

①

## Lecture: 1, Machine Learning

Date - 1-8-2024, P23CS0010

Q what is machine learning?

→ ML is a discipline of artificial intelligence (AI) that provides machines ability to automatically learn from data and past experiences.

• Any process by which system improves its performance by experience.

• ML is study of algorithms that :-

→ improves performance P

→ at some task T

→ with experience E.

→ A well defined learning task is given by  $\langle P, T, E \rangle$

e.g:- we get suggestion of next sentence's word while writing mail or while searching something on Google.

⇒ ML algorithms adaptively improves performance with an increase in number of samples during 'learning' process.

## • Why Machine Learning is used ?

(2)

We'll learn with the help of some examples!—

Example 1 + Let's see, mathematic digit 2 can be written in various way by human as:

→ handwritten digit 2:

2 2 2 2 2 2 2 <

goal :- identifying different digits from handwritten text.

• All 2 are in different shape, there is so much variation in data.

⇒ Shape alone is not sufficient to identify digits.

Example 2:- Let's take example of a sentence!—

original sentence :- I will watch football Match.

understood sentence , - I will watch fo~~ot~~ball Match.

• We need to eliminate observer based decision .

Because of so many rules and variation in data we need to learn from experience in ML.

Now let's see some examples of machine learning based day to day tasks -

- Self-driving cars and driver assistance features, automatic stopping, improve overall vehicle safety.
- A bank's fraud detection service automatically flag suspicious transactions.
- Recommendation engines suggest products, songs, television show such as Netflix, Amazon or Spotify.
- Speech recognition software that allow us to convert voice memos to text.

### • Applications of ML

- |                                     |                                   |
|-------------------------------------|-----------------------------------|
| (i) Image Recognition               | (vi) Email Spam/malware filtering |
| (ii) Automatic language Translation | (vii) Self driving cars           |
| (iii) Medical Diagnosis             | (viii) Product Recommendations    |
| (iv) Stock Market trading           | (ix) Traffic Prediction           |
| (v) Online Fraud Detection          | (x) Virtual Personal Assistance   |

## Lecture-2

Date : 06-Aug-2024

①

- why ML is used?
- ML is used for Automation, Speed, Observer independence.
  - Observer Independence is done through decision support systems.
  - e.g:- For disease Prediction different doctor diagnosis differently.
  - e.g:- Astrophysics uses ML to calibrate simulation parameters based on observational data.
  - e.g:- Stock market prediction.

Note :- Final decision will be made by human.  
↓  
ML system supports in making decisions.

- Definition :-
  - ML is field of study that gives computer ability to learn without being explicitly programmed.
- Rule based program fails because of variation in data.

Inductive Learning :- It is learning from experience, i.e!- refining our observation based on experience.

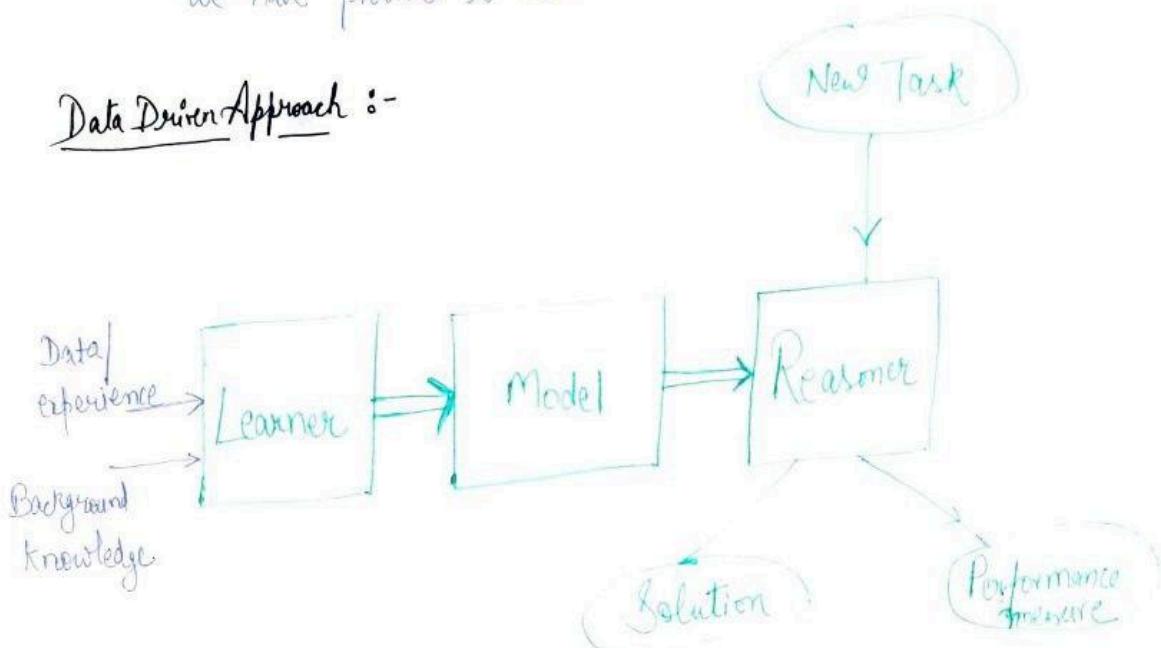
Traditional Programming :-



- In programming language :- behaviour of program don't depend on data.

In ML  $\stackrel{?}{\rightarrow}$  How computer comes up with program depends on type of data we have provided to learn.

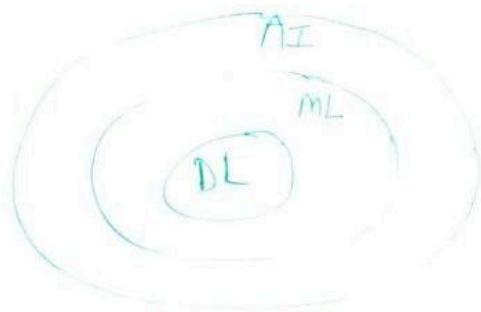
### Data Driven Approach :-



### Steps of creating a learner :-

- i) Choose the training data / experience
  - extract features from data.
  - different method may require different feature extraction.
- ii) choose target function (for classification between cat, dog some fctn give different set of values of cat, dog).
- iii) choose how to represent target function.  
eg:- we want target function to linear or circular or else?
- iv) choose a learning algorithm to infer target function.

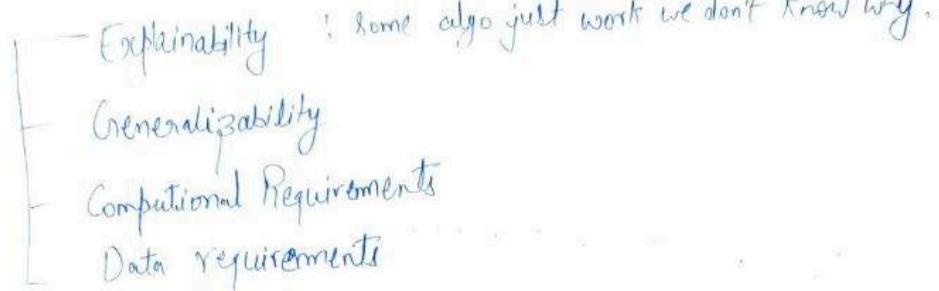
- How to represent target function ! -  
 → may be complex (but richer representation better result)  
 → difficult to learn
- Set of target function : Hypothesis.



History of ML :-

- ML is used to known as Pattern Recognition (PR).
- i) ML is used to known as Pattern Recognition (PR).
  - ii) 1950 : Samuel's checker playing algorithm.
  - iii) 1990 : support vector machine
  - iv) After 1990 :- Adaboost got tremendous success (ensemble learning)  
Bioengineering especially.
  - v) Popularity of this field in recent time :- new software, new hardware,  
  neural n/w  
  CPU  
→ cloud enabled = Availability of big data.
  - vi) Alexnet changed innovation = Deep learning come into existence.  
eg:- superresolution problem after ML Boom.

14



### Applications of ML

- Recommendation System
- Computer Vision
- Medical data Analysis
- Robotics, Biometrics
- Speech recognition
- Autocomplete.
- Security, banking, finance, physcs.

### Different types of ML

(i) Supervised

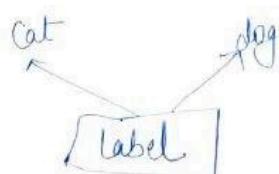
(ii) Unsupervised

(iii) Reinforcement Learning

### Types of ML learning

#### (i) Supervised learning :-

Input :- Image + Label = Training



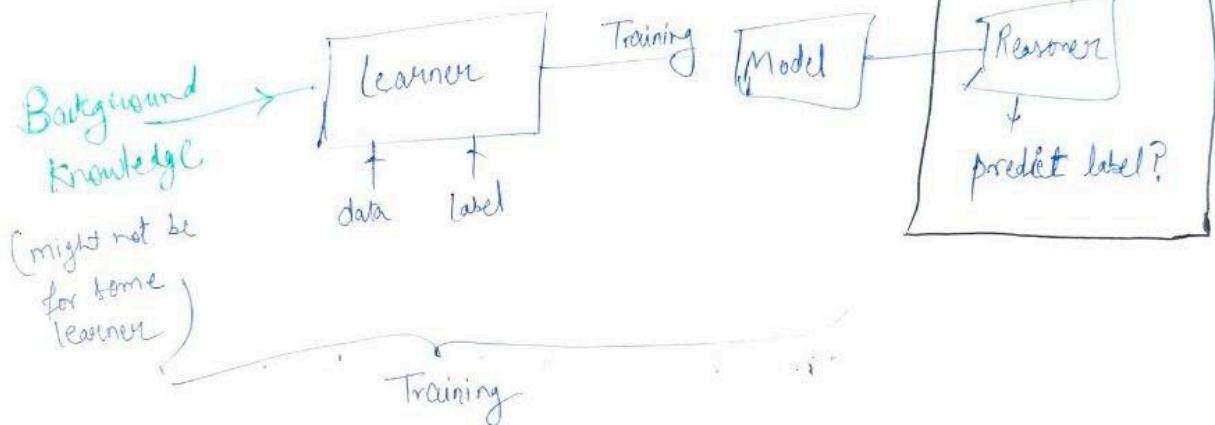
→ decide learning algorithm

• Adding Label of image → Cat or dog?

→ **It** is kind of ml where algorithm is trained with label data.

(5)

Task of supervised learning :-  
Predict label of unlabeled data;  
House price prediction.



- Validation :-
- Train model
  - With some amount of data validate

Validation is internal testing after training.

e.g. Training :- when company manufactures car.

Validation :- Based on internal checking deciding which thing to fit, tune, optimize change.

Hypothesis tuning / optimization.

Testing :- when customer desires the car | selling of car.

- During training Background knowledge might not be there for some learner.

- For some we may have separate Training, Testing. (7)
- eg: finding similar viruses based on genetic pattern recommendation system based on customer's previous choices.
- Making cohesive group of data input points.
- we can define performance measure but need not to have training data set.

eg:-

• Clustering	• Dimensionality reduction	• Association rule mining
• Generative model	• unsupervised feature extraction.	

↑  
extract feature from image.

- Dimensionality Reduction All features might not be meaningful, finding salient features in data.
- ↓  
Data compression.
- most of generative model is unsupervised.
- Discriminative model :- creates new image.

feature space The space defined by features is called feature space.

or

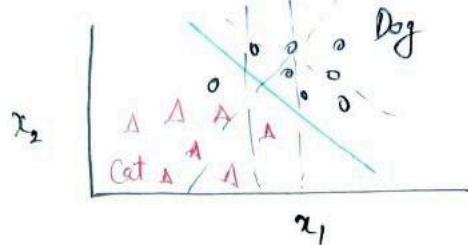
It is an Euclidean space defined by features.

- Any data / sample is a point in the (feature space).

need not to be always  
euclidean

we assume feature space  
to be euclidean space

Hypothesis :-



⑦

(full not approximates  
target function.)

→ How many such hypotheses are possible?

: Let's assume that I have  $N$  data points.

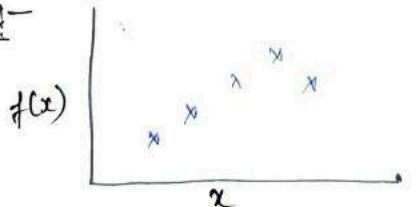
• Each data points can be classified in one of the two classes.

• So possible class combination for  $N$  data points is  $2^N$ .

→ 1 hypothesis for each class combination so,  $2^N$  hypothesis are possible.

It is very very huge.

Inductive Bias :-



• A regression problem with data samples  $x$  and regression value  $f(x)$

Inductive Learning (for supervised learning)

→ necessary as we have some amount of data not entire data.

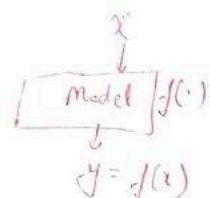
Given data points  $X$  and label  $Y$ .

or we may consider data  $X$  and a function  $f(x)$

→  $f(x)$ : Categorical for classification

→  $f(x)$ : Real no. for Regression

→  $f(x)$ : Probability of  $X$  for probability estimation



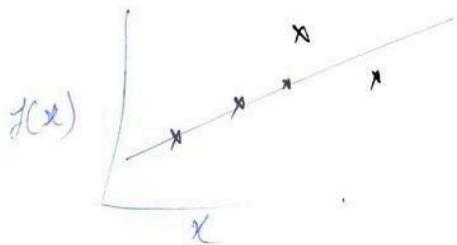
• Learning from experience.

## How to obtain features?

(10)

- Features may be obtained from raw data  
egt as texture, SIFT feature etc
- may be obtained through measurement at the time of data collection.  
eg + height, weight etc. of different animals.
- Feature Engineering
- Many NN methods can directly deal with raw data without explicitly requiring features at input.

## Hypothesis &



- Regression problem with data samples  $x$  and regression value  $f(x)$   
which curve fits best?

Decision Boundary → The boundary by which decision before and after changes. (1)

→ The (set of all legal hypothesis) which help to solve ML problem is known as Hypothesis Space.

Q: How many hypothesis can exist?

Sol<sup>n</sup>t Let's assume that we have ' $N$ ' data points.

If each data points can be classified into one of 2 classes,

→ So possible class combination for  $N!$  data points is  $2^N$ .

not possible as could be very large.

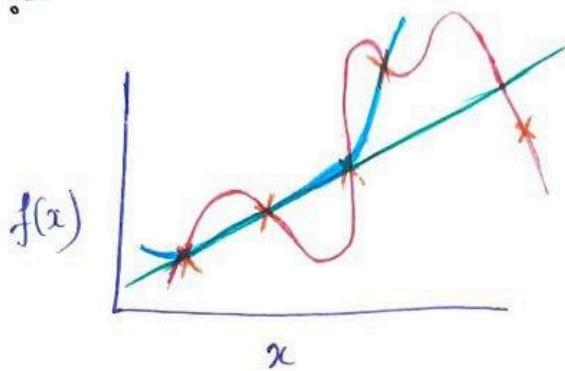
We have to make some restrictive expression.

Date : 13-08-2024

## Lecture - 4

①

Hypothesis :-



Let's consider regression problem with data samples  $x$  and regression value  $f(x)$ .

Q :- which curve fits best ?

Note :- our decision of fitting the curve depends on our assumption about problem.

[  
our assumption is :- Inductive bias  
↓  
(based on my Prior)  
Inductive bias reduces Search space]

→ Inductive bias reduces  $2^n$  assumptions to finite no. of hypothesis.

► No Free Lunch theorem :-

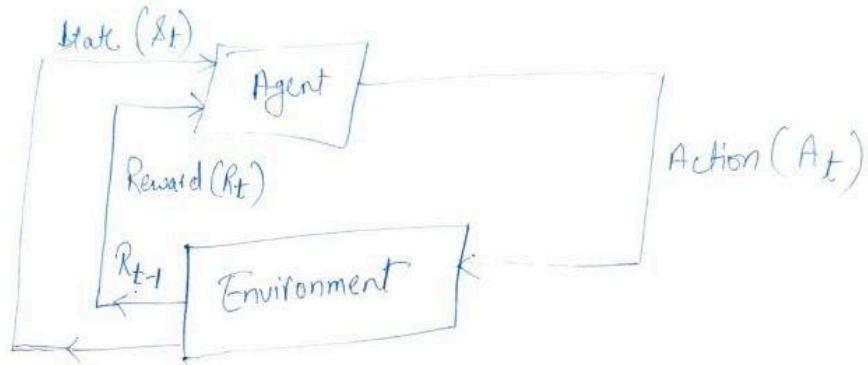
In essence, unless you make assumption about problem, you can't make solution.

[ There is issue in ML which is Generalizability. ]

Generalized over  
hypothesis that not only fits training data, also unseen data.

## Reinforcement learning

(3)



What do we learn?

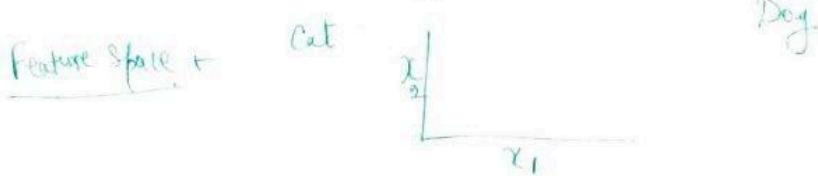
To control an agent (bicycle) from sequence of good decisions

How do we learn it?

From experiences through reward from the environment.

## Inductive learning (for supervised learning)

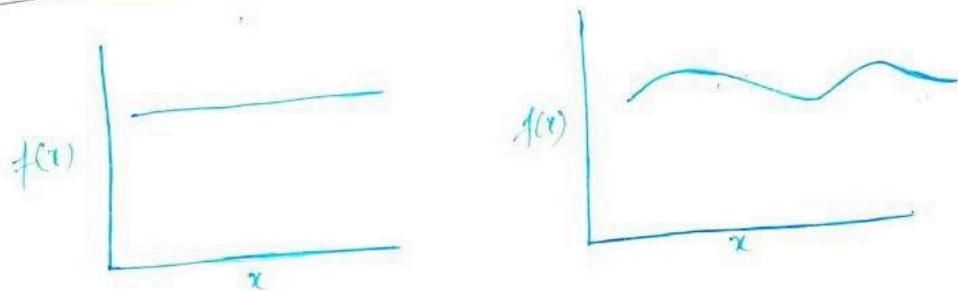
- Given data  $X$  and label  $Y$
- we may consider data  $X$  and a function  $f(x)$
- Return a function  $h$  that approximates  $f$ :
  - we don't know  $f$
  - But we assume that there is an  $f$ .
  - $h$  is called a hypothesis.



$x_1$ : Height  
 $x_2$ : Weight

- we want generalized hypothesis. ④
- ▲ Ockham's razor : As training data may fit to multiple curves,  
 → Prefer simpler hypothesis consistent with data.
- This is just an assumption and does not always work. (not universal truth)

eg:- Nature follows simple rule  
 i.e. mixing of gas in room.



+  
 more assumptions to get  
 this curve.  
 simple curve → more assumption = more bias

eg:- In order to drive on straight line  
 we consider some assumptions like good weather, road, good vehicle,  
 good cond<sup>n</sup>, expert rider etc.

- Realistic curve → less assumption → more complex curve

## Types of Inductive Bias

(a) Restriction / Absolute bias      (3)

(b) Relative / Preference bias

### (a) Restriction bias / Absolute bias:

- representational power of algorithm.

i.e.: set of hypothesis that algo will consider.

e.g.: my algo will consider only straight line.

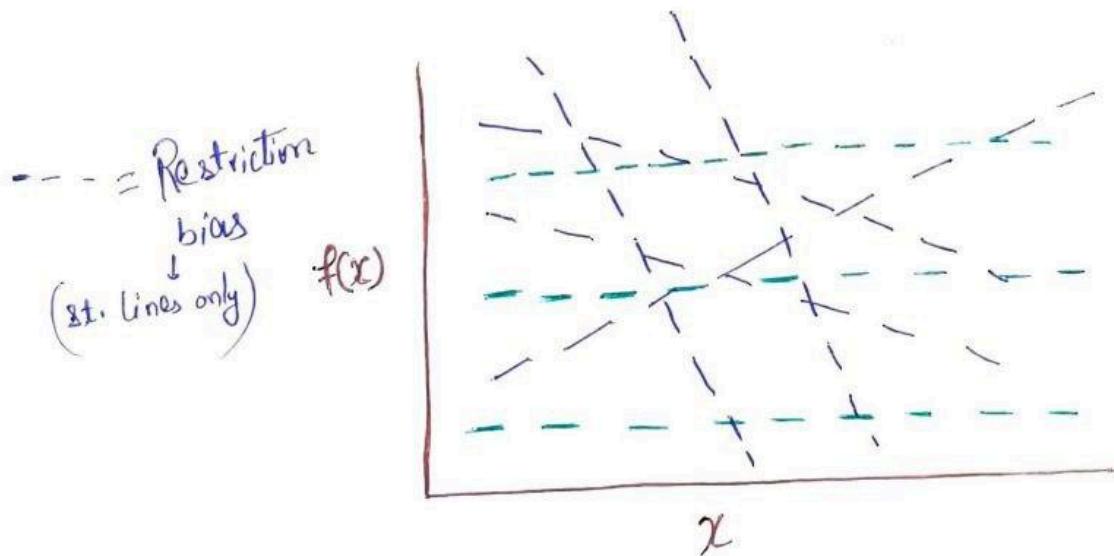
- even of this assumption many straight lines possible



### (b) Preference bias / Relative bias:

- what type of representation my algo prefers.

e.g.: my algo may prefer straight lines parallel to x-axis.



