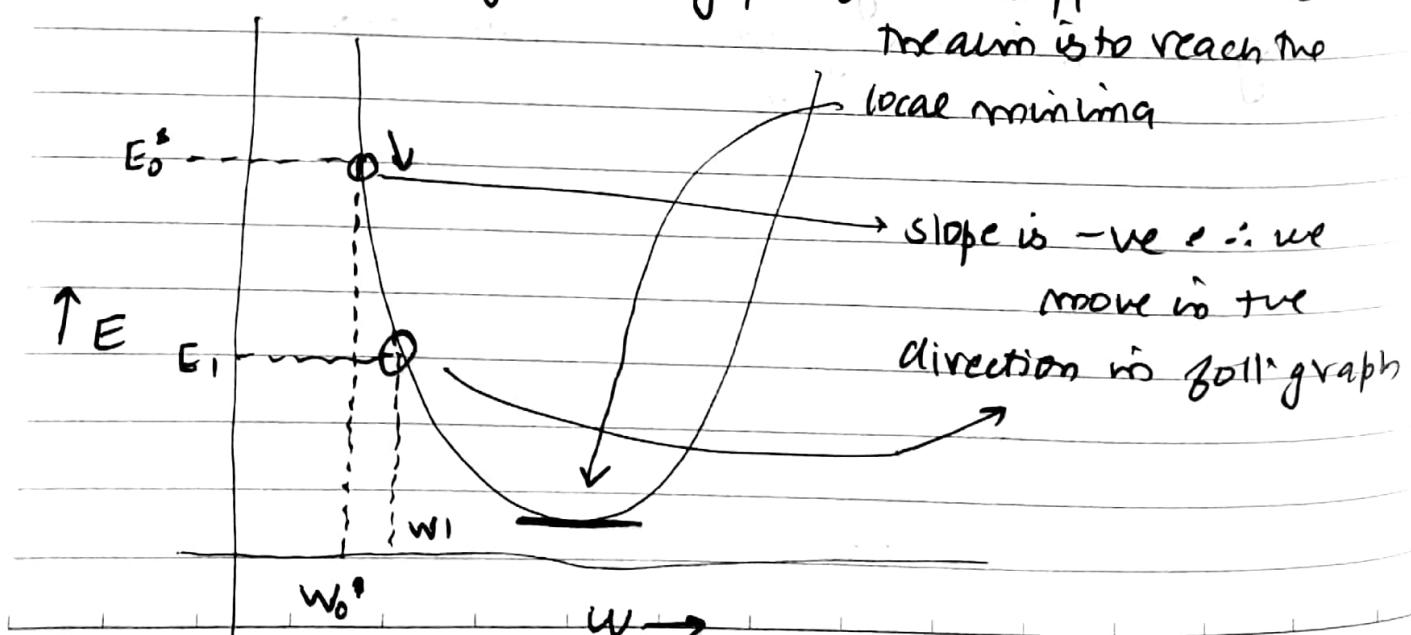
$w = [w_0, \dots, w_d]$  → dimensions of data