(--- = Relative bias = st. line  $\parallel$  to x-axis only)

## Inductive learning :

(4)

- i Let's take general fun from training examples

(a) Construct  $h$  that best approximates target fun  $f$ .

$$h \rightarrow f$$

(b)  $h$  is consistent if it agrees with all training examples.

(c)  $h$  is generalized well if correctly predicts label for new examples

(d) learning  $\rightarrow$  refining hypothesis space.

Inductive learning is an ill-posed problem,

• Ill-posed Problem: unless we see all possible examples related to problem, inductive learning algo can't find unique soln?

i.e. one hypothesis can't be generalized to all examples.

[it is not practical to have all examples]

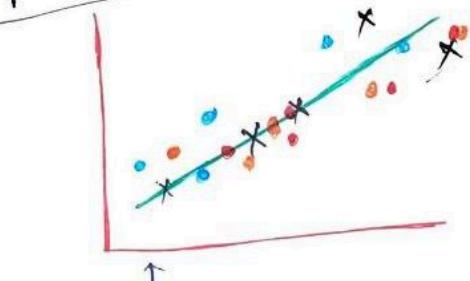
## \* Assumption in Inductive learning :

- A hypothesis  $h$  approximates  $f$  well in sufficiently large set of examples, will also approximate target fun well over unobserved examples.

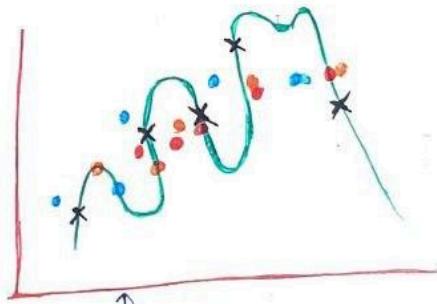
- Tasks in ML :-
- (a) come up with good hypothesis space ⑤
  - (b) find algo that works well in that hypothesis space
  - (c) generalizability in future data points
  - (d) Analyse confidence in result.
  - (e) Analyse computational tractability

[NN earlier wasn't computationally tractable.]

For Test data :-



It has more error



It has less error as curve fit to all points.

- similar performance on training, test data
- Underfitting didn't fit well to data = consistent performance
- overfitting: outstanding training data

[ almost same, similar type of error ] [ diff. Performance every time ]  
 ↑  
 to new data → inconsistent performance to testing data.

⑥

$T_1$  Testify data 2  $T_3$

Eg:-

$e_{\text{Green(st)}}$  3.4 2.9 3.1

$e_{\text{red(curve)}}$  1.7 5.7 10.9

↓

more error to curve which is more complicated.

Green Curve	Red curve
<ul style="list-style-type: none"> <li>High bias, underfitting</li> <li>but low variance</li> <li>consistent but poorer performance</li> </ul>	<ul style="list-style-type: none"> <li>Low bias, overfitting</li> <li>High variance, inconsistent performance across datasets.</li> </ul>

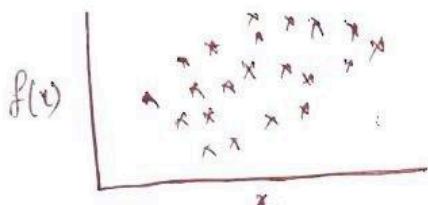
→ Bias = Avg. prediction of our model - correct value we want to predict.

- Model with high bias pays very less attention to training data and oversimplifies the model.

↑  
It leads to high error in training, test data.

- Variance is variability of model prediction across different sets of data.

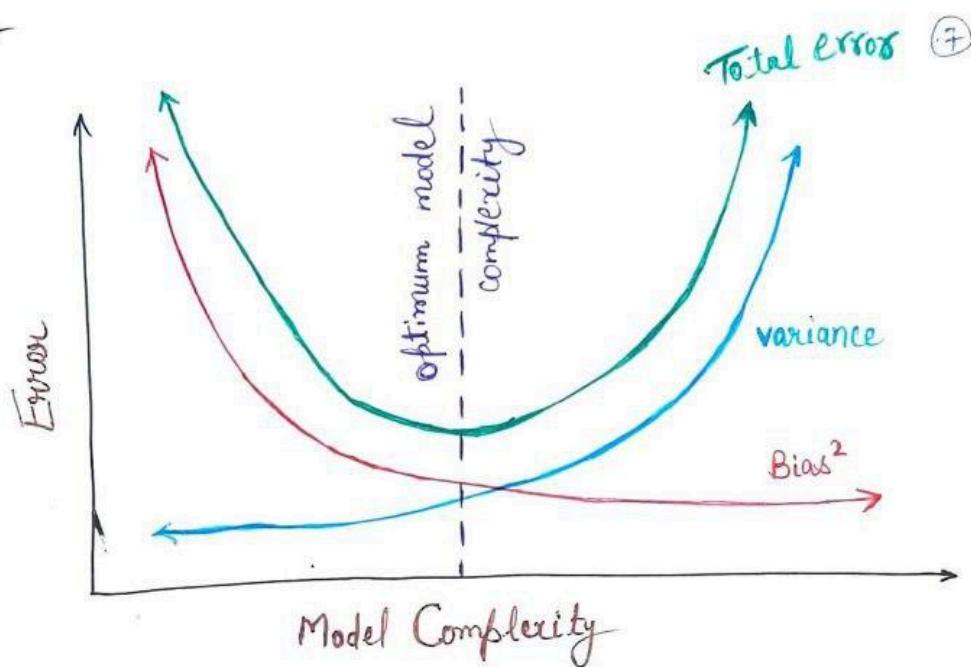
If we have a lot of training data of sufficient diversity



→ Training data of sufficient diversity

↓  
prevents overfitting

Errors :



- $\uparrow \uparrow$  model complexity  $\Rightarrow$   $\downarrow \downarrow$  Bias  
 $\uparrow \uparrow$  variance
- $\downarrow \downarrow$  model complexity  $\Rightarrow$   $\uparrow \uparrow$  Bias  
 $\downarrow \downarrow$  variance
- we need optimum model complexity  $\Rightarrow$   $\downarrow \downarrow$  Bias  
 $\downarrow \downarrow$  variance

→ keep lot of training data  
with sufficient diversity.

Date - 14/08/2024

## Probability

①

- Uncertainty: Real world is complex.

eg: there is a moment of particle, we can't find out absolute moment.  
Real world is uncertain.

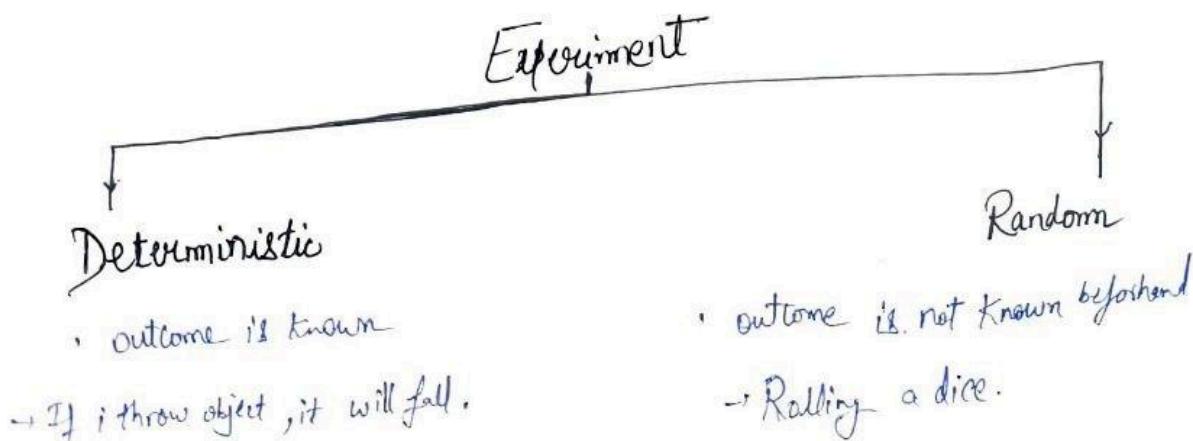
eg: Tossing of a coin = Certain

[Based on many factors like physics formulae, how much force gonna be on tossing of coin, calculate rotation of coin, properties of atmosphere etc  $\Rightarrow$  If I have too many parameters = I can predict right.]

but it is impossible to have microscopic details.

- In practice we often don't model exact system

$\hookrightarrow$  we choose some abstractions.  
 $\hookrightarrow$  result in the probabilistic nature of our experiments.  
use probability to model the system.



(2)

Sample Space :- Set of possible outcomes.

(a) Finite

e.g. :- Dice

$$S = \{1, 2, 3, 4, 5\}$$

(b) Infinite

e.g. Tomorrow's temperature

Amount of rainfall at Jodhpur in July

Event :- Any subset of the sample space.

e.g. if while rolling dice outcome is a no.  $< 3$

$$\text{Event } (A) = \{1, 2\}$$

Random Variable :- it is numerical description of outcome of a random experiment.

→ A fun<sup>n</sup> that assigns numerical values (real or Boolean) to each sample point.

Random Variable

Discrete

Continuous

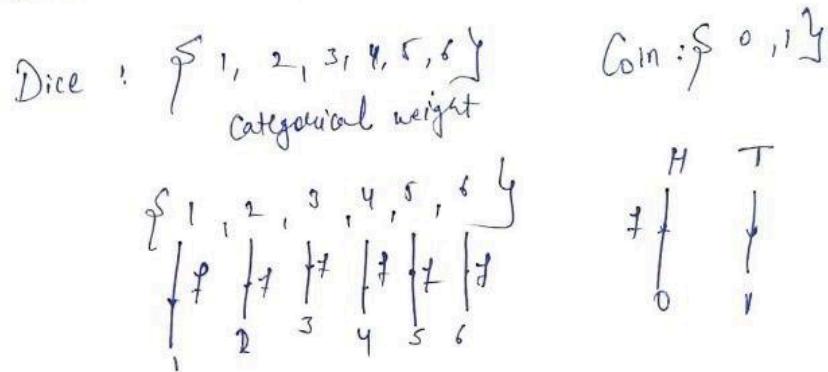
- Takes only countable no. of discrete values
- Sample space for weather cond<sup>n</sup> : {sunny, rainy, cloudy}

Takes uncountably infinite no. of possible values

e.g. Temperature at Jodhpur

$$\{67^\circ, 46.2^\circ\}$$

- R.V usually indicated by Capital letter. (eg:  $X$ ) (3)
- Outcome of Dice is not random variable. A fn 'f' assigns value to event.



### Probability of an Event :-

- $A$  = event
- Probability of event  $A$

$$P(A) = \frac{\text{No. of elements in set } A}{\text{No. of elements in sample space } S}.$$

$$P(A) = \frac{\text{No. of favourable outcome}}{\text{Total no. of outcome}}$$

$$P(\text{outcome} < 3) \rightarrow \text{event set} = \{1, 2\}$$

$$P(\text{outcome} < 3) = \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

$$\begin{aligned} P(\text{outcome odd}) &\quad \text{event set} = \{1, 3, 5\} \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2} \end{aligned}$$

Probability :- Sample space is not equally likely in realistic way. ①

Probability defn' applies if  
→ sample space is equally likely  
→ It should be Countable.

### Frequentist Approach of Probability

- If a experiment is done  $n$  times out of it event  $A$  occurs  $n(A)$  no. of times.
- Relative frequency of  $A$  is  $f_r(A) = \frac{n(A)}{n}$ .

for true probability outcome experiment should be done 'as no. of times.'

$P(A) = f_r(A)$  only if experiment done nearly '∞' times.

$$\boxed{\text{Probability of } A \text{ is } P(A) = \lim_{n \rightarrow \infty} f_r(A) = \lim_{n \rightarrow \infty} \frac{n(A)}{n}}$$

- Probability measure probability to an event.  
or

$P(\cdot)$  Probability function assigns probability to an event.

$P(A) = \text{chance that event } A \text{ occurs.}$

(5)

## Probability Properties

- $P(A) = \text{true no. i.e. } (0 \leq P(A) \leq 1)$
- if  $\emptyset = \text{null event, no outcome from an experiment, } P(\emptyset) = 0$   
It is impossible that experiment has no outcome.
- if  $S = \text{sample space}$   

$$P(S) = 1$$

probability that something happens is always 1.
- If  $A_1, A_2, \dots, A_n$  are countable sequence of disjoint events  

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n)$$

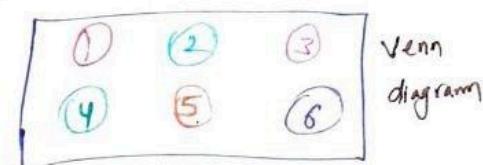
e.g. mutually exclusive of

$$\begin{aligned} P(2 \cup 4 \cup 5) &= P(2) + P(4) + P(5) \\ \text{getting either } 2, 4, 5 &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2} \end{aligned}$$

- Disjoint / Mutually Exclusive Two events are disjoint if they can't occur simultaneously.

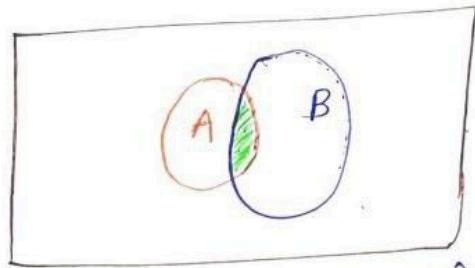
e.g. probability of getting either 2 or 4 or 5.

Rolling a dice.



## Probability + Addition Rule

6



$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

## Revisiting Random Variable

- $X$  is rv with  $k$  possible values  $x_1, x_2, \dots, x_k$

.  $P(X = x_j) = p_j$

Then  $\sum_{j=1}^k p_j = 1$

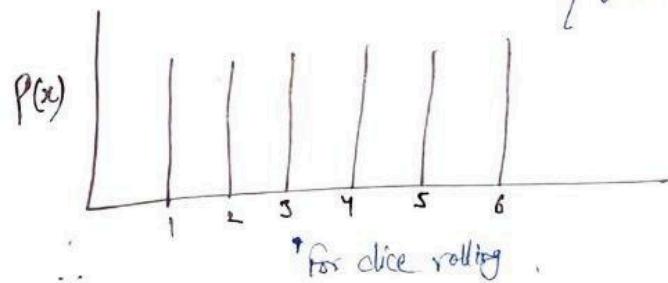
$\Rightarrow P(A \cup A^c) = P(\Omega) = 1$   
 $P(A) = 1 - P(A^c)$

- Probability distribution : A list containing  $p_1, p_2, \dots, p_k$

$$\left\{ \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6} \right\}$$

Often represented by

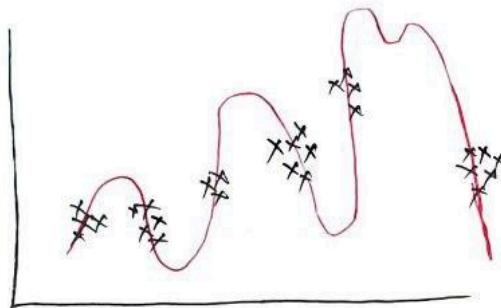
Histogram



Date : 16-08-2024

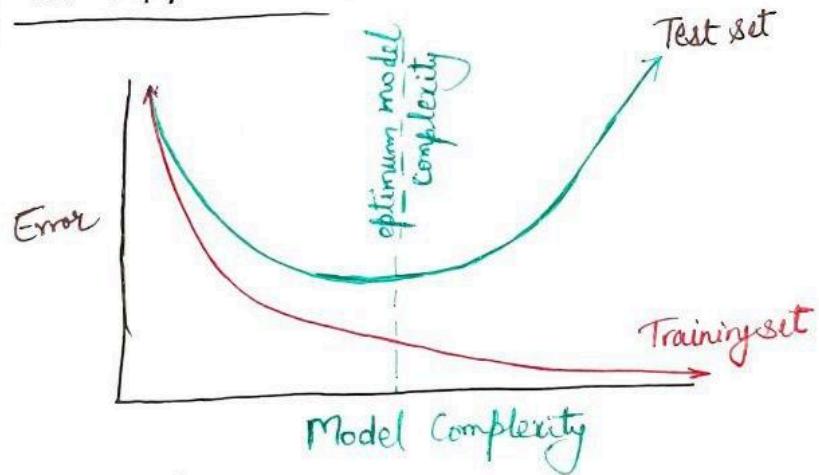
## Lecture - 6

①



- Now, we have a lot of training data but of similar types.

### Training Vs. Test set Error :



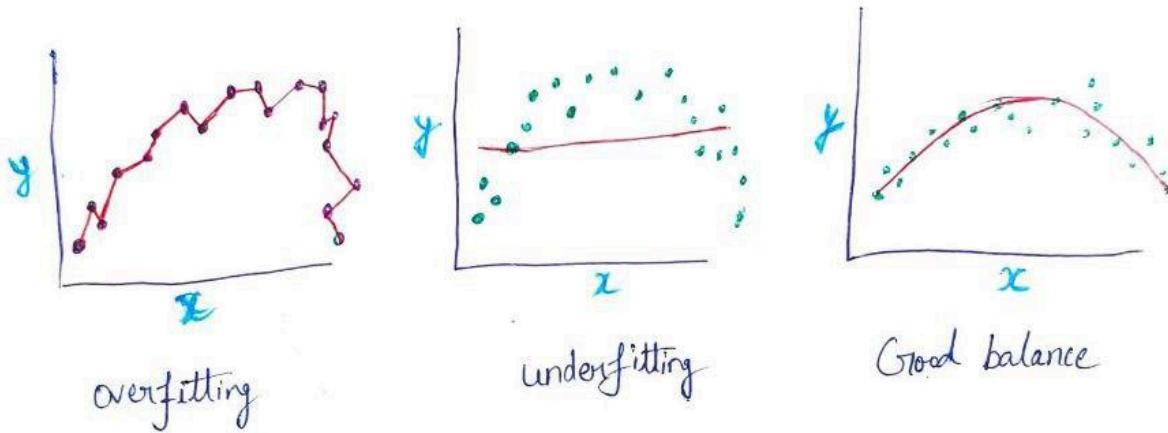
- For every line we can quantitatively calculate bias.  
Variance can't be calculated around multiple error.
- Qualitatively we can't say bias, variance.
- There is trade-off b/w bias and variance, based on my experience i need to choose reasonable limit.

Ideally we want low bias, low variance.

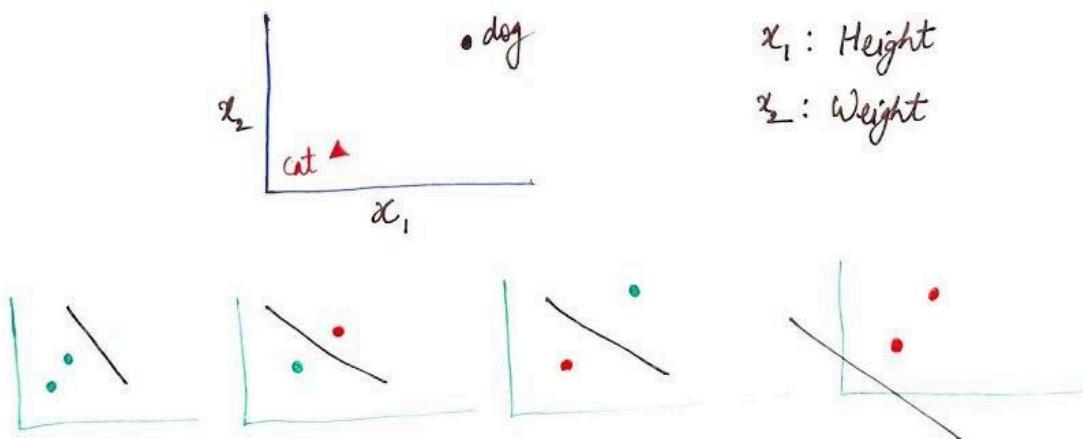
## Bias - Variance Tradeoff

(2)

- we want  $\downarrow \downarrow$  bias, variance.
- we need to find sweet spot b/w simple and complex model.



No. of Data points that a hypothesis can classify correctly :-



→ Hypothesis space of straight lines is expressive enough to classify two points in 2D plane.

③

Note

At least 1 configuration of 2 points st. line which can classify correctly.

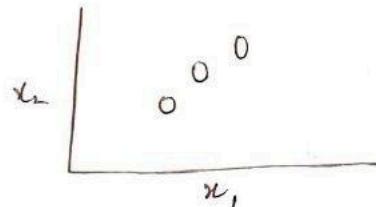
- Hypothesis space of st. line can shatter two points in a 2D plane.

Shattering :- A set of  $N$  points is said to be shattered by a hypothesis space  $H$  if:-

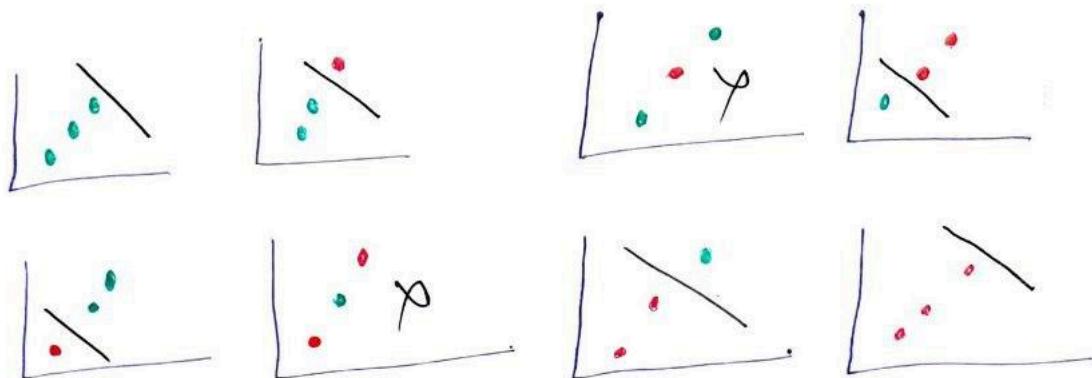


there are hypotheses in  $H$  that can separate +ve and -ve examples in all  $2^N$  possible ways for atleast one configuration of data points.