$\downarrow$  errors       $\downarrow$  data

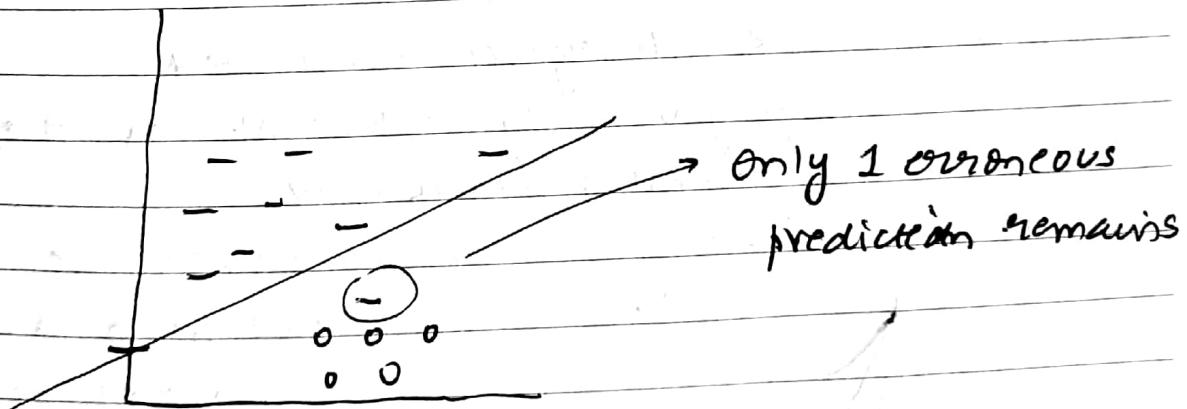
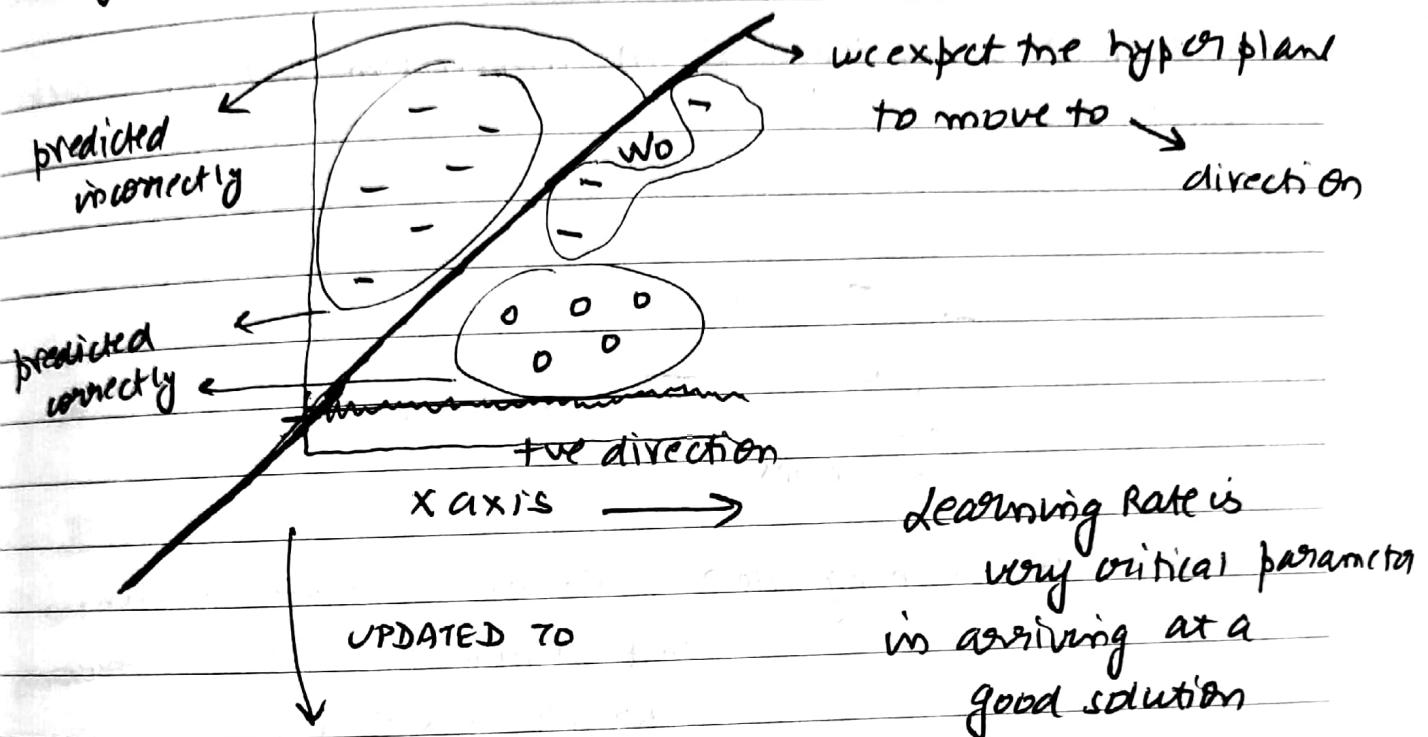
$$w^* = \underset{w}{\operatorname{argmin}} E(w|D)$$

$\downarrow$   $D \text{ or } X$

- Gradient descent is good for smooth convex objective fun<sup>n</sup>
- It reduces error at each step
- Need to determine step size (learning rate)
- GD is good for large 'n' w.r.t analytical method
- we start randomly at any pt. for  $w$ , suppose at  $w_0$



The graph for  $w_0$  in data space is.