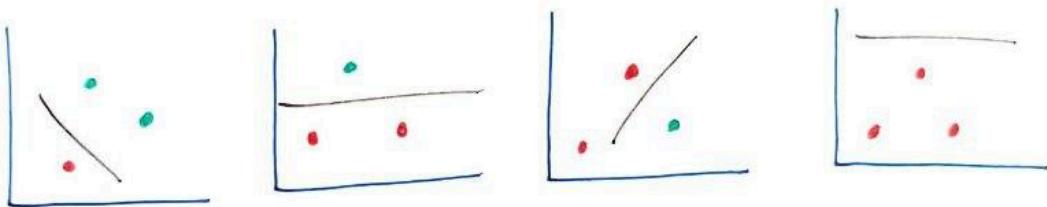
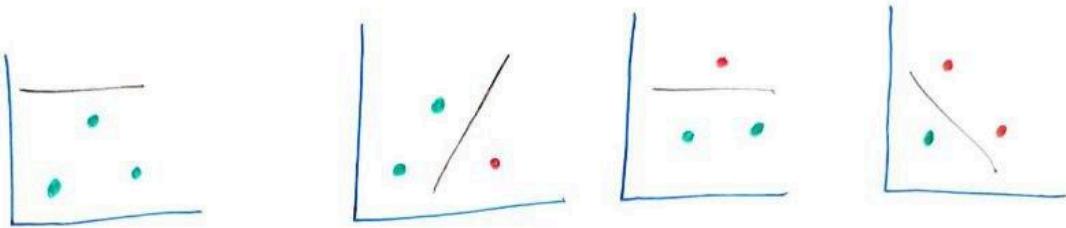
Three Points :-



Can st. line classify all possible combination of classes for 3 points?

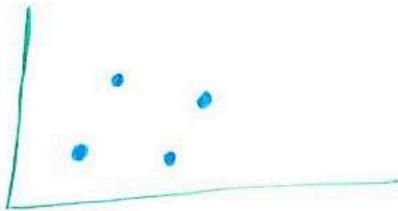


If 3 points are not co-linear → this is one of possible configurations ↗



Note: one classification by 1 configuration is okay.

4 Points ↗



Can st. line shatter any 4 points in 2D plane? No

Note :-

Maximum no. of points that can be shattered by st. line in 2D plane is 3.

eg:- 

A single 2D coordinate system with a horizontal blue line and a vertical blue line intersecting at the center. Three points are arranged in a triangle: top vertex is green, bottom-left vertex is red, and bottom-right vertex is green. A curved arrow points from the text "can be separated by st. line." to this diagram.

## Vapnik - Chervonenkis (VC-dimension) :

(5)

VC dimension of hypothesis space consisting of st. lines in 2D plane : 3.

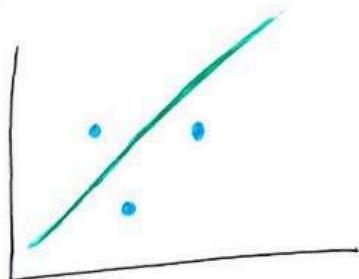
VC-dimension characterizes the expressive power or capacity of a hypothesis space.

- Maximum # of points that can be shattered by a hypothesis space is called VC dimension of the Hypothesis Space.

e.g. If classifcatn model can shatter D points atmost by changing its parameters. VC dimension of classification : D.

↑↑ expressive hypotheses = VC dimension ↑↑

### VC dimension

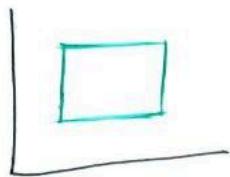


- Consider a classification model that creates a decision boundary →

$$y = mx + c$$

where  $m, c$  = parameters of model

By changing  $m, c$  of classification model we'll get diff hypothesis space.



VC-dimension of hypothesis (6)

: Space consisting of axis-aligned rectangle.

• Let's see classification model with VC-dimension D

- $N$  = no. of training data / samples
- $D$  = VC dimension

Then

$$P \left( \text{test error} < \text{training error} + \sqrt{\frac{1}{N} \left[ D \left( \log\left(\frac{2N}{D}\right) + 1 \right) - \log\left(\frac{n}{4}\right) \right]} \right) = 1 - \gamma$$

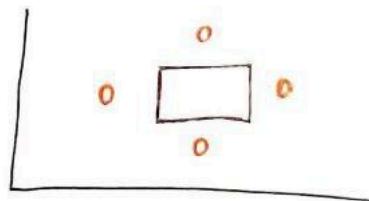
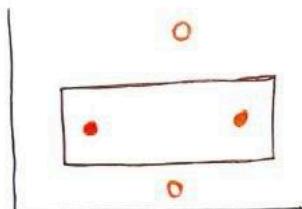
$$0 \leq \gamma \leq 1$$

This is applicable when  $D \ll N$ .

No. of training sample is way higher than VC-dimension.

If VC-dimension is high, model is more expressive.

- More complex  $\rightarrow$  more chance of  $\rightarrow$  ↑↑ test error.  
overfitting



VC=4

- We want same error for Testing, training. (7)

$\eta$  is Bound to train, test error

$$P\left(\frac{\text{test error} < \text{training error}}{\eta} + \sqrt{\frac{1}{N} \left[ D \left( \log\left(\frac{2N}{D}\right) + 1 \right) - \log\left(\frac{\eta}{4}\right) \right]} = 1 - \alpha\right)$$

Case 1: If  $\eta \downarrow \downarrow \downarrow$   $1 - \alpha = 1 \Rightarrow \text{Probability} = 1$

Probability of test error < training error = 1 confidence ↑↑

$\eta$  should be minimum.

Case 2:

if VC dimension is ↑↑  $\Rightarrow$  test error becomes ↑↑  
can't bound test error

if VC dimension ↑↑ = model is expressive  
↓  
more complex

→ more complex means  $\rightarrow$  more chance of overfitting.

So, higher test error.

⑧

## Challenges in Data for ML :-

- Good ML models can't be designed without good quality data.

Domain Knowledge is needed to handle such data.

- Various problems with most practical data →
  - Noise
  - Missing data
  - Duplicate data
  - Inconsistent data etc.
- Data Cleaning is done to →
  - Reduce noise
  - handle duplicate data
  - handle missing data
  - handle inconsistent data
  - May require domain knowledge to deal with these problems.

Date : 20/08/2024

## Lecture - 7

P23CS0010 ①

### Challenges in Data for ML

(i) Missing data eg :-

Student id

10012  
10013  
10024  
[ ]  
10098

(ii) Inconsistent data

Classification

Junior  
Senior  
Sophomore  
~~Sr~~  
Senior  
Jr

ID

2102  
2214  
AUC8  
2308

→ id starting with letter.

Age

21  
[200]  
24  
23  
[2.8]  
30

(iii) Noisy data

$$\boxed{\text{Quality of Model} = \text{Quality of data}}$$

↓  
performance of model depend on data

How to clean data, fill data?

Q:- Pattern can only be identified correctly by correct data.

① While coming up with a hypothesis :—

②

Step 1 :- Clean data

Step 2 Normalization

Step 3+ Coming up to a best hypothesis

Step 4+ Error Minimization

Step 1+ Clean Data + Data cleaning involves :

- Reducing noise
  - Handling duplicate data
  - Handling inconsistent data
  - Handling missing values
- } • Domain knowledge is required to deal with these problems.

e.g.: If a person has Hemoglobin = 30 ← incorrect data

↑  
need domain knowledge to fill/deal with data.

Dealing with cleaning of data

(i) Dealing with missing values :

• Data deletion : results in loss of data

• Imputation : predict data, through informed guess, by using regression.

(3)

## (ii) Dealing with noise :

(a) Binning :- Divide total data into bins.

Data →  $\boxed{\text{Data}} / \boxed{\text{Bin}}$

- replace value of noise data based on some characteristic of bin (mean, avg, min, max).

(b) Regression : based on other parameters → predict

(c) Outlier analysis : treat them separately or eliminate those.

Scale features :- suppose I want to find abnormalities:-

eg Avg Hb : 13500 mg/dL } Avg values = normal level  
 Avg Creatine : 0.7 mg/dL }

suppose if a person has

Hb	: 12500 mg/dL
Creatine	: 1.7 mg/dL

$Hb$   $13500 - 12500 = 1000$  based on huge change in  
 $Creatine$   $1.7 - 0.7 = 1$  Hb compared to Creatine  
 we might conclude person has  
 Hb disorder

→ This wrong conclusion may happen because:  
↓  
Normal values are different range.

- We want to bring different features at same range.
  - If we don't normalize
    - ↓  
features with lower value range may lose their meaning when combined response is created.
  - Large spread in feature value may mislead the ML model.
    - may mislead ML model
    - lower feature may lose their meaning when collectively taken.
- Normalization is solution
- data standardization
  - Min-max normalization

## (i) Data Standardization : Z-Score normalization ⑤

. Std<sup>n</sup> at input :- create feature value with mean=0,  
variance = 1/unit

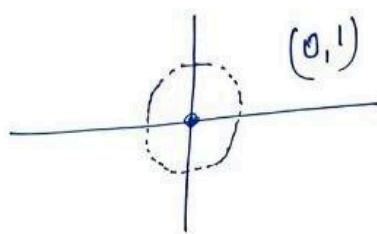
- . n data points.
- . feature values for n data points :  $f_1, f_2 \dots f_n$
- . Standardized feature value:

$$f'_i = \frac{f_i - \mu_f}{\sigma_f}$$

$\mu_f$  = mean of  $f_1, f_2 \dots f_n$

$\sigma_f$  = std of  $f_1, f_2 \dots f_n$

- . Std. value will have (mean 0, unit variance)



## (ii) Min-Max Normalization :-

- . n data points
- . feature values for n data points :  $f_1, f_2 \dots f_n$
- . normalized feature value

$$f'_i = \frac{f_i - \min(f_1, f_2 \dots f_n)}{\max(f_1, f_2 \dots f_n) - \min(f_1, f_2 \dots f_n)}$$

(6)

$$f'_i = \frac{f_i - \min(f_1, f_2, \dots, f_n)}{\max(f_1, f_2, \dots, f_n) - \min(f_1, \dots, f_n)}$$

$\max(f_1, f_2, \dots, f_n) - \min(f_1, \dots, f_n)$  = denominator

$f_i - \min(f_1, \dots, f_n)$  = numerator

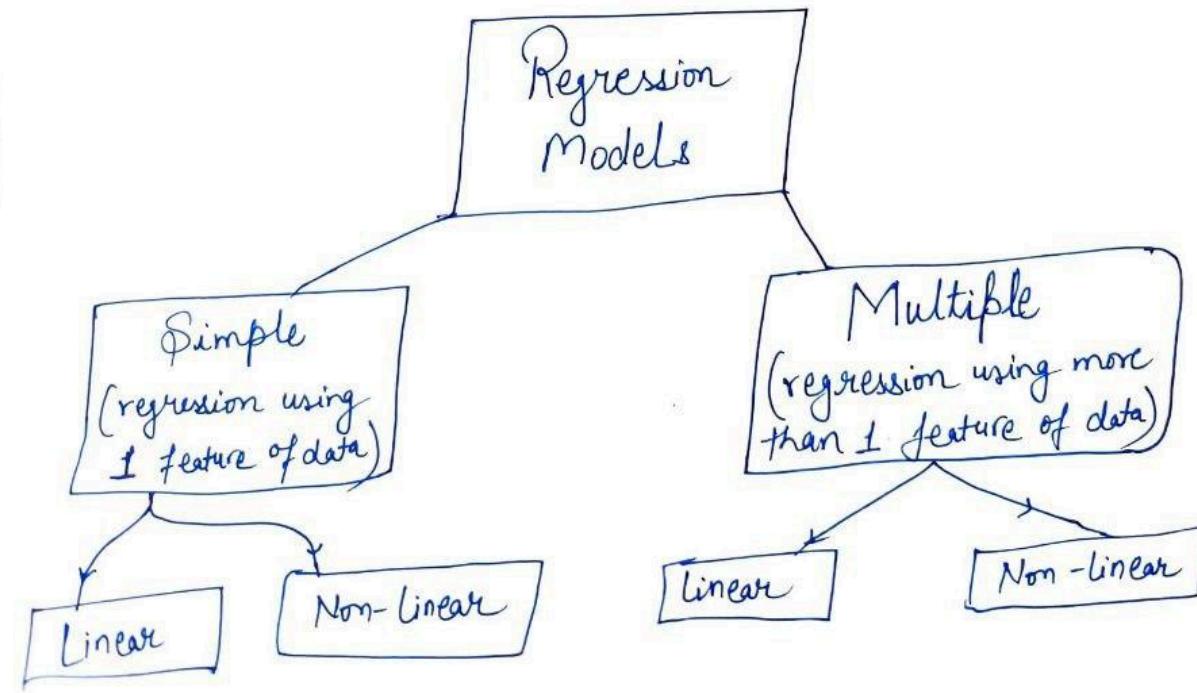
deno.  $\geq$  numerator

Q :- Given Hb values, find min-max normalization ↴

Hb
11.7
13.4
14.5
12.2

$$f'_1 = \frac{11.7 - 11.7}{14.5 - 11.7} = 0$$

$$f'_2 = \frac{13.4 - 11.7}{14.5 - 11.7} = 0.6$$

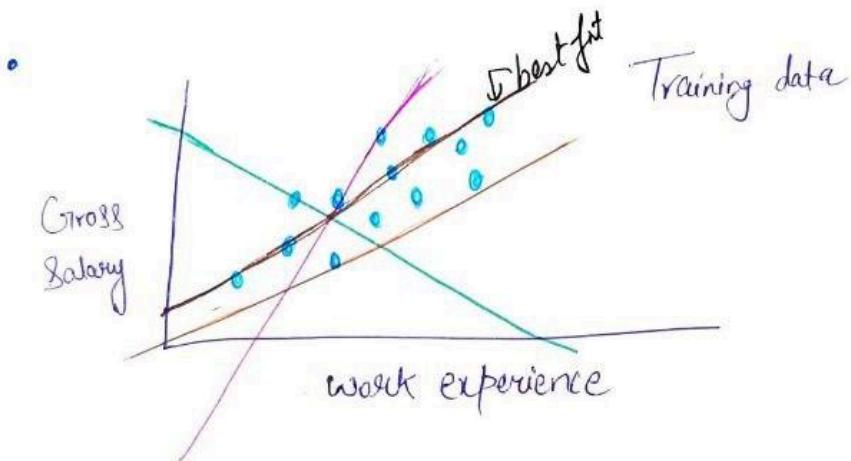


Regression : eg: predicting price of house.

(7)

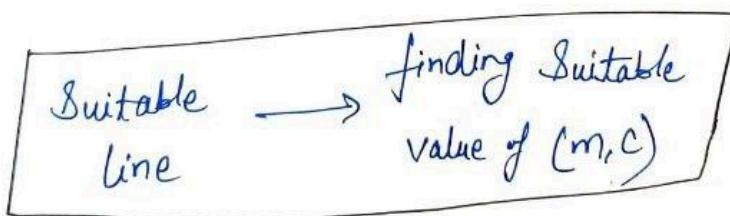
- ML introduces Inductive Bias:
  - ④ restrictive
  - ② Reference

e.g. restrictive inductive bias  
is used for St. line.



- Straight line best fits on given data.

St. line  $\Rightarrow$  explain more points in good way.



- eq<sup>n</sup> of St. line

$$y = mx + c$$

parameter of eq<sup>n</sup> / regression model

- different values of  $(m, c)$   $\rightarrow$  different straight lines

(8)

- Different st. lines, which to prefer?

use **Error**

[Brown line has least total error]

- find st. line ( $m, c$  parameter) that minimizes total error in training data.

For multiple Linear Regression ↗

$$y = m_1 x_1 + m_2 x_2 + \dots + m_p x_p + c$$

parameters of eqtn / regression model

- find st. line ( $m_1, m_2, \dots, m_p$  and  $c$ ) that minimizes total error in training data.

e.g.  $y$  = credit score depend on multiple factor:

$x_1$  = how many loan

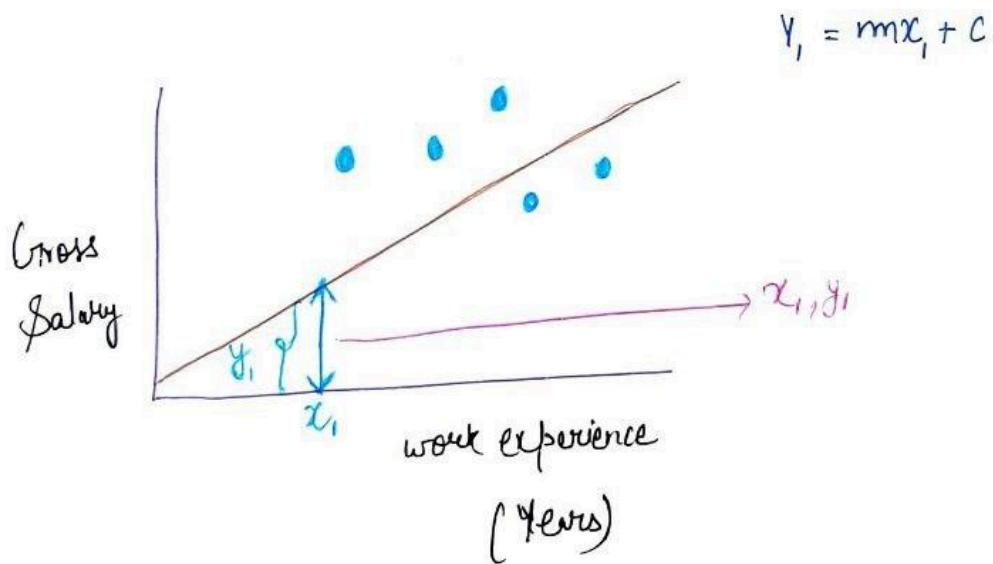
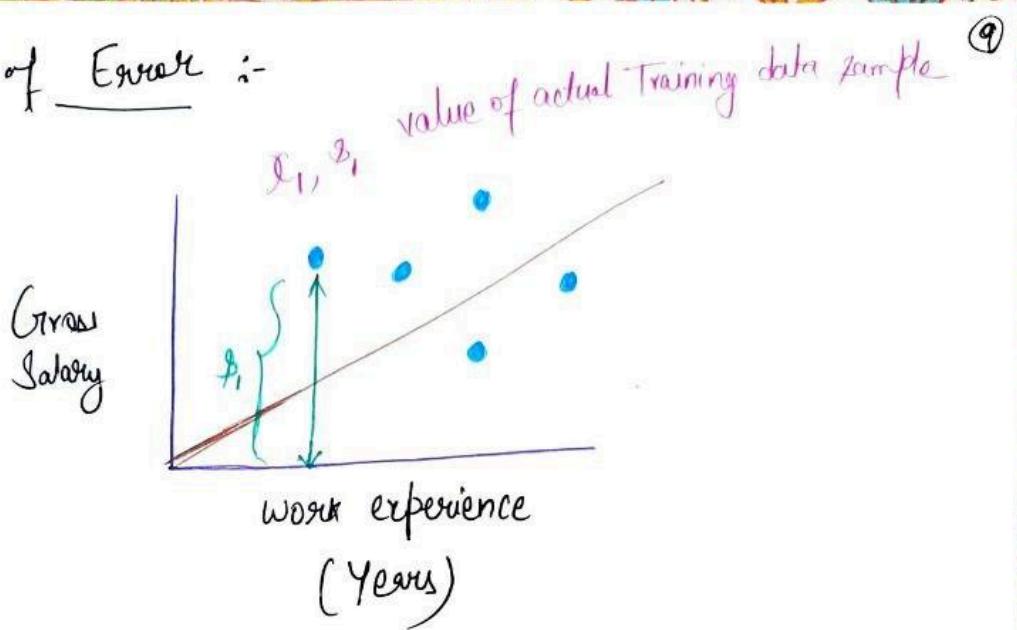
$x_4$  = #. of experience in bank account

$x_2$  = how many bank account

etc.

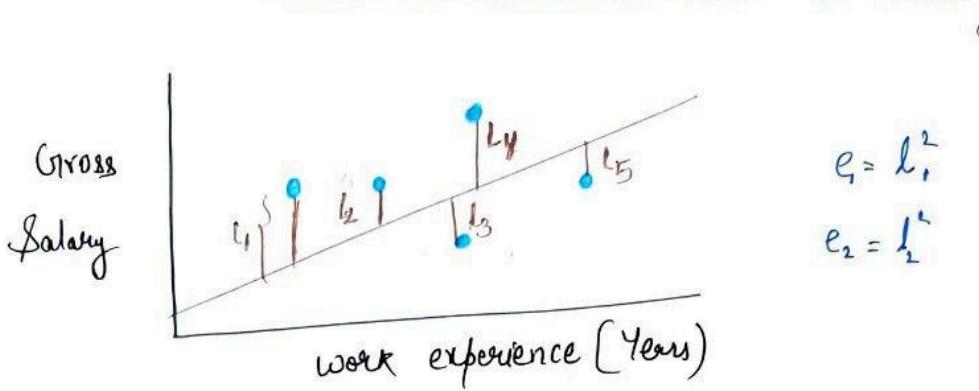
$x_3$  = previous installment history

Use of Error :-



Error for this data point w.r.t prediction

$$\begin{aligned}e_1 &= (y_1 - \hat{y}_1)^2 + (x_1 - \bar{x})^2 \\&= (y_1 - \hat{y}_1)^2\end{aligned}$$



$$E = e_1 + e_2 + \dots + e_N$$

- We can calculate error for training data point and sum them upto get total error.

$$E = (y_1 - s_1)^2 + (y_2 - s_2)^2 + \dots + (y_N - s_N)^2$$

$$E = \sum_{n=1}^N (y_n - s_n)^2$$

- Goal :- To find a hypothesis such that error  $E$  is minimized

$$J(m, c) = \frac{1}{2} \sum_{n=1}^N (y_n - s_n)^2$$

[find  $m, c$  such that above is minimized]

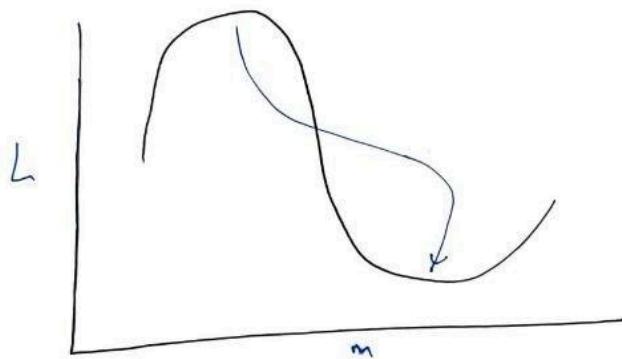
- minimize squared error = Least-square regression.