GRADIENT VECTOR  $\nabla_w E = \left[ \frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_d} \right]^T$

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

learning Rate

opposite direction

Updated weight

$w_i = w_i + \Delta w_i$

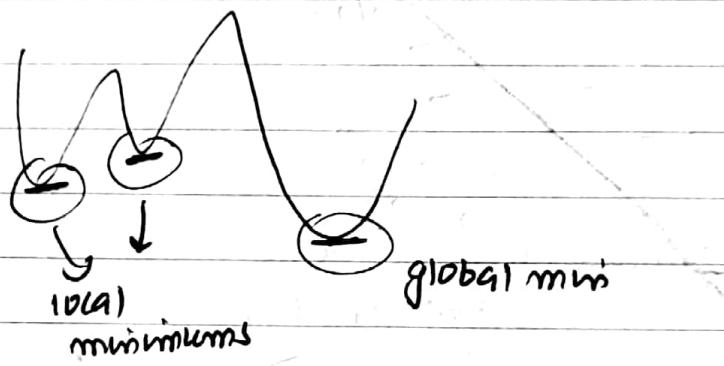
slope - ve  
→ move in the to ↓ error

$w_0$  remains same :- we assumed initial line passing thru  $w_0$

Date : \_\_\_\_\_

Random initialization akin to random separating hyperplane

No guarantee of finding global minimum for non-convex function

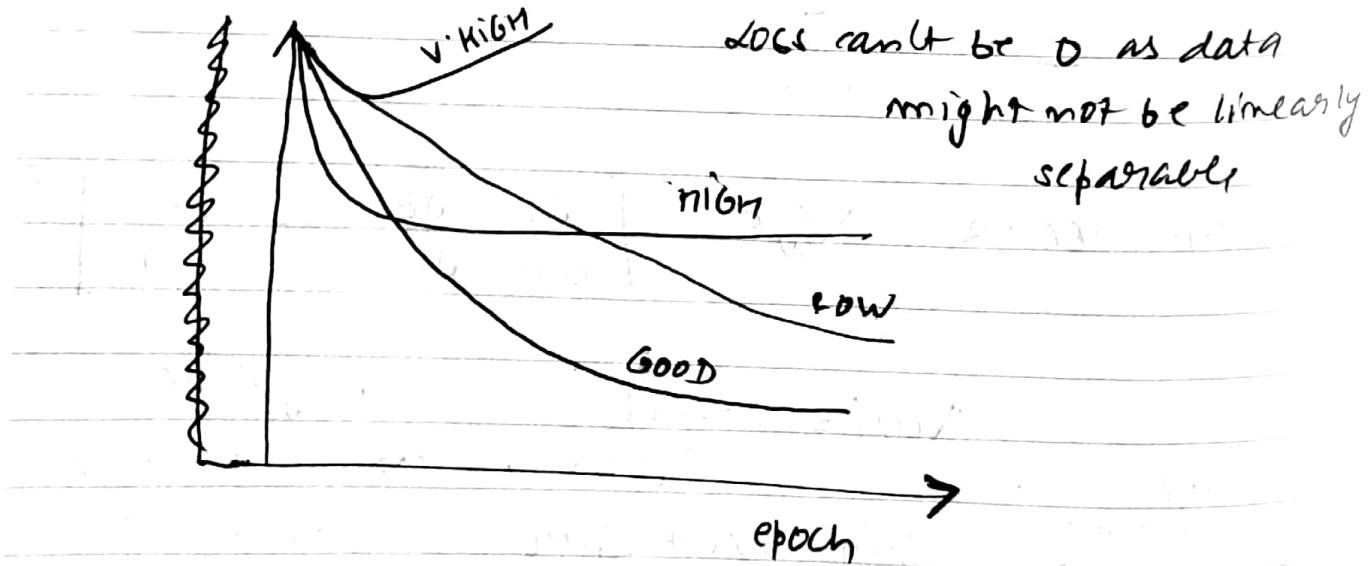


How to decide learning rate? : Critical since determines how fast or slow we will move towards optimal weights

Very high learn LR: divergence is indicated

Low LR: converging very slowly, not feasible as lot of computational power is used.

High LR: might become constant & doesn't reduce



## popular Heuristics for computing LR

i) Exact line search - bisection method

ii) LB                          UB  
find mid pt

iii) LB                          MP                          UB

depending on how loss changes choose  
LB                          MP

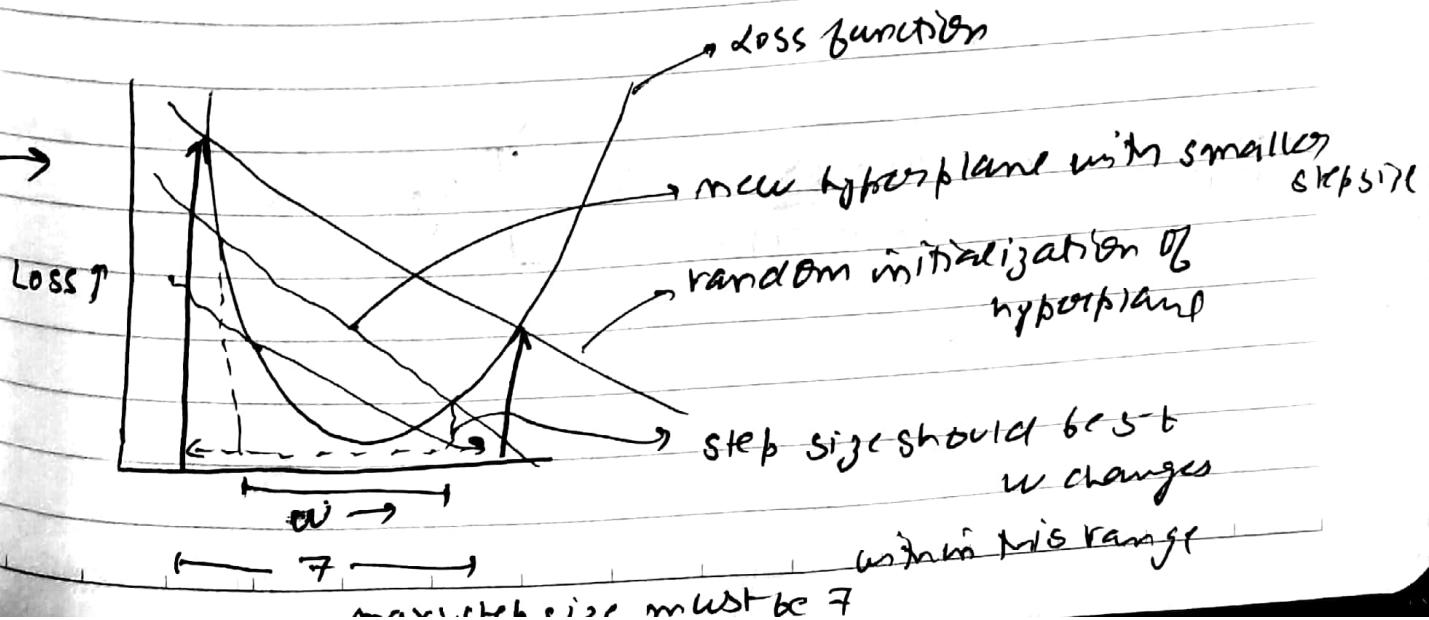
OR  
MP                          UB

2) Inexact line search

3) Decaying learning rate can be  $\frac{1}{K}$ ,  $\frac{1}{\sqrt{K}}$ , etc. depends on data

function should be inverse

4) Fixed learning rate



Date : \_\_\_\_\_

- assumptions:
- 1) class conditional densities are Gaussian
  - 2) and share covariance matrix

## Probabilistic classification as discriminant Based classification

$$y = h(\mathbf{x}) = \begin{cases} 1 & \text{if } P(C_1 | \mathbf{x}) > P(C_2 | \mathbf{x}) \text{ or} \\ & \bullet P(C_1 | \mathbf{x}) > 0.5 \text{ or} \\ & \bullet \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})} > 1 \text{ or} \\ & \bullet \log\left(\frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})}\right) > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\log \left[ \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})} \right] = \log \left[ \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} \right]$$