(1)

$$J(m, c) = \frac{1}{2} \sum_{n=1}^N (m x_n + c - y_n)^2$$

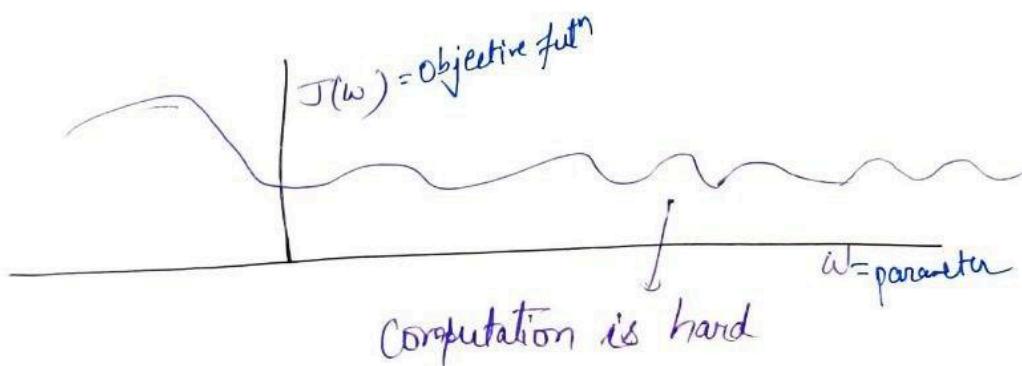
↓  
minimize Square regression

How to minimize  $J(m, c)$  :-



Goal :- of training is to get  $m, c$   
↑  
to minimize objective function

- goal is to reach global minima
- Double derivative finds min, max values.



## Lecture -8

①

### Revisiting Random Variable

- $X$  is a discrete rv with  $K$  possible values  $x_1, x_2 \dots x_K$ .

eg:-  $X$  contain  $\{0, 1\}$  while tossing coin.

$$P(X=x_j) = p_j \quad \text{Then} \quad \sum_{i=1}^K p_i = 1$$

- Probability Distribution: list containing  $p_1, p_2 \dots p_K$

↳ often represented as histogram.

Cond'n :- Sample Space

finite	T
Equally likely	

- If probability isn't equally likely ? Frequentist approach

do experiment 'os' no. of times.

(It's discrete)

Skewed off if experiment is less no. of times.	T
Equally Probable :- if experiment done many no. of times.	

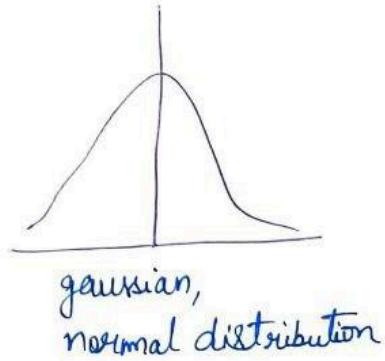
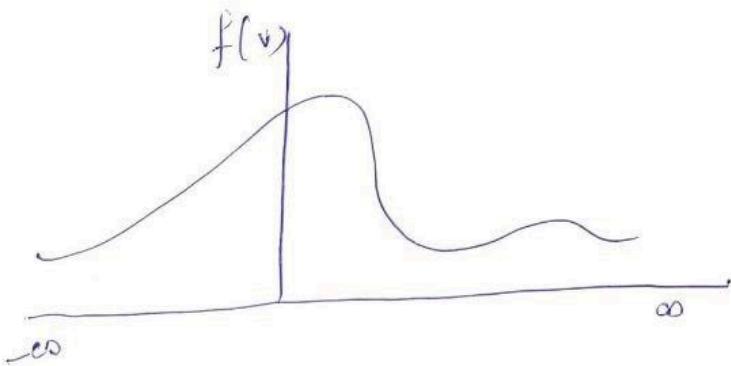
(b) If S.S is continuous :- for probability density function (3)  
 probability range is needed.

$$P(a < X \leq b) = \int_a^b f(v) dv$$

- If  $X$  is continuous r.v:-  
 probability of value of  $X$  are represented by a curve  $f(v)$ .
  - $f(v) \geq 0$  at all points.
  - Area under curve = 1.

$$\int_{-\infty}^{\infty} f(v) dv = 1$$

Probability range is needed for continuous data.



- In ML we want to learn  
 distribution of data.

## Conditional Probability

(3)

In a box, Samsung phones = 40, MI phones = 20

[Out of these 10 Samsung, 2 MI phone are defective/not working]

↳ You pick up a phone and find probability of (phone isn't working) :- ?  
↳ Is it of Samsung Company ?

Given :- picked up phone is not working.

Now - if so probability of it's being Samsung ?

$$P\left(\frac{SP}{NW}\right) = \frac{\# \text{ Samsung phone that are not working}}{\# \text{ phones that aren't working}}$$

$$= \frac{10}{12} \quad \text{multiplying, } \div \text{ by } 40, 60$$

$$= \frac{\frac{10}{40} \times \frac{40}{60}}{\frac{12}{60}}$$

$\frac{10}{40}$  : Given Samsung phone that is NW

$\frac{40}{60}$  : Probability of SP in box

$\frac{12}{60}$  : prob. of NW phones

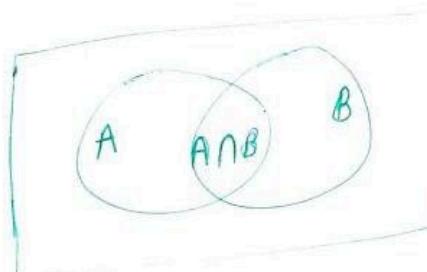
②

$$P\left(\frac{SP}{NW}\right) = \frac{P(NW \neq P) \cdot P(SP)}{P(NW)}$$

$$P\left(\frac{SP}{NW}\right) = \frac{P(SP \cap NW)}{P(NW)}$$

Conditional probability ↗

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P\left(\frac{B}{A}\right) \cdot P(A)}{P(B)} \xrightarrow{\text{Prior probability}}$$

Baye's TheoremConditional/  
Posterior  
probability

$$P\left(\frac{A}{B}\right) = \frac{P\left(\frac{B}{A}\right) \cdot P(A)}{P(B)} \xrightarrow{\text{prior}} \begin{array}{l} \text{indicate my} \\ \text{belief about} \\ \text{occurrence} \\ \text{of } A \end{array}$$

↓  
evidence  
that probability that B occurs

②

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause}) P(\text{Cause})}{P(\text{Effect})}$$

- If a system with causes  $C_1, C_2$  and  $e_1, e_2$

Q. If we want to find probability of diff. causes given effect  $e_2$ ?

$$P(C_1|e_2) = \frac{P(e_2|C_1) P(C_1)}{P(e_2)}$$

$$P(C_2|e_2) = \frac{P(e_2|C_2) P(C_2)}{P(e_2)}$$

Since  $e_2$  may be caused either by  $C_1, C_2$

$P(C_1|e_2) + P(C_2|e_2) = 1$

$\text{Posterior} = \frac{\text{Prior} * \text{Likelihood}}{\text{Evidence}}$

Similarly  $P(C_1|e_1) = \frac{P(e_1|C_1) P(C_1)}{P(e_1)}$

$$P(C_2|e_1) = \frac{P(e_1|C_2) P(C_2)}{P(e_1)}$$

- Evidence are fewer

(6)

also  $P\left(\frac{C_1}{e_1}\right) + P\left(\frac{C_2}{e_1}\right) = 1$

Posterior & likelihood \* prior

e.g. given i get Cough/fever (A)  $\rightarrow$  A doesn't depend on B.  
There is power cut (B)

$P(\text{fever})$  = only depends on A.

Q1- If among two person what is probability of having birthday  
 on same month ?

$$P(A) = \frac{1}{12} \quad P(B) = \frac{1}{12}$$

$$P(A) \cdot P(B) = \frac{1}{12} * \frac{1}{12}$$

If A, B born in Jan =  $\frac{1}{144}$

in Feb =  $\frac{1}{144} - - -$

$$\begin{aligned} P(X) &= P(\text{Jan}) + P(\text{Feb}) - - - + P(\text{Dec}) \\ &= \frac{1}{144} + \cdot - - - - \frac{1}{144} = \frac{12}{144} = \frac{1}{12} \end{aligned}$$

• Since  $e_1$  may be caused by  $C_1$  or  $C_2$ .

(7)

$$P(C_1/e_1) + P(C_2/e_1) = 1$$

• For both  $P\left(\frac{C_1}{e_1}\right)$  and  $P\left(\frac{C_2}{e_1}\right)$  denominator is same :-

Let's consider  $\frac{1}{P(e_1)} = \alpha$

Then  $P(C_1/e_1) = \frac{P(e_1/C_1) P(C_1)}{\alpha} = \alpha P\left(\frac{e_1}{C_1}\right) P(C_1)$

$$\therefore P(C_2/e_1) = \alpha P\left(\frac{e_1}{C_2}\right) P(C_2)$$

Posterior  $\propto$  (likelihood  $\times$  prior)

Independent Events :- Two events A, B are independent if

$$- P\left(\frac{A}{B}\right) = P(A) \quad (A \text{ doesn't depend on } B)$$

$$P\left(\frac{A}{B}\right) = \frac{P(A \text{ AND } B)}{P(B)} = P(A)$$

$$- P\left(\frac{B}{A}\right) = P(B)$$

$P(A \text{ AND } B) = P(A) \cdot P(B)$

## Lecture - 9

( 22 - 08 - 2024 )

①

### How to reach Global Minima : 2 D ex. :-

Important point to note :-

[ if superscript : sample  
 ] if subscript : variables

e.g.:-  $x^1, x^2$  1st person, 2nd person ... so on.

$x_1, x_2, \dots$  :  $y = mx + c$  = independent data points.

- Minimization can be done in by many approaches

↓  
best is Gradient Descent

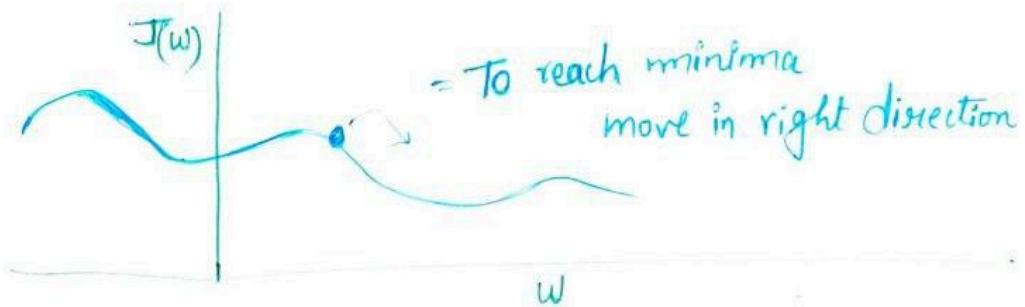
↓  
minimizing  $(m, c) \Rightarrow$  changing  $(m, c) =$  diff. st. lines

$$J(m, c) = \frac{1}{2} \sum_{n=1}^N (mx^n + c - s^n)^2$$

minimize  $J(m, c)$

- if multiple minimization of linear regression :-

$$y = m_1x_1 + m_2x_2 + \dots + m_p x_p + c$$



if objective fun" in right side = error decreases.

- moving in right side :-

$\therefore$  [ increasing value of  $w$   
reduces value of  $L(w)$  ]

$$\frac{\partial J(w)}{\partial w} < 0$$

$$\text{It concludes as :- } J(w) \leftarrow \begin{cases} \partial J(w) & \downarrow \\ \partial w \uparrow & = +ve \end{cases}$$

$$\left( \frac{\partial J(w)}{\partial w} - re \right) \text{ overall its } -re$$

moving towards right  $\rightarrow$  need new  $w$  value.

$$\omega_{\text{new}} = \omega_{\text{old}} + \eta (\Delta \omega) \xrightarrow{\substack{\uparrow \\ \downarrow \\ \omega_{\star}(\text{-ve})}}$$

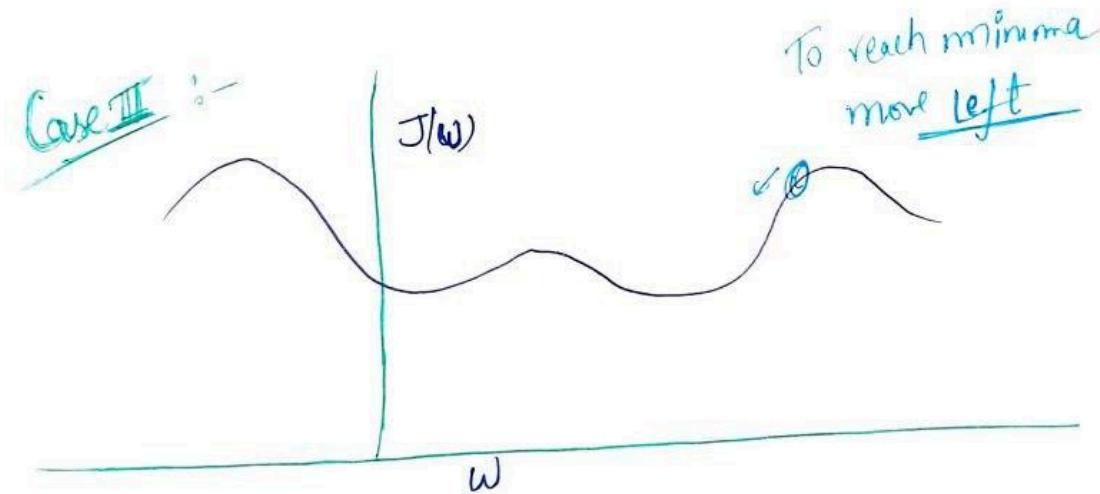
$$w_{\text{new}} = w_{\text{old}} + \eta (\Delta w)$$

(3)

$\frac{\partial J(w)}{\partial w} = -\text{ve} \Rightarrow$  to make it +ve, multiply by -ve value.

$$w_{\text{new}} = w_{\text{old}} - \eta \underbrace{\left( \frac{\partial J(w)}{\partial w} \right)}_{\text{overall +ve}} \rightarrow -\text{ve}$$

$\eta$  = +ve constant



To decrease  $w$  :-

$$w \leftarrow w - \eta \frac{\partial J(w)}{\partial w}$$

$$\left( \frac{\partial J(w)}{\partial w} \rightarrow -\text{ve} \right) \rightarrow +\text{ve}$$

. After multiple checking of  $\eta$  values, testing on:-

$\eta = 0.1, 0.2 \dots$  we conclude at particular value of  $\eta$ , we are getting less error, best result.

$\eta$  = learning rate, generally = 0.01

How to decide  $\eta$ ? Test performance of model  
on validation data.

How to know minima is reached: Training error will reduce.

e.g I may fix epoch 1000 times and i assume minima will be reached and return particular value.

$$\begin{aligned}x^n &= n^{\text{th}} \text{ data point} \\x_p &= p^{\text{th}} \# \text{ of variable}\end{aligned}$$

$x_1, x_2, \dots, x_p = P \# \text{ of parameters}$

• we is updating at every step.

(5)

### Gradient Descent ↴

- (i) Initialize parameter  $w$
- (ii) loop until convergence
  - (a) compute gradient  $\frac{\partial J(w)}{\partial w}$

$$(b) \text{ update parameters } w \leftarrow w - \eta \frac{\partial J(w)}{\partial w}$$

- (iii) Return parameters

$\eta$  = hyperparameter whose value should be set using validation performance.

### for many parameters ↴

- (i) Initialize parameter  $m = \{m_1, m_2, \dots, m_p\}$

- (ii) loop until convergence

$$(a) \text{ Compute gradient } \frac{\partial J(m)}{\partial m}$$

$$(b) \text{ update parameters } \theta_i \leftarrow \theta_i - \eta \frac{\partial J(m)}{\partial m_i}$$

- (iii) return parameters.

$$J(m) = \frac{1}{2} \sum_{n=1}^N (y^n - s^n)^2$$

find expression for update ↴

$$\theta_i \leftarrow \theta_i - \eta \frac{\partial J(m)}{\partial m_i}$$

$$m_i \leftarrow m_i - \eta \sum_{n=1}^N (y^n - s^n) x_i^n$$

(6)

Gradient descent for all  
training samples  $\rightarrow$  Batch G.D

Challenges  $\rightarrow$  As it computes gradient over all data, computing  
gradient over millions of training data is

- computationally hard
- need huge storage

In Batch G.D :- N training samples (dat points)

$\uparrow$   
we compute total objective function considering all N training  
samples.

$$J(m, c) = \frac{1}{2} \sum_{n=1}^N (m x_n + c - s_n)^2$$

$$\frac{\partial J(w)}{\partial w} = \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial w} (m x_n + c - s_n)^2$$

$\rightarrow$  for each update : it will take long time if we have huge  
training samples

$\rightarrow$  so we need to compute and store N gradients for each update.

## Batch GD $\rightarrow$ Stochastic Gradient Descent (7)

- In Batch GD with  $N$  samples:

$$J(m, c) = \sum_{n=1}^N \frac{1}{2} (mx_n + c - y_n)^2 = \sum_{n=1}^N J_n(m, c)$$

Contra: If I calculate objective function for 1 sample and do parameter update.  
 Do same process for 2nd points. [For Batch GD: Need to wait till  $n$  data points]

Stochastic GD: each update would require only one gradient computation

less time  
less storage for gradient values

e.g.: Going from jodhpur to mehrangarh fort

multiple paths are there

reach at particular point  $\rightarrow$  again discuss  $\rightarrow$  decide.

[In particular step = will take significant time,  
 Collective opinion  $\rightarrow$  lots of Noise]

going by discussion  
 based decision has as diff person has  
 diff opinions.

Q

SGD →

1. Initialize parameters
2. Loop until convergence
  - (i) Randomly shuffle samples in training dataset
  - (ii) For training sample  $n=1, 2, \dots, N$ 
    - (a) Compute gradient  $\frac{\partial J_n(w)}{\partial w}$
    - (b) Update parameters  $w \leftarrow w - \eta \frac{\partial J_n(w)}{\partial w}$
3. Return parameters.

Challenges in Stochastic and Batch GD →

- SGD → not data efficient when data is very similar.
- BGD → high computational time for each step.

SGD → not computationally efficient

- noisy results since at a time, one data point is considered.

Solution → Minibatch GD.

Minibatch GD → neither to

single points

nor to whole group.

It is more likely to be correct., less noisy, fast moment.

## Minibatch GD

(i) Initialize parameters

⑦

(ii) Loop until convergence

(a) Randomly shuffle samples in training dataset

(b) for training dataset, create  $b$  minibatch of size  $n$

(c) for each minibatches  $n=1, 2 \dots b$

- Compute gradient  $\frac{\partial J(w)}{\partial w}$

- Update parameters  $w \leftarrow w - \eta \frac{\partial J(w)}{\partial w}$

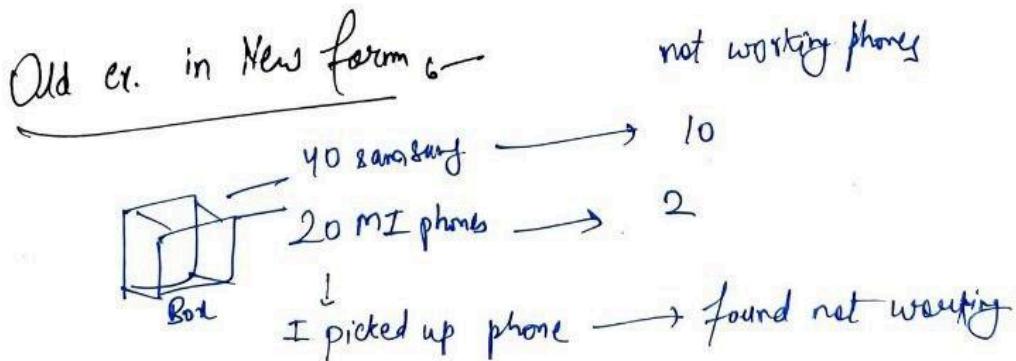
(iii) Return Parameters

How to initialize  $w$

can be all zeros?

can be random?

. Performance depends on initialization.



Tell whether phone is Samsung or MI?

Find  $P\left(\frac{SP}{NW}\right)$  and  $P\left(\frac{MIP}{NW}\right)$ ?

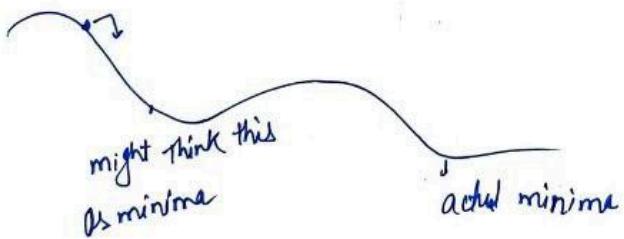
If  $P\left(\frac{8P}{NW}\right) > P\left(\frac{M1P}{NW}\right)$  (10)  
 $\Rightarrow$  i had picked Samsung phone.  
 else MI

$$P\left(\frac{M1P}{NW}\right) = \frac{\# \text{ MI phones that are not working}}{\# \text{ phones that are not working}} = \frac{2}{12}$$

$$P\left(\frac{8P}{NW}\right) = \frac{10}{12}$$

$P\left(\frac{8P}{NW}\right) > P\left(\frac{M1P}{NW}\right) \Rightarrow$  concludes that i picked up Samsung phone

Q: Eventual minima is not effected by linear regression?



In linear regression: gradient is constant because relationship is a straight line. This means it always leads directly to same minimum — point where error is smallest regardless of where we start.

In curves: gradient changes as we move along curve, which lead to different paths and potentially different minima depending on where we start.

Baye's Decision Theory :-Joint Probability Distribution :

Consider 2 random variables  $(X, Y)$ :

$X$ : Corresponding to weather { sunny, rainy, cloudy }

$$P(X) = \{0.6, 0.1, 0.3\}$$

$Y$ : Corresponding to power cut { power cut, no power cut }

$$P(Y) = \{0.15, 0.85\}$$

A joint probability distribution of  $X$  and  $Y$  (pairs considered)  
 ↓  
 probability distribution on all possible pairs of qps.

$$\therefore 3 \times 2 = 6 \text{ possible cases}$$

- $\sum$  of all possible outcome = 1

$\rightarrow 3 \times 2$  matrix of values

- Joint probability is necessary in daily life.

e.g. AC temp. depend on weather, our mood etc.

Chain Rule :-

(2)

$$P(A \cap B) = P(A \cap B) = P(A|B) P(B)$$

If  $A_1, A_2, \dots, A_n$  are n events, then

$$\bullet P(A_n \cap A_{n-1} \cap \dots \cap A_1) = P\left(\frac{A_n}{A_{n-1} \cap \dots \cap A_1}\right) P(A_{n-1} \cap \dots \cap A_1)$$

Similarly :-

$$P\left(\frac{A_n}{A_{n-1} \cap \dots \cap A_1}\right) P\left(\frac{A_{n-1}}{A_{n-2} \cap \dots \cap A_1}\right) P(A_{n-2} \cap \dots \cap A_1)$$

e.g. ... so on

If 4 boxes :-

$$P(A_4 \cap A_3 \cap A_2 \cap A_1) =$$

$$P\left(\frac{A_4}{A_3 \cap A_2 \cap A_1}\right) P\left(\frac{A_3}{A_2 \cap A_1}\right) P\left(\frac{A_2}{A_1}\right) P(A_1)$$

(3)

	Fever		$\neg$ Fever	
	Cough	$\neg$ Cough	Cough	$\neg$ Cough
Covid	0.21	0.10	0.11	0.08
$\neg$ Covid	0.11	0.07	0.09	0.23

$$P(\text{fever}) = \frac{\text{cough - covid possible pairs}}{\text{all exhaustive combination of variable}}$$

$$P(\text{fever}) = 0.21 + 0.10 + 0.11 + 0.07 = 0.49$$

$$\begin{aligned} P(f \vee c) &= P(f) \vee P(c) \vee P(f \wedge c) \\ &= 0.21 + 0.11 + 0.10 + 0.07 + 0.11 + 0.08 = 0.68 \end{aligned}$$

$$P(\neg \text{covid} | \neg \text{fever}) = \frac{P(\neg \text{covid} \cap \neg \text{fever})}{P(\neg \text{fever})} = \frac{0.09 + 0.23}{0.11 + 0.09 + 0.08 + 0.23} = 0.627$$

Conditional probability :- tells joint prediction

if know whole distribution  $\downarrow$  can answer any query

If  $P\left(\frac{8P}{NW}\right) = P\left(\frac{M1}{NW}\right) \Rightarrow$  can't make any decision (y)

Sample space corresponding to joint distribution is →

$$S_J = \{ (s, pc), (s, npc), (u, pc), (u, npc), (c, pc), (c, npc) \}$$

eg 3r if  $P\left(\frac{SP}{NW}\right) > P\left(\frac{MI}{NW}\right)$  : phone picked up is Samsung.  
 if  $P\left(\frac{MI}{NW}\right) > P\left(\frac{SP}{NW}\right)$  : picked up phone is MI

Eg 4: Let's say 2 class classification problem with  $C_1, C_2$

$c_1 = \text{cat}$ ,  $c_2 = \text{tiger}$

If i have <sup>a</sup> data point  $x$  as feature that i need to classify to one of these 2 classes

$x$ : weight of animal

evidence : it tells proportion of animal there in 500.

Priest : wt

Likelihood : animal with particular weight is  $x, y, z$  etc

(5)

Using Bayes' decision rule, I will find out:-

$$P(Cl_1/x) \text{ and } P(Cl_2/x)$$

$$P\left(\frac{Cl_1}{x}\right) > P\left(\frac{Cl_2}{x}\right) \Rightarrow Cl_1$$

$$P\left(\frac{Cl_1}{x}\right) = \frac{P(x/Cl_1) * P(Cl_1)}{P(x)} \quad P\left(\frac{Cl_2}{x}\right) = \frac{P(x/Cl_2) * P(Cl_2)}{P(x)}$$

Let's use Bayes theorem:-

$$P(Cl_i/x) = \frac{P(x/Cl_i) P(Cl_i)}{P(x)}$$

Class conditional probability

$$P\left(\frac{Cl_2}{x}\right) = \frac{P(x/Cl_2) P(Cl_2)}{P(x)}$$

Prior probability (my belief about existence  
of particular class)

Total no of animal  
in 200  $\rightarrow$  fixed  
Particular ratio of  
animal in 200

probability depend on  
weight is picked up  
 $P\left(\frac{x}{Cl_i}\right)$   
of particular animal.

$$P(C_1|x) > P(C_2|x) \Rightarrow C_1 \quad \textcircled{⑥}$$

After expanding :—

$$P\left(\frac{x}{C_1}\right) P(C_1) > P\left(\frac{x}{C_2}\right) P(C_2) \Rightarrow C_1$$

$$P\left(\frac{x}{C_i}\right) = \text{Posterior} = \frac{\text{Prior} \times \text{Likelihood}}{\text{Evidence}}$$

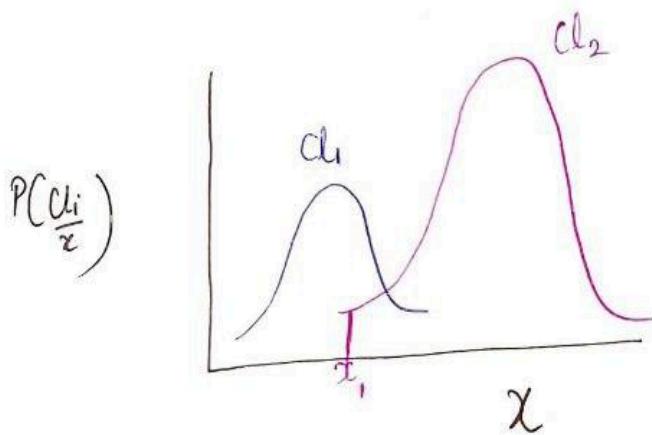
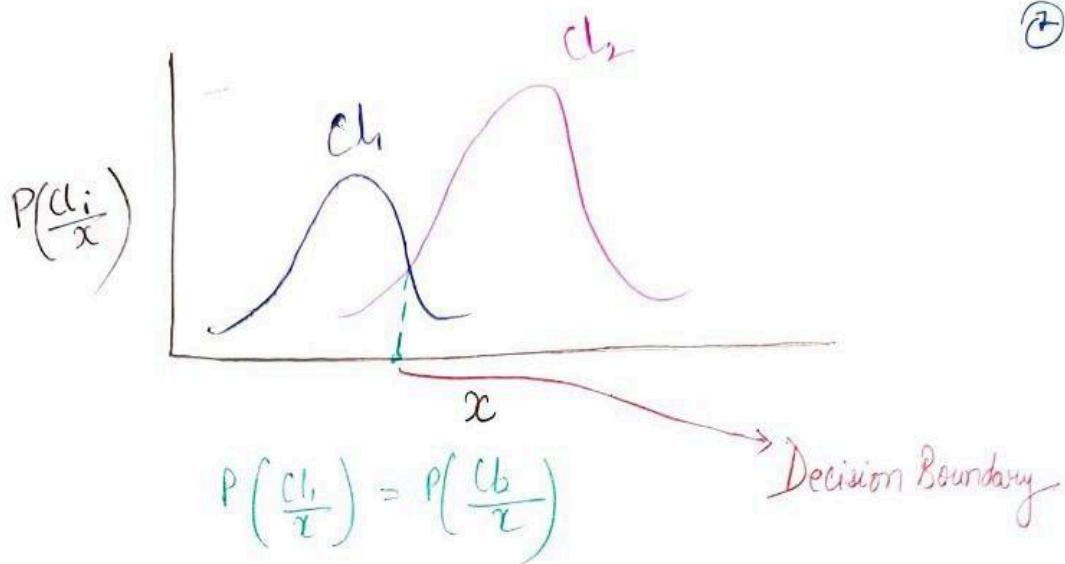
$$P\left(\frac{x}{C_i}\right) = \text{Changes : Class conditional probability.}$$

$P(x)$  = fixed , common in both classes.  
 Before coming to Zoo inside, already presence of no. of animal in it fixed]

So Posterior depends on :—

Posterior  $\propto$  likelihood  $\times$  prior

②



Error in decision =  
 $\min \left\{ P\left(\frac{Cl_1}{x}\right), P\left(\frac{Cl_2}{x}\right) \right\}$

i.e.  $P\left(\frac{Cl_1}{x}\right) = 0.6$        $P\left(\frac{Cl_2}{x}\right) = 0.4$   
↑  
error of this much

Class  $Cl_1$       correct probability  $P\left(\frac{Cl_1}{x}\right)$

⑧

Correct decision is obtained only when we consider minimum error probability.

$$\left[ \begin{array}{l} \text{if I take decision} = C_2 \\ \text{Error probability } P\left(\frac{C_1}{x}\right) \end{array} \right]$$

$$\min \text{ error probab} = P\left(\frac{C_2}{x}\right)$$

o Probability 1  $\longrightarrow$  Consider error of other class i.e  $C_2$

, Probability 2  $\longrightarrow$  error,

check :  $\{ \text{error}_1, \text{error}_2 \}$  min  
 $\downarrow$   
that is class.

①

Loss functionLecture-11

Date : 29/8/2024

- Consider there are  $K$  no. of classes :-

$\rightarrow C_1, C_2, \dots, C_K$   
 States of nature

- Consider that there are ' $q$ ' no. of actions :-

$\alpha_1, \alpha_2, \dots, \alpha_q$   
 Action can be assigning one class to the data  
 Action can also be assigning no class when  
there is tie.

- Loss function +

$\lambda \left( \frac{\alpha_i}{C_j} \right)$  : Loss incurred for taking action  $\alpha_i$  when state of nature is  $C_j$

- We consider a data point  $x$  to be d-dimensional future vector.

- Generic error in terms of loss.

Action : assigning class  $\alpha_i$  to any class

State of Nature :  $C_1, C_2, \dots, C_K$   $\left[ \begin{array}{c} K \\ \text{Classes} \end{array} \right]$

- State of nature = all classes
- If classes are 'equal': might not take action

(2)

~ Given state of nature,

Actual class label  $C_{L_2}$ .

For error = ' $\lambda$ ' function

e.g. State of nature {Cat, Dog}

[Data point € to Dog and we say 'Cat'  
Acc. to this two state of nature loss  $\lambda$   
if told correct classification, loss = 0]

- Earlier with example of zoo, Animal weight, 1 feature = 1D

was considered

now 'D-dimension' is considered.

e.g. weight, color, height, eyes ... etc = D-dimensional.

Expected Loss :-

$$X \rightarrow z$$

↓

$$x_1, x_2, \dots, x_N$$

$$\text{Expected value} \leftarrow \sum_{i=1}^N x_i P(x_i) \quad (3)$$

$x_i$  = don't know which class it belongs to

$a_i$  = Action taken ,  $K = \# \text{ of classes}$ .

$Cl_j$  = True state.

$$\text{Posterior probability} = P\left(\frac{Cl_i}{x}\right)$$

- Expected loss for taking action  $a_i$ , when we observe data  $x$

$$\sum_{j=1}^K \left( \begin{array}{l} \text{loss of being in } Cl_i \\ \text{considering class } j \\ \text{as true} \end{array} \right) * \text{Posterior probability}$$

Expected loss | Risk function | conditional Risk :-

$$R\left(\frac{a_i}{x}\right) = \sum_{j=1}^K \lambda\left(\frac{a_i}{Cl_j}\right) P\left(\frac{Cl_j}{x}\right)$$

↑

expected loss for taking action  $a_i$ , when we observe data  $x$ .

$$\Rightarrow \lambda\left(\frac{a_i}{Cl_1}\right) \cdot P\left(\frac{Cl_1}{x}\right) + \lambda\left(\frac{a_i}{Cl_2}\right) P\left(\frac{Cl_2}{x}\right) \dots \dots$$

$\boxed{\lambda_{ij} = \lambda\left(\frac{a_i}{Cl_j}\right)}$

- We want to take an action which minimizes the risk. \*

### Minimum Risk Classifier.

e.g If state of Nature

$$X = \left\{ \begin{array}{l} \text{pneumonia} \\ \text{viral} \\ \text{cardiac} \\ \text{None} \end{array} \right\}$$

If doctor predicts for pneumonia, but actually it's cardiac.

If doctor starts medicine, risk of health, time.

Goal is = minimum error

$$\lambda \left( \frac{\alpha_i}{C_{lj}} \right) = \lambda_{ij}$$

Classification Problem :- Assigning point to class, given risk (observation).

- Minimum Risk Classifier = do classification that minimizes risk.

$\lambda_{ij}$  = loss of taking action  $\alpha_i$ ,  
given true class  $C_j$

$\boxed{\text{if } i == j} = \text{correct action}$   
 $\hookrightarrow \text{loss should be minimum.}$

(5)

Two class Classification :-

$$R\left(\frac{\alpha_i}{x}\right) = \sum_{j=1}^2 \lambda_{ij} \cdot P\left(\frac{C_j}{x}\right)$$

$\alpha_{11}$  :- True class 1,  $\boxed{\text{True}}$ , error should be small.  
Action  $a_1$

$\alpha_{21}$  :- action  $\rightarrow$  class 2  $\rightarrow$  error  
actual class 1

• There might be no decision too.

eg:- Action may not be related to classification  
e.g. Doctor might call ambulance not classified in any disease

•  $\alpha_j$  = I may have more than  $K$  no. of actions  
 $K = \# \text{ of class}$

It may be  $K \geq K$  for not taking action.

$$R(\alpha_1/x) = \sum_{j=1}^2 \lambda_{1j} \cdot P\left(\frac{C_j}{x}\right) = \underbrace{\lambda_{11} P\left(\frac{C_1}{x}\right)}_{\text{True}} + \lambda_{12} P\left(\frac{C_2}{x}\right)$$

$$R(\alpha_2/x) = \sum_{j=1}^2 \lambda_{2j} \cdot P\left(\frac{C_j}{x}\right) = \underbrace{\lambda_{21} P\left(\frac{C_1}{x}\right)}_{\text{True}} + \underbrace{\lambda_{22} P\left(\frac{C_2}{x}\right)}_{\text{True}}$$

If Risk of assigning  $\alpha$  to class 1 is b.  
we want to assign to class 1.

④

$$\cdot R\left(\frac{\alpha_1}{x}\right) < R\left(\frac{\alpha_2}{x}\right)$$

$$\Rightarrow \left( \alpha_{11} P\left(\frac{C_1}{x}\right) + \alpha_{12} P\left(\frac{C_2}{x}\right) \right) < \left( \alpha_{21} P\left(\frac{C_1}{x}\right) + \alpha_{22} P\left(\frac{C_2}{x}\right) \right)$$

$$\Rightarrow (\alpha_{11} - \alpha_{21}) P\left(\frac{C_1}{x}\right) < (\alpha_{22} - \alpha_{12}) P\left(\frac{C_2}{x}\right)$$

$$\cdot (\alpha_{21} - \alpha_{11}) P\left(\frac{C_1}{x}\right) > (\alpha_{12} - \alpha_{22}) P\left(\frac{C_2}{x}\right)$$

From Baye's decision rule :-

we need to have following to assign class 1 to i/p

data  $x$  &

$$\cdot P\left(\frac{C_1}{x}\right) > P\left(\frac{C_2}{x}\right)$$

If we want to assign  $\alpha$  to class 1, we want:-

$$\cdot (\alpha_{21} - \alpha_{11}) P\left(\frac{C_1}{x}\right) > (\alpha_{12} - \alpha_{22}) P\left(\frac{C_2}{x}\right)$$

↑  
↑

weight

- $(\lambda_{11} \text{ or } \lambda_{22}) \Rightarrow$  is always lower as they are correctly classified) ⑦
- $(\lambda_{12}, \lambda_{21}) \Rightarrow$  is always higher as they are misclassified.

$$\lambda_1 \leq \lambda_2$$

$$(\lambda_{21} - \lambda_{11}) P\left(\frac{C_1}{x}\right) > (\lambda_{12} - \lambda_{22}) P\left(\frac{C_2}{x}\right)$$

↑  
 wrong  
 high value      ↓  
 lower value      (implies)

$\lambda_1$  = true action error so might be zero  
 $\therefore \lambda_1 \leq \lambda_2$ , or zero.

$$\left. \begin{array}{l} \lambda_{21} - \lambda_{22} \\ \lambda_{12} - \lambda_{11} \end{array} \right\} \text{Both are always positive}$$

$\lambda_{21}$  = error  $\approx$  approximately equal to 1  
 $\lambda_{11}$  = true " " " " 0

$$\left( \frac{1}{\lambda_{21}} - \frac{0}{\lambda_{11}} \right) P\left(\frac{C_1}{x}\right) > \left( \frac{1}{\lambda_{12}} - \frac{0}{\lambda_{22}} \right) P\left(\frac{C_2}{x}\right)$$

$\left. \int P\left(\frac{C_1}{x}\right) > P\left(\frac{C_2}{x}\right) \right\} \text{- Bayes' Theorem.}$

Minimum Risk Classifier : For Two Class Problem

$$(d_{21} - d_{11}) P\left(\frac{C_1}{x}\right) > (d_{12} - d_{22}) P\left(\frac{C_2}{x}\right)$$

weight

- From Bayes decision rule, to assign class 1 to Input data x :-

$$P\left(\frac{C_1}{x}\right) > P\left(\frac{C_2}{x}\right)$$

Here, to assign class 1 to  $\frac{C_1}{x}$  :-

$$(d_{21} - d_{11}) P\left(\frac{C_1}{x}\right) > (d_{12} - d_{22}) P\left(\frac{C_2}{x}\right)$$

All the losses ( $d_{ij}$ ) are pre-defined depending on problem.

- $d_{ii} = 0$  : indicates correct decision

minimum error risk classifier is more general version of  
Baye's theorem

It is difficult in term of loss, error, risk.

How to minimize error :-

$$\lambda_{ij} = \lambda(\alpha_i | c_j) = \begin{cases} 0 & \text{when } i=j \\ 1 & \text{when } i \neq j \end{cases} \quad \text{(ii)}$$

we also have :-

$$R\left(\frac{\alpha_i}{x}\right) = \sum_{j=1}^k \lambda\left(\frac{\alpha_i}{c_j}\right) P\left(\frac{c_j}{x}\right) \quad \text{(iii)}$$

combining (ii) and (iii) :-

$$R\left(\frac{\alpha_i}{x}\right) = \sum_{i \neq j} P\left(\frac{c_j}{x}\right)$$

Calculation of risk :-

Let i = 2 and there are 3 classes  $c_1, c_2, c_3$

So,  $\sum_{j=1}^3 P\left(\frac{c_j}{x}\right) = 1$

$$\sum_{j=1}^3 P\left(\frac{c_j}{x}\right) = P\left(\frac{c_1}{x}\right) + \underbrace{P\left(\frac{c_2}{x}\right)}_{\text{if } i=2} + P\left(\frac{c_3}{x}\right)$$

$$\text{So, } \sum_{i \neq j} P\left(\frac{c_j}{x}\right) = P\left(\frac{c_1}{x}\right) + P\left(\frac{c_3}{x}\right) = \boxed{P\left(\frac{c_1}{x}\right) + 1 - P\left(\frac{c_2}{x}\right)}$$

- Every incorrect decision doesn't have same impact.
- Some incorrect decision are more incorrect.  
 ↓  
 Ideally we don't assign same assignment to all incorrect.  
 e.g. Tomorrow is my exam  $\rightarrow$  fail  $\rightarrow$  copy <sup>incorrect decision</sup>  
 $\qquad\qquad\qquad$  watched movie  
 $\qquad\qquad\qquad$  rather than preparing

Limit changes acc to level of incorrect condition.

$$\sum_{i \neq j} P\left(\frac{C_j}{x}\right).$$

here let's assume all as (incorrect  $A=1$ )

$$R\left(\frac{x_i}{x}\right) = \sum_{i \neq j} P\left(\frac{C_j}{x}\right) = 1 - P\left(\frac{C_i}{x}\right)$$

(Posterior probability)

I want to minimize  $x_i$  risk.

$$\text{Maximize } P\left(\frac{C_i}{x}\right)$$

whichever class maximize Posterior probability  
 take decision for that.