————— logit transformation —————  
 ↪ log odds of posterior prob of C

$$y = \frac{P(C_1)}{P(C_2)} \quad \boxed{\text{odds}}$$

$$= \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})}$$

$$\log y = \log \left( \frac{P(C_1 | \mathbf{x})}{1 - P(C_1 | \mathbf{x})} \right) \quad \boxed{\text{log odds}}$$

$$\log \left[ \frac{P(C_1 | X)}{P(C_2 | X)} \right] = \log \left[ \frac{P(X | C_1) \cdot P(C_1)}{P(X | C_2) \cdot P(C_2)} \right]$$

$$= \log \left[ \frac{P(X | C_1)}{P(X | C_2)} \right] + \log \left[ \frac{P(C_1)}{P(C_2)} \right]$$

to derive this  $\Rightarrow w^T x + w_0$  i.e. can we interpret

bayesian classifier as "discriminant function"

not augmented vector

~~Bayesian classifier~~

$$\begin{aligned} w_1 x_1 + w_2 x_2 \\ w_1 x_1 - \frac{1}{2} w_2 x_2 \end{aligned}$$

$$C_1 \sim (\mu_1, \Sigma)$$

$$C_2 \sim (\mu_2, \Sigma)$$

$$\log \left[ \frac{P(C_1 | X)}{P(C_2 | X)} \right]$$

We need to find  $P(X | C_1)$  &  $P(X | C_2)$

$\rightarrow$  for d-dimensional  $X$

$$(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu_1)^T \Sigma^{-1} (\mathbf{x} - \mu_1) \right]$$

$$\text{Solving } \log\left(\frac{y}{1-y}\right) = \log\left(\frac{P(X|C_1)}{P(X|C_2)}\right) + \log\frac{P(C_1)}{P(C_2)}$$

Date: \_\_\_\_\_  
 $G = w^T x + w_0$

For 1 dimensional : solving this

$$= (2\pi)^{-1/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right)$$

$$P(X|C_1) = (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right)$$

$$P(X|C_2) = (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(x-\mu_2)^T \Sigma^{-1} (x-\mu_2)\right)$$

$$\log\left[\frac{P(X|C_1)}{P(X|C_2)}\right] = \log P(X|C_1) - \log P(X|C_2)$$

$$\mu_1 = \mu_{1s} \quad \text{Mean of } i^{\text{th}} \text{ feature of class 1}$$

a vector to store mean of all attributes

$$\left(\frac{x-\mu}{\sigma}\right)^2 = (x-\mu_1)^T \Sigma^{-1} (x-\mu_1)$$

$\downarrow$   
 $(x_1-\mu_{11}), (x_2-\mu_{12}), \dots, (x_d-\mu_{1d})$

$$[(x_1-\mu_{11}), (x_2-\mu_{12}), \dots, (x_d-\mu_{1d})]_{1 \times d}$$

$$\begin{bmatrix} \sigma_1^2 & 0 & 0 & \dots & 0 \\ 0 & \ddots & 0 & \dots & 0 \\ 0 & 0 & \ddots & \ddots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \dots & \sigma_d^2 \end{bmatrix}_{d \times d}$$

$$\begin{bmatrix} x_1-\mu_{11} \\ x_2-\mu_{12} \\ \vdots \\ \vdots \\ x_d-\mu_{1d} \end{bmatrix}$$

Date: \_\_\_\_\_

$$\rightarrow = \log \left( \frac{\exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right)}{\exp\left(-\frac{1}{2}(x-\mu_2)^T \Sigma^{-1} (x-\mu_2)\right)} \right)$$

$$= \log \left( \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right) \right) - \log \left( \exp\left(-\frac{1}{2}(x-\mu_2)^T \Sigma^{-1} (x-\mu_2)\right) \right)$$

$$= -\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1) + \frac{1}{2}(x-\mu_2)^T \Sigma^{-1} (x-\mu_2)$$

$$\downarrow$$

$$\left[ -\frac{1}{2} \left( \frac{(x_1-\mu_{11})^2}{\sigma_1^2} + \frac{(x_2-\mu_{12})^2}{\sigma_2^2} + \dots + \frac{(x_d-\mu_{1d})^2}{\sigma_d^2} \right) \right]$$

$$+ \frac{1}{2} \left[ \left( \frac{(x_1-\mu_{11})^2}{\sigma_1^2} + \dots + \frac{(x_d-\mu_{1d})^2}{\sigma_d^2} \right) \right]$$

open all brackets to get

$$-\frac{1}{2} \left[ \left( \frac{\mu_{11}^2 - \mu_{12}^2}{\sigma_1^2} + \dots + \frac{\mu_{1d_1}^2 - \mu_{1d_2}^2}{\sigma_d^2} \right) \right]$$

$$+ \left[ \left( \frac{(\mu_{11}-\mu_{12})x_1}{\sigma_1^2} + \dots + \frac{(\mu_{1d_1}-\mu_{1d_2})x_d}{\sigma_d^2} \right) \right]$$