Minimum error rate classifier

This is similar to Bayes decision rule.

4

$$P\left(\frac{C_1}{x}\right) > P\left(\frac{C_2}{x}\right) \rightarrow C_1$$

Posterior probability  $\rightarrow$  Bayes' decision rule.

[Every incorrect  $\rightarrow$  different weight,  $\alpha_{ij}$  diff. for diff. incorrect cond.]

Conditional Probability Table & assume, we can tell any probability answer from this.

Drawback  $\rightarrow$  Storage.

for a system with many causes and effects

we have to maintain a large set of values and operate on those

e.g. If 3 random variables

$$2^3 = 8 \text{ entries}$$

↑ storage [• for  $D$  possible answers on  $N$  random variable  $= D^N$   
• 3 answers of 10 question  $= 3^{10}$

## Independent Events

(5)

Two events A and B are said to be independent if

$$P\left(\frac{A}{B}\right) = P(A) \quad \text{--- (1)}$$

we already have

$$P\left(\frac{A}{B}\right) = P\left(\frac{A \cap B}{B}\right) = \frac{P(B|A) P(A)}{P(B)} \quad \text{--- (2)}$$

From (1), (2) ←

$$P\left(\frac{B}{A}\right) = P(B) \quad \boxed{P(A \cap B) = P(A) P(B)}$$

If I want to deal with fever, cold, cough and Internet speed

(Internet speed doesn't depend on others)

So if we want to find out Joint Distribution :—  $P(AB) = P(A) \cdot P(B)$

$$P(\text{fever, cough, covid, I-speed}) = \underbrace{P(\text{fever, cough, covid})}_{A} \cdot \underbrace{P(\text{I-speed})}_{\text{independent B}}$$

If sample space for Internet Speed (I-speed) is

{slow, medium, fast, very fast}

		fever	¬fever
covid	slow medium fast very fast	cough → rough a <sub>1</sub> a <sub>2</sub>	cough → rough a <sub>1</sub> a <sub>2</sub>
¬covid	slow medium fast very fast	⋮      ⋮	⋮      ⋮
			a <sub>32</sub>

Total entries : 32 (31 parameters) (6)

Now I know :-

$$P(\text{fever, cough, covid, } I_{sp}) = P(\text{fever, cough, covid}) P(I_{sp})$$

$$\sum p_1 + \dots + p_{32} = 1$$

$$p = 32 - 1$$

$$\boxed{\text{no. of parameters needed} = 1 - \frac{\text{one variable left}}{31}}$$

$$P(\text{internet}) = 4 \text{ entries}$$

slow	fast	medium	very fast
------	------	--------	-----------

Table  
for  $I_{sp}$

↑  
4 entries, 3 parameters ( $4 \times 3$ )

Before entries required was :- 32 (31 parameters)

After performing factorization :-

$$P(\text{f, cough, covid, } I_{sp}) = P(f, \text{cough, covid}) P(I_{sp})$$

$$\text{Total entries} = \frac{8}{2} + 4 = 12 \quad (7+3=10 \text{ parameters})$$

Complete Independence is extremely powerful.

+ we will not include such factors (useless for our modelling purpose)

• only by using notion of Independent events = Reduced Storage.

complete independent is not feasible.

Some factor completely or may not completely, partially independent to others.

partially + weather cond<sup>n</sup>, food  
partially dependent.

### Conditional Independence ↗

I want to calculate probability of fever given a person has covid, cough.

$$P\left(\frac{\text{fever}}{\text{covid, cough}}\right)$$

Knowledge about cough doesn't give additional info to predict fever.

↙ fever is outcome of covid  
fever is not outcome of cough

(Given covid, fever is conditional independent of cough.)

$$\text{eg 2: } P\left(\frac{\text{flood}}{\text{water leakage}}, \frac{\text{heavy rainfall}}{\text{rainfall}}\right)$$

(8)

flood is conditional independent of water leakage

So we can say

$$P\left(\frac{\text{fever}}{\text{covid, cough}}\right) = P\left(\frac{\text{fever}}{\text{covid}}\right)$$

also,  $P\left(\frac{\text{fever}}{\neg \text{covid}, \text{cough}}\right) = P\left(\frac{\text{fever}}{\neg \text{covid}}\right)$

∴ fever is conditional independent of cough, given covid.

To model,  $P(\text{fever, covid, cough}) = 8 \text{ entries, 7 parameters}$

$$P(f, \text{cough}, \text{covid}) = P\left(\frac{\text{fever}}{\text{covid, cough}}\right) P(\text{covid, cough})$$

$$= P\left(\frac{f}{\text{covid}}\right) P(\text{covid, cough})$$

$$= P\left(\frac{\text{fever}}{\text{covid}}\right) P\left(\frac{\text{covid}}{\text{cough}}\right) P(\text{cough})$$

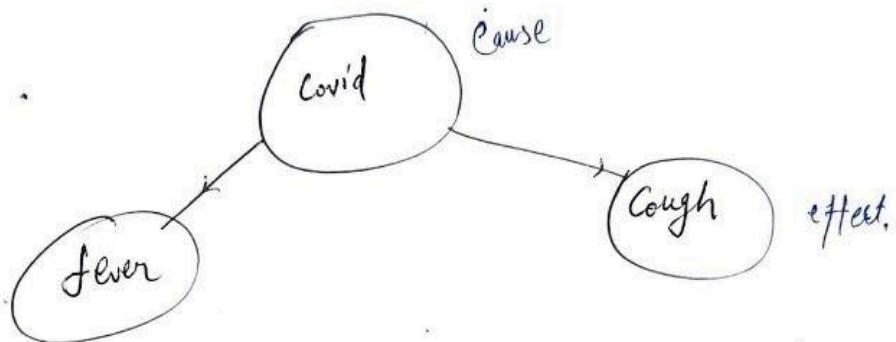
$\uparrow$   
2 parameter

$\uparrow$   
2  
 $\uparrow$   
1

= 5 parameter.

①

## Graphical Representation :-



one cause resulting in multiple independent effects.

Also known as Cause-Effect Graph.

Independent events don't interact with each other.

eg:- Genetic disease disorder  
for blood, kidney if Hemoglobin is low, kidney effect  
not independent.

Genetic disease disorder  
don't follow conditional independence.

eg2r Heavy rainfall : cause [power cut, water cut]

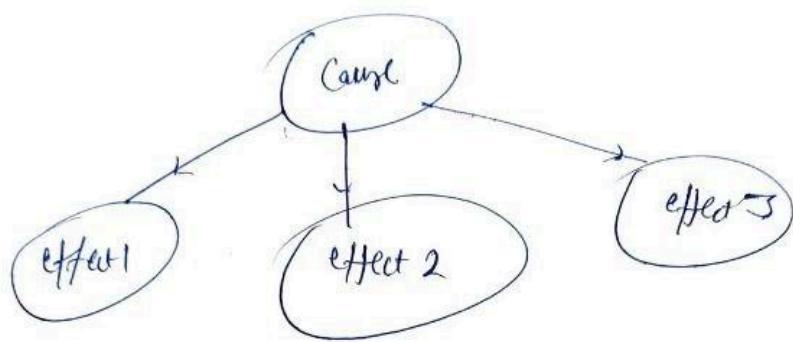
for some causes

[ ]

effects are independent

effects are dependent

10



Naïve assumption [one cause resulting in multiple independent effects.]

Naïve assumption : effects are independent, in Cause-effect

## Lecture - 13

5-Sept-2024

①

Problem 1: At a certain university? —

4.1. men  $\rightarrow$  over 6 feet tall.

1.1. women  $\rightarrow$  over 6 feet tall.

Total students are in ratio  $\neq 3:2$  in favor of women.

If a student is selected at random from all those  
6 feet tall.

what is probability that student is woman?

Solution:

Probability that student is male  $P(M) = \frac{2}{5}$

" " " female  $P(W) = \frac{3}{5}$

Probability that a male that is over 6 feet tall:  $P\left(\frac{H}{M}\right) = 0.04$

" " " female " " " " "  $P\left(\frac{H}{W}\right) = 0.01$

Total probability of being over 6 feet tall: —

$$P(H) = P\left(\frac{H}{M}\right) \cdot P(M) + P\left(\frac{H}{W}\right) \cdot P(W)$$

$$P(H) = (0.04 \times 0.4) + (0.01 \times 0.6) = 0.016 + 0.006 \\ = 0.022$$

Using Bayes' rule: —

$$P\left(\frac{W}{H}\right) = \frac{P(H|W) \cdot P(W)}{P(H)} = \frac{0.01 \times 0.6}{0.022} \approx 0.273$$

Problem-2

(2)

Probability of certain disease = 0.01

Probability of test positive given person is infected with  
disease = 0.95

Probability of test positive given person is not infected  
with disease = 0.05

(a) Calculate probability of test positive

(b) Using Baye's rule, calculate probability of being infected with  
disease given that test is +ve.

Soln =

$$P(D) = 0.01 \quad P\left(\frac{T}{D}\right) = 0.95$$

(a)

$$P\left(\frac{T}{\neg D}\right) = 0.05$$

$$P(T) = P\left(\frac{T}{D}\right) \cdot P(D) + P\left(\frac{T}{\neg D}\right) P(\neg D)$$

$$(0.95 \times 0.01) + (0.05 \times 0.99) =$$

$$0.0095 + 0.0495 = 0.059$$

$$(b) P\left(\frac{D}{T}\right) = \frac{P\left(\frac{T}{D}\right) \cdot P(D)}{P(T)} = \frac{0.95 \times 0.01}{0.059} = 0.161$$

Problem 3: A factory manufactures bolts using A, B, C machines. (3)  
Defect rates are :—

- . Machine A : 25%. proportion, Defect rate : 5%.
- . Machine B : 35%. , Defect : 4%.
- . Machine C : 40%. , Defect rate : 2%.

A bolt is chosen at random, found defective.

Probability it came from +

- (a) machine A
- (b) machine B
- (c) machine C

Soln: Total probability of Defect :—

$$P(D) = P\left(\frac{D}{A}\right) \cdot P(A) + P\left(\frac{D}{B}\right) \cdot P(B) + P\left(\frac{D}{C}\right) \cdot P(C)$$

$$P(D) = (0.05 \times 0.25) + (0.04 \times 0.35) + (0.02 \times 0.40)$$

$$= 0.0345$$

using Bayes' rule

$$P\left(\frac{A}{D}\right) = \frac{P\left(\frac{D}{A}\right) \cdot P(A)}{P(D)} = \frac{0.05 \times 0.25}{0.0345} = 0.3623$$

$$P\left(\frac{C}{D}\right) = \frac{P\left(\frac{D}{C}\right) \cdot P(C)}{P(D)} = \frac{0.02 \times 0.40}{0.0345} = 0.2319$$

$$P\left(\frac{B}{D}\right) = \frac{P\left(\frac{D}{B}\right) \cdot P(B)}{P(D)} = \frac{0.04 \times 0.35}{0.0345}$$

Q4

Q. A garage machine has box with

good springs : 5

faulty springs : 3

Machine picks from box = 2 springs.

Find probability of (a) 1st spring drawn is faulty  
 (b) 2nd " "

Soln Probability of 1st spring drawn is faulty :-

$$P(\text{First faulty}) = \frac{\text{No of faulty spring}}{\text{Total no. of spring}} = \frac{3}{8}$$

Total probability that 2nd spring is faulty :-

Case 1:- First spring is faulty :-

$$P\left(\frac{\text{2nd faulty}}{\text{first faulty}}\right) = \frac{2}{7}$$

Case 2:- first spring is not faulty :-

$$P\left(\frac{\text{second faulty}}{\text{first not faulty}}\right) = \frac{3}{7}$$

Overall probability :-

$$P(\text{second faulty}) = \frac{3}{8} \times \frac{2}{7} + \frac{5}{8} \times \frac{3}{7} = \frac{21}{56} = \frac{3}{8}$$

Or manufacturing plant produces three types of widgets → (S)

Type	<u>Proportions</u>	<u>Defects rate</u>
A	40%	5%
B	35%	4%
C	25%	3%

(a) what is probability that defective widget is of Type A?

$$P(D) = P\left(\frac{D}{A}\right) \cdot P(A) + P\left(\frac{D}{B}\right) \cdot P(B) + P\left(\frac{D}{C}\right) \cdot P(C)$$

$$= (0.05 + 0.40) + (0.04 \times 0.35) + (0.03 \times 0.25)$$

$$= 0.0415$$

using Bayes' rule

$$P\left(\frac{A}{D}\right) = \frac{P\left(\frac{D}{A}\right) \cdot P(A)}{P(D)} = \frac{0.05 \times 0.40}{0.0415} = 0.4819$$

(b) Recalculate if defect rate of Type A increases to 10%.

New defect rate for Type A = 10%.

Recalculate Total probability of Defect →

$$P(D) = (0.10 \times 0.40) + (0.04 \times 0.35) + (0.03 \times 0.25)$$

$$= 0.0615$$

Using Bayes' Rule

$$P\left(\frac{A}{D}\right) = \frac{0.10 \times 0.40}{0.0615} = 0.6504$$

(c) Similarly recalculate if defect rate for type B is reduced by 1%. (6)

<u>Quiz 11</u>	$x_0$	$x_1$	$x_2$	$t$	$w_0 = -0.3$	$\eta = 0.5$
	0	0	0		$w_1 = w_2 = 0.5$	
	0	1	1		$-x_0 = 1$	Default
	1	0	1			
	1	1	0		Calculate gradient :-	

Soln :- Model Setup :-  
 $y_{pred} = \text{step}(w_0 x_0 + w_1 x_1 + w_2 x_2)$

i) for  $x = (0, 0)$

$$\begin{aligned} Z &= w_0 x_0 + w_1 x_1 + w_2 x_2 \\ &= (-0.3 * 1) + (0.5 * 0) + (0.5 * 0) \\ &= -0.3 \end{aligned}$$

$$\Delta w_0 = \Delta w_0 + \eta \left( t - y_{predicted} \right) x_0$$

true value      predicted

$$= 0 + 0.5 \left( 0 - (-0.3) \right) = 0.15$$

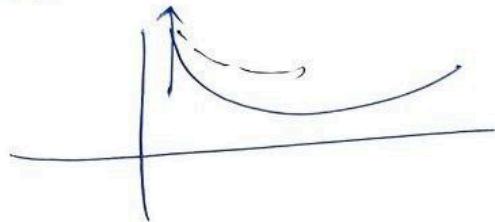
$$y_{pred} = \text{step}(-0.3) = 0$$

$$\begin{aligned} \Delta w_0 &= -0.3 + 0.15 \\ &> 0.15 \end{aligned}$$

$$\begin{aligned}\sum w_i x_i &= w_0 + w_1 x_1 + w_2 x_2 \\ &= -0.3 + 0.5(0) + 0.5(1) \\ &= 0.2\end{aligned}$$

(1)

only for 1st update it (error) spikes then decreases



$$x_0 = 0, x_1 = 1$$

$$\begin{aligned}\sum w_i x_i &= w_0 x_0 + w_1 x_1 + w_2 x_2 \\ &= -0.3 \times 1 + 0.5 \times 0 + 0.5 \times 1 \\ &= 0.2\end{aligned}$$

$$\begin{aligned}\Delta w_0 &= \Delta w_0 + \eta(t-0)x_0 \\ &= 0.15 + 0.5(1 - 0.2) \\ &\approx 0.185\end{aligned}$$

$$x_1 = 1, x_2 = 0$$

$$\begin{aligned}\Delta w_0 &= \Delta w_0 + \eta(t-0)x_0 \\ &= 0.15 + 0.5(0 - 0.7) \times 1 = 0.6\end{aligned}$$

$$\begin{aligned}\sum w_i x_i &= w_0 x_0 + w_1 x_1 + w_2 x_2 \\ &= 0.3 + 0.5(1) + 0.5(0) \\ &= 0.7\end{aligned}$$

$$w = w_0 + \Delta w_0 = -0.3 + 0.6 = 0.3$$

⑧

$$\text{loss} = \frac{1}{2} (\text{predicted} - \text{actual})^2$$

$$\frac{\partial \text{loss}}{\partial w} = w_1 - \eta (t - o)^2$$

$$\text{cost} = J = \sum (y - y_{\text{pred}})^2$$

ML-14

Date : 9/09/2024 ①

### Naïve Bayes Model

- Bayes's decision rule :  $P(Cl_1/x) > P(Cl_2/x) \Rightarrow Cl_1$

Now let's use Bayes theorem :-

$$P\left(\frac{Cl_1}{x}\right) = \frac{P(x/Cl_1)P(Cl_1)}{P(x)}$$

Let's consider  $x$  to be d-dimensional feature vector  $\{x_1, x_2, \dots, x_d\}$

- All features are conditionally independent of each other given the class

Now Let's use Bayes theorem :-

$$\begin{aligned} P\left(\frac{Cl_1}{x}\right) &= \frac{P(x, Cl_1)}{P(x)} \\ &= \frac{P(x_1, x_2, \dots, x_d, Cl_1)}{P(x)} \\ &= \underbrace{P\left(\frac{x_1}{Cl_1}\right) P\left(\frac{x_2}{Cl_1}\right) \dots P\left(\frac{x_d}{Cl_1}\right)}_{P(x)} P(Cl_1) \end{aligned}$$

$$P(Effect_1, effect_2, \dots, Effect_d, Cause)$$

$$= P(Cause) \prod_{k=1}^d P\left(\frac{Effect_k}{Cause}\right)$$

(2)

Similarly :-

$$P\left(\frac{C_{l_2}}{x}\right) = \frac{P\left(\frac{x_1}{C_{l_2}}\right) \cdot P\left(\frac{x_2}{C_{l_2}}\right) \cdots P\left(\frac{x_d}{C_{l_2}}\right) P(U_2)}{P(x)}$$

• Baye's decision rule :-

$$P\left(\frac{U_1}{x}\right) > P\left(\frac{U_2/x}{x}\right) \rightarrow C_1$$

 $P(U_2)$ 

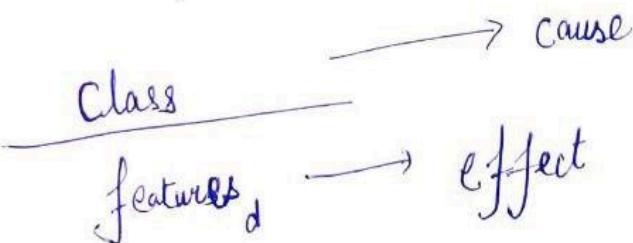
$$\underbrace{P(x_1/U_1) P(x_2/U_1) \cdots P(x_d/U_1) P(U_1)}_{P(x)} > \underbrace{\frac{P(x_1/C_{l_2}) P(x_2/C_{l_2}) \cdots P(x_d/C_{l_2}) P(C_{l_2})}{P(x)}}_{P(U_2)}$$

c)  $C_{l_1}$ 

also  $P(x_1/U_1) P\left(\frac{x_2}{C_{l_1}}\right) \cdots P\left(\frac{x_d}{C_{l_1}}\right) > \underbrace{P\left(\frac{x_1}{C_{l_1}}\right) \cdot P\left(\frac{x_2}{C_{l_1}}\right) \cdots P\left(\frac{x_d}{C_{l_1}}\right) \cdot P(C_{l_1})}_{P(U_1)}$

Since from denominator  $\underline{P(x)}$  is not there, we'll  
not get exact answer  $\rightarrow$  Probable answer.

e.g. for a classification problem of Tiger, Lion classification.



As features are conditionally independent given class.

(5)

e.g. -  $\frac{\text{Body color}}{\text{Lion}}, \frac{\text{size of Lion}}{\text{Lion}}, \text{Lion}$   
Conditionally Independent

e.g. - If a person may or may not have Anemia if either

BP weight Appetite  
high low higher lower high low  
2 values for each feature  
given Anemia, features BP, weight, appetite are conditionally independent.

→ It follows chain rule.

→  $P_i \rightarrow \prod \rightarrow \text{It is for product}$

\* Here feature values are discrete.

If feature value is continuous  $\rightarrow$  It becomes complex

→ can't construct Table values of continuous values ( $\infty$  value)  
infinite

• what if conditional probability becomes 0

(4)

completely out of context as whole output is zero.

→ Let's assign some small value to this case

↓  
else output might end up = 0

Example of Naïve Bayes :-

Sample no	Appetite	Weight	BP	Class
1	low	Normal	low	No Anemia
2	low	low	low	Anemia
3	normal	low	low	Anemia
4	low	low	Normal	No Anemia
5	normal	low	Normal	Anemia
6	normal	Normal	low	Anemia
7	normal	Normal	Normal	No Anemia

If we see a new sample who has normal appetite, normal weight, and low BP. Find if person has Anemia?

Soln:-

Step 1 :- Calculate Prior :-

$$P(\text{Probability of having Anemia}) = \frac{4}{7} \quad P(\overline{\text{Anemia}}) = \frac{3}{7}$$

(5)

Ques

Step 2: Calculate  $P\left(\frac{\text{feature}}{\text{Anemia}}\right)$ ,  $P\left(\frac{\text{feature}}{\text{No Anemia}}\right)$  from

training data:

Appetite	Anemia		$\overline{\text{Anemia}}$
	Low	High	Low
Normal	$\frac{1}{4}$	$\frac{3}{4}$	$\frac{2}{3}$

$$P(\text{Appetite} = \frac{\text{Low}}{\text{Anemia}})$$

BP	Anemia		$\overline{\text{Anemia}}$
	Low	High	Low
Normal	$\frac{1}{4}$	$\frac{3}{4}$	$\frac{2}{3}$

Weight	Anemia		$\overline{\text{Anemia}}$
	Low	High	Low
Normal	$\frac{1}{4}$	$\frac{3}{4}$	$\frac{2}{3}$

$$P(\text{Anemia}) = \frac{4}{7} \quad P(\text{no anemia}) = \frac{3}{7}$$

Step 3: Testing:-

Suppose, I observe a test data with normal appetite, normal weight, and low BP.

- To predict if person has Anemia:-

$$P\left(\text{appetite} = \frac{\text{normal}}{\text{anemia}}\right) P\left(\text{weight} = \frac{\text{normal}}{\text{anemia}}\right) P\left(\text{BP} = \frac{\text{low}}{\text{Anemia}}\right) \cdot P(\text{anemia}) \quad (6)$$

$$= \frac{3}{4} \times \frac{1}{4} \times \frac{3}{4} \times \frac{4}{7} = 0.08$$

$$\frac{P(\text{no anemia})}{P(\text{anemia})} = P\left(\text{appetite} = \frac{\text{normal}}{\text{Anemia}}\right) \cdot P\left(\text{weight} = \frac{\text{normal}}{\text{Anemia}}\right) \cdot P\left(\text{BP} = \frac{\text{low}}{\text{Anemia}}\right) \cdot P(\text{no anemia})$$

$$= \frac{2}{3} \times \frac{2}{3} \times \frac{1}{3} \times \frac{3}{7} = 0.06$$

$$P(C_1) \cdot P\left(\frac{x_1}{C_1}\right) \cdot P\left(\frac{x_2}{C_1}\right) \cdots P\left(\frac{x_d}{C_1}\right) > P\left(\frac{x_1}{C_2}\right) \cdot P\left(\frac{x_2}{C_2}\right) \cdots P\left(\frac{x_d}{C_2}\right) \cdot P(C_2)$$

$$\Rightarrow C_1$$

Conclusion: Person has anemia as per Naive Bayes Classifier

$$\text{i.e. } P(\text{anemia}) > P(\text{no anemia})$$

(7)

### Training data

- Assume  $N$  training samples and class label pairs  $(x^1, y^1), (x^2, y^2), \dots, (x^N, y^N)$ .
- Each training sample  $x^i$  is a  $D$  dimensional feature vector  $\{x_1^{(i)}, x_2^{(i)}, \dots, x_D^{(i)}\}$   
e.g.  $3$  is value for  $\frac{w_1}{x_1}, \frac{appetite}{x_2}, \frac{BP}{x_3}$
- Assume that attribute  $x_K$  can take values:-  
 $x_{K_1}, x_{K_2}, \dots, x_{K_N}$   
e.g.  $BP$  [ High  $K=2$  here for each there is  $K$ .  
Low ]
- The values of attributes are discrete.
- We have total of  $C$  number of classes  $c_1, c_2, \dots, c_C$ .

### Training Method

for  $j = 1 : C$   
calculate  $P(c_j)$  from training data

for  $d = 1 : D$

for  $q = 1 : V$

Calculate  $P(x_{dq} | c_j)$

## Naïve Bayes' classifier: The algorithm

(8)

Inference for a new test sample  $x^{(new)}$ , feature vector is  $\{x_1^{new}, x_2^{new}, \dots, x_d^{new}\}$

for  $j = 1 : c$

→ get  $P(C_j)$  compute during training

→ get  $P\left(\frac{x_1^{new}}{C_j}\right) \cdot P\left(\frac{x_2^{new}}{C_j}\right) \cdots P\left(\frac{x_d^{new}}{C_j}\right)$  computed during train

. Calculate

$$\pi_j = P\left(\frac{x_1^{new}}{C_j}\right) P\left(\frac{x_2^{new}}{C_j}\right) \cdots P\left(\frac{x_d^{new}}{C_j}\right) \cdot P(C_j)$$

. find out  $j$  for which  $\pi_j$  is maximum.

$$j^* = \underset{j}{\operatorname{argmax}} \pi_j$$

, Assign class label  $j^*$  to test data.

- Dependence of features in Algo is complex problem.
- "Zero problem" completely discard a particular class.

### Life without Naive Bayes'

, consider  $x$  to be a  $d$ -dimensional feature vector

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

If all features are not conditionally independent of each other

given the class



No Naive Baye = make solution very very complex.

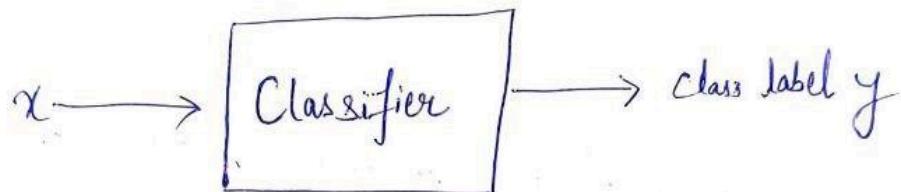
Inductive Bias: making assumption of Condition Independent makes life easier.

Bayes theorem would be :-

$$\begin{aligned}
 P\left(\frac{C_1}{x}\right) &= \frac{P(x, C_1)}{P(x)} \\
 &= \frac{P(x_1, x_2, \dots, x_d, C_1)}{P(x)} \\
 &= \frac{P\left(\underbrace{x_1, x_2, \dots, x_d}_{x}, C_1\right) \cdot P\left(\frac{x_2}{x_1, \dots, C_1}\right) \cdots P\left(\frac{x_d}{x_1, \dots, C_1}\right) P(C_1)}{P(x)}
 \end{aligned}$$

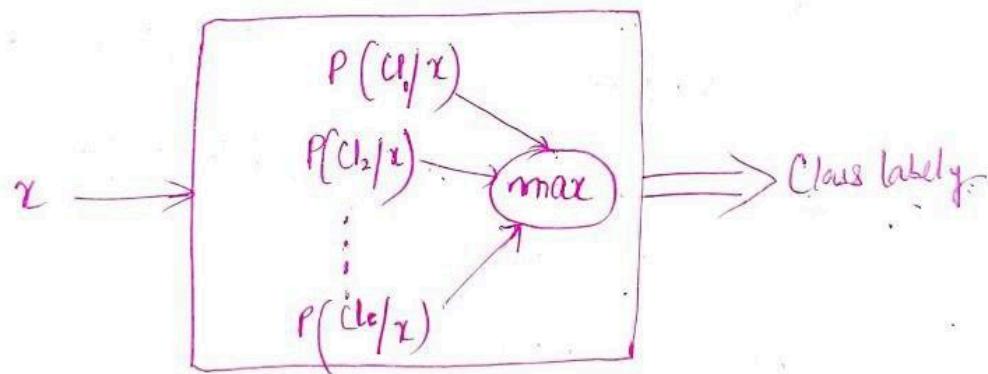
Computing these conditional probabilities would be  
 extreme difficult.

Decision rule: try to find probability which gives maximum value discriminant fun<sup>n</sup> for  $x$ . ②

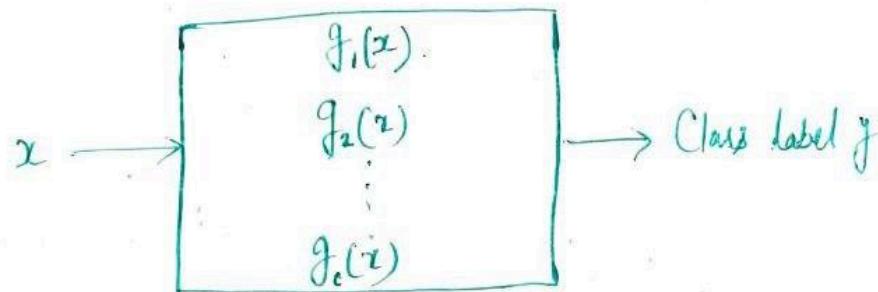


How does the classifier do it?

Suppose I have C number of classes. I consider Baye's decision rule



$g_i(\cdot)$  is called discriminant function.



If  $g_i(x) > g_j(x)$ ,  $\forall j \neq i$  we assign class label  $i$  to the input data  $x$ . (3)

### Nature of Discriminant functions-

For the minimum risk classifier:

We assign the class label corresponding to the action of minimum risk.

If the risk of action  $\alpha_i$  is  $R(\alpha_i/x)$ .

- Minimum risk  $\rightarrow$  in order to minimize risk  
= negative of risk.

$$g_i(x) = -R\left(\frac{\alpha_i}{x}\right)$$

- we assign class label corresponding to maximum posterior probability.

$$\boxed{g_i(x) \propto P\left(\frac{C_i}{x}\right)}$$

discriminant fun<sup>n</sup>  $\propto \begin{pmatrix} \text{posterior} \\ \text{probability} \end{pmatrix} \uparrow \uparrow$

- choice of discriminant function is not unique.

$$\text{eg } -R, h = \frac{1}{e^x}, R = \frac{1}{\text{softmax}} \dots \dots$$

$-R\left(\frac{\alpha_i}{x}\right)^2$  = valid fun<sup>n</sup>, also  $-2\left[R\left(\frac{\alpha_i}{x}\right)\right]$  → also valid fun<sup>n</sup>.

If  $g_i(x)$  is a discriminant fun<sup>n</sup> and  $f(\cdot)$  is monotonically increasing function ,  $f(g_i(x))$  is also a valid discriminant fun<sup>n</sup>. (4)

If  $g_i(x) > g_j(x)$ ,  $\forall j \neq i$ , we assign class label  $i$  to input data  $x$ .

$$g_i(x) \propto P\left(\frac{U_i}{x}\right)$$

$$\propto \frac{P\left(\frac{x}{U_i}\right) \cdot P(U_i)}{P(x)}$$

Since denominator is common for every class  $i$  :-

we take

$$g_i(x) = P\left(\frac{x}{U_i}\right) \cdot P(U_i)$$

fun<sup>n</sup> monotonically increasing means :-

with increasing value of argument, monotonically increasing function gives higher values.

if  $x_1 \rightarrow x_2 \rightarrow x_3$

$$f(x_1) < f(x_2) < f(x_3)$$

monotonically increasing function don't have equality. (<sup>negative</sup>  
<sub>means have</sub>)

(5)

i.e.  $\hat{f}(g_i(x)) \neq \hat{f}(g_j(x))$   
 ↓  
 can't be equal

we take

$$g_i(x) = P\left(\frac{x}{C_{li}}\right) P(C_{li}) \rightarrow$$

multiplication is difficult

↓

addition is better.  
(log likelihood)

we can also take :—

$$\begin{aligned} g_i(x) &= \ln P\left(\frac{x}{C_{li}}\right) \cdot P(C_{li}) \\ &= \ln P\left(\frac{x}{C_{li}}\right) + \ln P(C_{li}) \end{aligned}$$

As log is monotonically increasing function.

- If we take a two-class problem, 2 discriminant func's  $g_1(x), g_2(x)$ .
  - If  $g_1(x) > g_2(x)$ , we conclude that  $x$  belongs to class 1.
  - If  $g_1(x) < g_2(x)$ , we conclude that  $x$  belongs to class 2.

Decision Boundary

$$\boxed{g_1(x) = g_2(x)}$$

$$g_1(x) > 0 : \text{class 1} \quad \left[ g^{(1)} = \underbrace{g_1(x)}_{\substack{\text{+ve} \\ \text{-ve}}} - \underbrace{g_2(x)}_{\substack{\text{-ve} \\ \text{+ve}}} \right]$$

class 2

$$\left[ g^{(1)} = \underbrace{g_1(x)}_{\substack{\text{-ve} \\ \text{+ve}}} - \underbrace{g_2(x)}_{\substack{\text{+ve} \\ \text{-ve}}} \right]$$

-ve value

As numerator is good enough to do classification so  
 removed denominator  $P(x)$  :-

$$g_i(x) = P\left(\frac{x}{U_i}\right) = P\left(\frac{x}{U_i}\right) P(U_i)$$

Taking ln of prior, posterior to make it easy as addition is easier than multiplication :-

$$\begin{aligned} g(x) &\Rightarrow g_1(x) - g_2(x) \\ &= \ln(g_1(x)) - \ln(g_2(x)) \\ &= \ln\left[P\left(\frac{x}{U_1}\right) \cdot P(U_1)\right] - \ln\left[P\left(\frac{x}{U_2}\right) \cdot P(U_2)\right] \\ &= \ln\left[P\left(\frac{x}{U_1}\right) \cdot P(U_1)\right] - \ln\left[P\left(\frac{x}{U_2}\right) \cdot P(U_2)\right] \\ &= \ln P\left(\frac{x}{U_1}\right) + \ln P(U_1) - \ln P\left(\frac{x}{U_2}\right) - \ln P(U_2) \\ &= \underbrace{\ln\left[\frac{P(x|U_1)}{P(x|U_2)}\right]}_{\text{Posterior}} + \underbrace{\ln\left[\frac{P(U_1)}{P(U_2)}\right]}_{\text{Prior}} \end{aligned}$$

here at outside we  
 can't write as  $\frac{\ln P(U_1)}{\ln P(U_2)} x, \ln\left(\frac{P(U_1)}{P(U_2)}\right)$

$$g_i(x) = \ln P\left(\frac{x}{c_{l_i}}\right) \quad (7)$$

Talking about probability distribution :-

$P\left(\frac{x}{c_{l_i}}\right)$  is probability distribution which can take many forms.

- Assume that  $P\left(\frac{x}{c_{l_i}}\right)$  follows Gaussian or Normal distribution in d-dimensions.

So,

$$P\left(\frac{x}{c_{l_i}}\right) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp\left[-\frac{(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)}{2}\right]$$

$\mu_i$  = Expected value of  $x$  (samples) given that  $x$  belongs to class  $c_{l_i}$ .

$\Sigma_i$  = Covariance matrix computed from  $x$  (samples) given that  $x$  belongs to class  $c_{l_i}$

$d$  dimensional  $\Rightarrow d$  features

- Talking about probability distribution : Continuous term.  
↓  
value can be find out in range.

Expected value :-  $\frac{1}{n} \sum x_i$

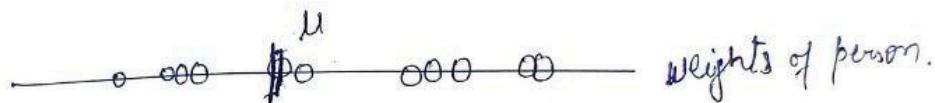
Variance :- Scattering measure from mean.

$$p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2} \right] \quad (8)$$

$\mu_i, \mu_j$  are two diff. mean, (may or may not)

Variance :-

$$\text{var}(x) = E[(x-\mu)^2]$$



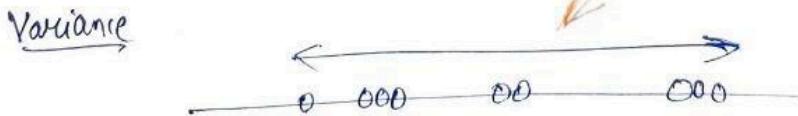
Suppose, I want to measure mean weights

$$\text{Mean } \mu = E(x) = \int x \cdot p(x) dx$$

for  $N$  distinct (equally likely) data points  $x_1, \dots, x_N$

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i$$

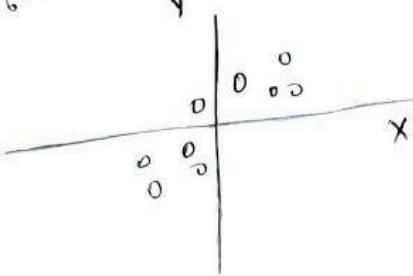
Variance indicates spread



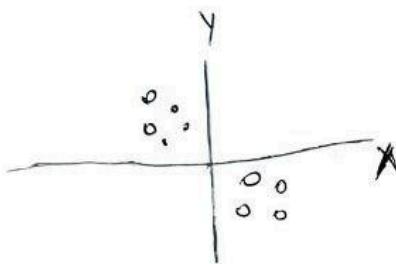
for  $N$  distinct (equally likely) data points  $x_1, x_2, \dots, x_N$

$$\text{var}(x) = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

### Covariance :-



we observe that as  $X \uparrow$   
increases,  $Y$  also increases.



we observe that as  $X$  increases,  
 $Y$  decreases

If we look at  $X, Y$  spread can't differentiate b/w  $X, Y$  data  
individually.

→ we need to see  $X, Y$  together

↓  
Joint Variance : Covariance.

$$\boxed{\text{cov}(X, Y) = E \left[ (X - \mu)(Y - \mu)^T \right]}$$

$\mu_X, \mu_Y$  are different

$$\text{cov}(X, Y) = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$$

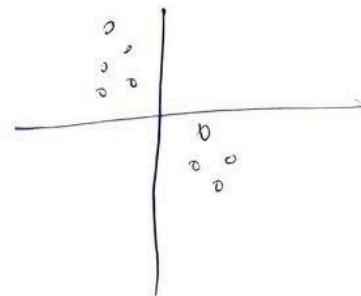
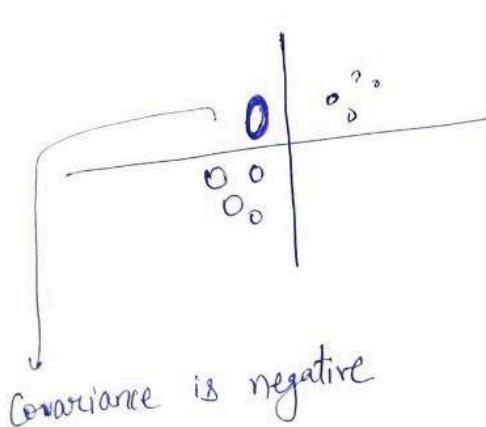
pointwise product of differences from mean.

shows a joint trend of different variables.

(10)

[How much difference)  
in X direction], [How much it differs)  
in Y direction]

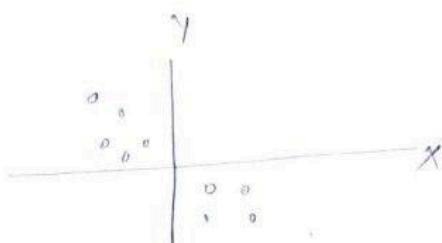
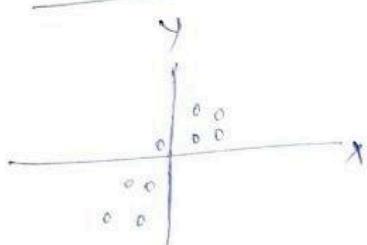
Tell variance of data along X,Y direction = correlated.



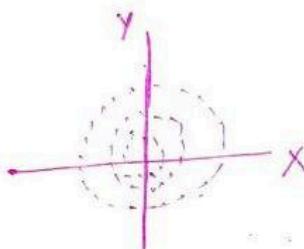
Covariance is +ve.

$$\Sigma = \begin{pmatrix} \text{var}(X) & \text{cov}(X, Y) \\ \text{cov}(X, Y) & \text{var}(Y) \end{pmatrix}$$

Covariance Matrix :



$$\Sigma = \begin{pmatrix} \text{var}(x) & \text{cov}(x,y) \\ \text{cov}(x,y) & \text{var}(y) \end{pmatrix}$$



$\Rightarrow$  Covariance = 0

"No trend/trade of data"

-ve product	+ve product
+ve product	-ve product

• Trade +ve = covariance

• Trade -ve = -ve covariance.

$$\Sigma = \begin{pmatrix} \text{variance off diagonal} \end{pmatrix} \rightarrow \text{variance off diagonal}$$

$\rightarrow$  off diagonal = co-variance

- Covariance is calculated for pair of random variable.

$\Downarrow$   
only for pair

- Covariance Matrix of any no of variable.

(2)

$$\text{Cov}(x, x) = \text{var}(x)$$

for same variable its covariance is another variable is  $x$  itself.

- Covariance for three random variables  $x_1, x_2, x_3$  :-

$$\Sigma = \begin{pmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) & \text{cov}(x_1, x_3) \\ \text{cov}(x_2, x_1) & \text{var}(x_2) & \text{cov}(x_2, x_3) \\ \text{cov}(x_3, x_1) & \text{cov}(x_3, x_2) & \text{var}(x_3) \end{pmatrix}$$

## Probability Density and the Discriminant function

$$\begin{aligned} g_i(x) &= \ln \left( P\left(\frac{x}{C_{li}}\right) \cdot P(C_{li}) \right) \\ &= \ln P\left(\frac{x}{C_{li}}\right) + \ln P(C_{li}) \end{aligned}$$

$$P\left(\frac{x}{C_{li}}\right) = \text{likelihood}, \quad \ln P\left(\frac{x}{C_{li}}\right) : \text{log likelihood}$$

$$P(C_{li}) = \text{prior}, \quad \ln P(C_{li}) = \text{log prior}$$

$P\left(\frac{x}{C_{li}}\right)$   $x_i$  may be any of feature :- weight, appetite, BP...

$x$  = observed data

$P\left(\frac{x}{C_{li}}\right)$  = likelihood of observed data given a particular class.

(3)

given data for a particular class  $i$ .

e.g. all data given for Anemia circumstances :-

Calculate  $P\left(\frac{x_i}{c_{li}}\right)$  e.g.  $P\left(\frac{\text{weight} = \text{low}}{\text{anemia}}\right) \xrightarrow{?} c_i$

- Likelihood can be at any distribution.

→ we only have subset of data i.e. Small sample size  
so whole distribution can be at any Distribution.

We need to have inductive bias i.e. need to assume any particular distribution of data.

e.g. For this small subset of data in d dimension i.e. d no. of features  
given distribution is Normal distribution.

$$P\left(\frac{x}{c_{li}}\right) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp\left[-\frac{(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)}{2}\right]$$

ignoring  $P(x)$  as it was common for both classes.

$|\Sigma_i|$  = determinant of covariance matrix

↳ corresponding to class  $i$

e.g. only for anemia class patient  $|\Sigma|$

$\mu_i$  = mean of class  $i$  = mean is feature vector.

$\mu_i = \frac{1}{N} \sum w_i t_i$ ,  $\frac{1}{N} \sum \text{appetite}$ ,  $\frac{1}{N} \sum \text{BP} \dots$  d times  
 vector quantity d dimensional.

(9)

$\Sigma_i$  = d x d matrix ,  $x$  = d dimension

→ In order to solve for  $i$ , need to take Log of that term.

$$g_i(x) = -\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(\Sigma_i) + \ln P(C_i)$$

$$g_i(x) = -\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(\Sigma_i) + \ln P(C_i)$$

Parameters :  $\mu_i, \Sigma_i$

independent of class label

[So we may ignore this term in constructing discriminant function.]

→ If  $\mu_i, \Sigma_i$  are given, the Gaussian pdf can be uniquely identified.

$\mu_i, \Sigma_i$  can draw any distribution.

•  $-\frac{d}{2} \ln(2\pi)$  : it is constant  
 doesn't depend on class label.

(5)

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) - \frac{1}{2} \ln(\Sigma_i) + \ln P(c_i)$$

If we want to assign a particular class label as  $i$

then  $g_i(x) > g_j(x)$

$$-\frac{1}{2} \ln(2\pi) = \text{ignore: common term}$$

- To compare, making computational easy, considering terms which are not constant w.r.t variability.

## Normal Density and Discriminant Function

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) - \frac{1}{2} \ln(\Sigma_i) + \ln P(c_i)$$

Assume :-

$$\Sigma_i = \sigma^2 I + \nu_i$$

*constant*

$I \Rightarrow$  only diagonal term  $\perp$

$\rightarrow$  off diagonal = 0

Covariance b/w features = 0

⑥

Not Correlated.

Similar to Naive Bayes: one feature doesn't depend on other,  
given class i.

$|\Sigma|_{d \times d}$  : determinant of  $\sigma^2 I$

$|\Sigma| = \sigma^{2d}$  : determinant product of diagonal term d.  
d no. of terms product diagonally.

In 3D+

$$I = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

→ only diagonal elements (variance terms)  
are non-zero.

→ As off diagonal not dependent on each  
other, off diag. covariance: zero

→ It means features of data are statistically independent of each  
other, i.e. no pair of features show any trend.

$$\sum_i \Rightarrow \sigma^2 I$$

$|\Sigma_i| = \sigma^{2d}$  (assuming I to be a  $d \times d$  identity matrix)

$$\Sigma_i^{-1} = \frac{1}{\sigma^2} I$$

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) - \frac{1}{2} \ln(\sigma^2 I) + \ln P(C_i)$$

↑  
constant      independent of class label

→ so we may ignore this term in constructing the discriminant function.

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) + \ln P(C_i)$$

$$\Sigma^{-1} = \frac{1}{\sigma^2} I$$

$$g_i(x) = -\frac{\|x - \mu_i\|^2}{2\sigma^2} + \ln P(C_i)$$

$$= -\frac{1}{2\sigma^2} [x^T x - 2\mu_i^T x + \mu_i^T \mu_i] + \ln P(C_i)$$

↑  
independent of class label

$$x^T x = \|x\|^2$$

As  $x$  is input data, my observation independent of class.  
When I calculated  $x$  didn't know which particular class does it

belongs to.

$x \leftarrow$  not dependent on  $i$ . = observation doesn't depend on class  
. may ignore this term.

$$g_i(x) = -\frac{1}{2\sigma^2} [x^T \mu_i + \mu_i^T \mu_i] + \ln P(C_i) \quad (8)$$

Can't ignore  $\mu_i^T x$  as one term is dependent on  $i$ .

$$\Rightarrow \frac{\mu_i^T x}{\sigma^2} + \ln P(C_i) - \frac{1}{2\sigma^2} \mu_i^T \mu_i$$

$$= g_i(x) \boxed{= w_i^T x + w_{i0}}$$

$$\begin{aligned} w_i^T &= \frac{\mu_i^T}{\sigma^2} \\ w_{i0} &= \ln P(C_i) \\ &= -\frac{1}{2\sigma^2} \mu_i^T \mu_i \end{aligned}$$

→ Discriminant function is a linear function of input.

This is called Linear Machine.

- LDA assumption summary
- (i) distribution is gaussian
  - (ii) features are equally weighted
  - (iii) features doesn't depend on each other.

(9)

### Naïve Baye's Example

$$\begin{matrix} \pi_{\text{no anemia}}, \pi_{\text{anemia}} \\ 0.06 & 0.09 \end{matrix}$$

$$\pi_{\text{no anemia}} > \pi_{\text{anemia}}$$

Conclusion is that person doesn't have anemia as per Naïve Baye's Classifier.

### Naïve Baye's: Slightly Different Example

Sample number	Appetite	Weight (kg)	BP	Class
1	low	61	low	No Anemia
2	low	49	low	Anemia
3	normal	47	low	Anemia
4	low	51	Normal	No anemia
5	normal	50	Normal	Anemia
6	normal	63	low	Anemia
7	normal	58	Normal	No Anemia

Step 1: Calculate prior probability

$$P(A) = \frac{4}{7} \quad P(\bar{A}) = \frac{3}{7}$$

Step 2+ Calculate  $P\left(\frac{\text{feature}}{\text{Anemia}}\right)$ ,  $P\left(\frac{\text{feature}}{\text{No Anemia}}\right)$  from (10)  
the training data.

		Anemia	No Anemia
Appetite	Low	$\frac{2}{4}$	$\frac{1}{3}$
	Normal	$\frac{1}{4}$	$\frac{2}{3}$

$$P\left(\text{Appetite} = \frac{\text{low}}{\text{Anemia}}\right), P\left(A = \frac{\text{normal}}{A}\right), P\left(A = \frac{\text{low}}{A}\right), P\left(\text{Appetite} = \frac{\text{normal}}{\text{no anemia}}\right)$$

All these are calculated

Similarly for BP all possible values are calculated.

		Anemia	No Anemia
BP	Low	$\frac{3}{4}$	$\frac{1}{3}$
	Normal	$\frac{1}{4}$	$\frac{2}{3}$

$$\text{All these are calculated. e.g. } P\left(BP = \frac{\text{low}}{\text{Anemia}}\right) \dots$$

What should I put in weight table entries ?

	Anemia	No Anemia
Weight	?	?

(11)

		Anemia	No Anemia
Weight	61	0	1
	49	1	0
	47	1	0
	51	0	1
	50	1	0
	63	1	0
	58	0	1

Viz Bayes  
fails here

Step 3 + Suppose, I observe a test data with normal appetite, 56 kg weight, and low BP. predict if person has anemia.

Let's first evaluate choice of Anemia:-

$$\pi_{\text{Anemia}} = P(\text{appetite} = \frac{\text{normal}}{\text{Anemia}}) P(\text{weight} = \frac{56 \text{ kg}}{\text{Anemia}}) P(\text{BP} = \text{low} \mid \text{Anemia})$$

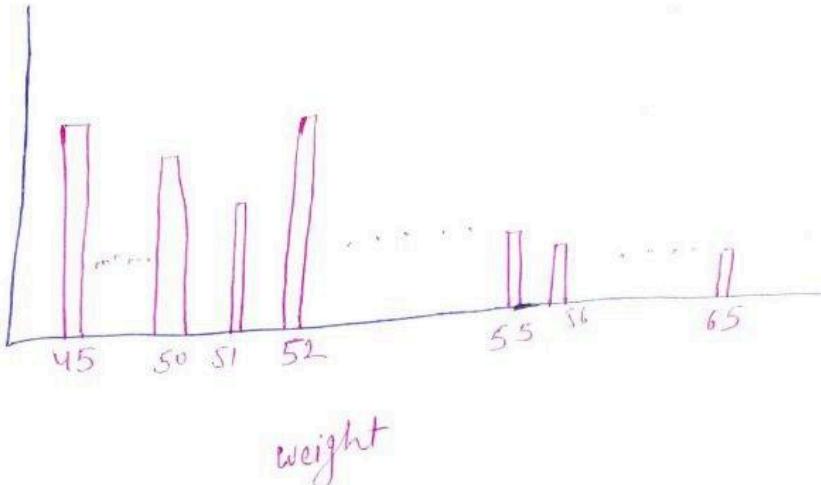
$$\frac{1}{4} \times \frac{1}{4} \times \frac{4}{7}$$

From my training data, I can't find  $P(\text{weight} = 56 \text{ kg} \mid \text{Anemia})$

If see, Table has no entry at  $P(\text{wt} = \frac{56 \text{ kg}}{\text{anemia}})$ .

(12)

$p(\text{weight}$   
anemia)



Accuracy compromised if we consider range.

| if we are given Probability / training data.

Goal:- to find distribution that maximally fits data.

\* See distance which lies nearby to which value.

e.g. given random variable value = weight,  
also disease is Anemia.

Can calculate Likelihood.

How to find probability density function?

Q:- Question is:- How to find probability density from training data.

so that we can estimate likelihood  $P(\text{weight} = 56 \text{ kg} / \text{anemia})$ ?

$$P(A_1) = \frac{P(z/A_1) P(A_1)}{P(z)}$$

## Maximum Likelihood Estimation :-

(13)

If I have c number of classes  $C_1, C_2, \dots, C_c$  in my dataset

$S_1$ : Training samples of class  $C_1$

$S_2$ : Training samples of class  $C_2$

$\vdots$   
 $S_c$ : Training sample of class  $C_c$ .

I have to find :-  $P\left(\frac{x}{C_{l_i}}\right)$  such that  $P\left(\frac{x}{C_{l_i}}\right)$  is maximized  
when I use training data of class  $C_{l_i}$ .

→ We assume that  $P\left(\frac{x}{C_{l_i}}\right)$  :- Known parametric form

• If  $P\left(\frac{x}{C_{l_i}}\right)$  is Gaussian, it is completely specified by

mean  $\mu_i$  and covariance matrix  $\Sigma_i$

parameter =  $\mu, \Sigma$  of gaussian dist<sup>n</sup>

If I can find out parameters  $\Theta_i$ , I can find  
out  $P(x|C_{l_i})$

- Find  $\theta_i$  which maximally fits data distribution.

(14)

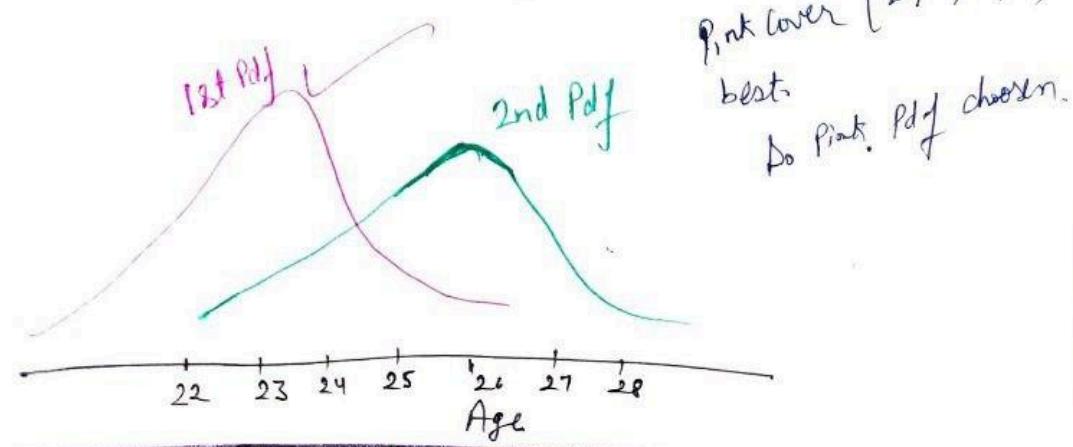
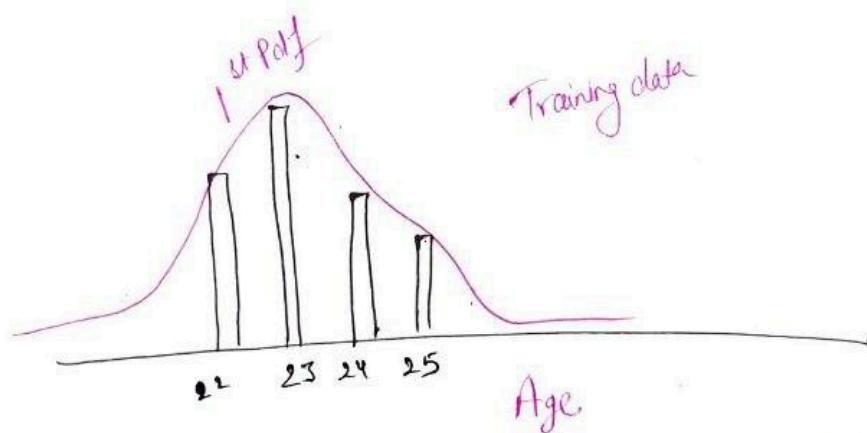
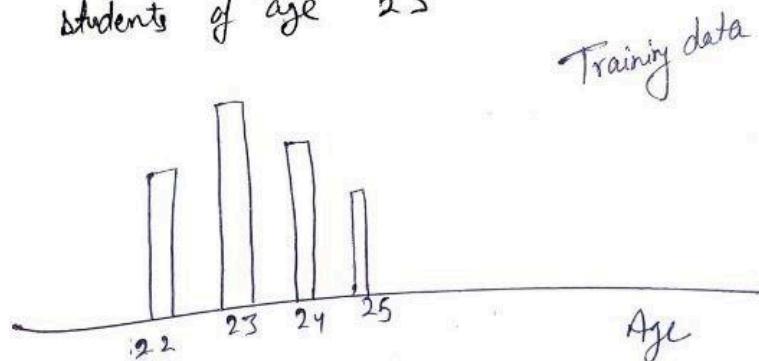
e.g. In the classroom

30 students of age 22

35 students of age 23

30 students of age 24

25 students of age 25



$P\left(\frac{\text{age} = 22}{s}\right) + \left(\frac{\text{age} = 23}{s}\right) P\left(\frac{\text{age} = 24}{s}\right) + \left(\frac{\text{age} = 25}{s}\right)$  is maximum for  
 which curve.  
 (15)  
 likelihood values.

- We assume that  $P\left(\frac{x}{C_i}\right)$  has a known parametric form.
- If  $P\left(\frac{x}{C_i}\right)$  Gaussian, it is completely specified by mean  $\mu_i$  and covariance matrix  $\Sigma_i$  (call them parameters  $\theta_i$  of distribution).
- Create → Find parameters  $\theta_i$ ,  
 Given → information of Training sample in sample i.e  $S_i$   
 such that  $P\left(\frac{x}{C_i}\right)$  is maximized.  
 → using training data of class  $C_i$ .

We assume information in  $S_j$  doesn't affect  $\theta_i$   
 if  $i \neq j$ .

$(\lambda)(\lambda)(\lambda)$   
 Likelihood multiplication

Find curve = product of terms that maximizes

(16)

- our goal is to use information from training sample in  $S_i$
- To find parameters  $\theta_i$  such that  $P\left(\frac{x}{C_{l_i}}\right)$  is maximized when I use training data of class  $i$ .
- We have to maximize  $P\left(\frac{x}{C_{l_i}}, \theta_i\right) \forall i$
- Since  $\theta_i$  is parameters corresponding to  $C_{l_i}$  only, we can say we have to find maximize  $P\left(\frac{x}{\theta_i}\right) \forall i$ .

Our goal is to use information from training sample in  $\mathcal{S}_i$  to find parameters  $\theta_i$  such that  $P\left(\frac{x}{C_i}\right)$  is maximized when we use training data of class  $C_i$ .

→ we have to maximize  $P(x/C_i, \theta_i) + i$

→ Since  $\theta_i$  are parameters corresponding to  $C_i$  only

we have to maximize  $P\left(\frac{x}{\theta_i}\right) + i$

- Since data from another class doesn't effect  $\theta_i$  distribution  
So considering whole data distribution will not effect  $\theta_i$  of class  $i$ .

$$P\left(\frac{x}{\theta}\right)$$

- In this, we assume information in  $\mathcal{S}_j$  doesn't affect  $\theta_i$  if  $i \neq j$ .

So even if we consider entire training dataset  $\mathcal{S}$  instead of just  $\mathcal{S}_i$ , parameter  $\theta_i$  will only be influenced by  $\mathcal{S}_i$ .

(2)

Likelihood is == probability distribution  
 { It can be binomial, lognormal ... }

→ A well defined distribution can write statistical in Composed form.

{ If 1 dimensional data = variance  
 multiple dimensional = Co-variance }

$\theta_i$  = parameter of pdf of particular class i.

[ In  $\theta$  parameter is not necessarily two, it varies,  
 diff. pdf has diff no. of parameters in  $\theta$  ]

e.g. for gaussian:  $\mu, \Sigma$

we can interchangeably say  $C_{li} \sim \theta_i$

→ Every training data is independent.

## Maximum likelihood Estimation

(3)

If we have  $N$  training samples  $x_1, x_2, \dots, x_N$   
goal :- We want to maximize likelihood of each of training samples.

So we want to maximize :-

$$P\left(\frac{x_1}{\theta}\right), P\left(\frac{x_2}{\theta}\right), \dots, P\left(\frac{x_N}{\theta}\right)$$

We have to find  $\theta$  that maximizes the product:-

$$\boxed{P\left(\frac{x_1}{\theta}\right) P\left(\frac{x_2}{\theta}\right) \dots P\left(\frac{x_N}{\theta}\right) = \prod_{k=1}^N P\left(\frac{x_k}{\theta}\right)}$$

We have to find  $\theta$  that maximizes product  $\propto \theta$  :-

Step 1 :- Take log ln

Step 2 :- Equate to  $\theta$  To maximize find gradient

Step 3 :- Equate gradients to zero

Step 4 :- find value of ' $\theta$ '

It is not necessary that while find  $\theta$  all times  $\Sigma$

is calculated,

If given  $\Sigma$  we need find  $\theta$  only.

Find  $\theta$  that maximizes product :-

④

$$P\left(\frac{x_1}{\theta}\right) \cdot P\left(\frac{x_2}{\theta}\right) \cdots P\left(\frac{x_N}{\theta}\right) = \prod_{k=1}^N P\left(\frac{x_k}{\theta}\right)$$

Step 1:- It is equivalent to finding  $\theta$  that maximizes :-

$$l(\theta) = \ln \left( \prod_{k=1}^N P\left(\frac{x_k}{\theta}\right) \right) = \sum_{k=1}^N \ln P\left(\frac{x_k}{\theta}\right)$$

Step 2:- To maximize, find gradient of  $l(\theta)$  and equate to zero

$$\nabla_{\theta} l(\theta) = 0$$

If it gives min, max ✓

else calculate double derivative.

→ If multiple max, choose maximum among all maxima.  
Find highest maxima.

- If we have m no. of parameters  $\theta_1, \theta_2, \dots, \theta_m$

we have to do :-

$$\begin{bmatrix} \frac{\partial l(\theta)}{\partial \theta_1} \\ \frac{\partial l(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial l(\theta)}{\partial \theta_m} \end{bmatrix} = 0$$

Let's assume pdf has Gaussian distribution with mean  $\mu$  and covariance matrix  $\Sigma$ . ⑤

- Assume  $\Sigma$  is known, we have to find  $\mu$

$$P\left(\frac{x_k}{\mu}\right) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2}\right]$$

$\rightarrow$  Take  $\boxed{\mu=0}$  parameter to be found

$$\ln P\left(\frac{x_k}{0}\right) = -\frac{1}{2} \ln(2\pi)^d |\Sigma| - \frac{1}{2} (x_k - 0)^T \Sigma^{-1} (x_k - 0)$$

$x_k$  = Training data  
feature vector containing all dimension

[Naïve Bayes consider all dimension separately]

$\rightarrow$  we have to find  $\theta_i$  that maximizes  $P\left(\frac{x}{\theta_i}\right) \cdot s_i$   
We assume that information in  $s_j$  doesn't affect  $\theta_i$  if

So even if we consider entire training dataset  $S$  instead of  $s_i$ ,

parameter  $\theta_i$  will only be influenced by  $s_i$

we can say that we have to find  $\theta$  that maximizes  $P(\frac{x_k}{\theta})$

(6)

$$P\left(\frac{x_k}{\theta}\right) = \left(\frac{1}{2\pi}\right)^{\frac{N}{2}} |\Sigma|^{\frac{1}{2}} \exp\left(-\frac{(x_k - \mu)^T \Sigma^{-1} (x_k - \mu)}{2}\right)$$

Take  $\mu = \theta$

$$l_k(\theta) = \ln P\left(\frac{x_k}{\theta}\right) = -\frac{1}{2} \ln(2\pi) |\Sigma| - \frac{1}{2} (x_k - \theta)^T \Sigma^{-1} (x_k - \theta)$$

$$l(\theta) = \sum_{k=1}^N \ln P\left(\frac{x_k}{\theta}\right) = \sum_{k=1}^N l_k(\theta)$$

It can be shown as to take derivative : 1st term independent of  $\theta$

$$\nabla_\theta l_k(\theta) = \Sigma^{-1} (x_k - \theta)$$

$$\nabla_\theta l(\theta) = \sum_{k=1}^N \nabla_\theta l_k(\theta) = \sum_{k=1}^N \Sigma^{-1} (x_k - \theta)$$

For MLE :- equate to zero +

$$\nabla_\theta l(\theta) = \sum_{k=1}^N \nabla_\theta l_k(\theta) = \sum_{k=1}^N \Sigma^{-1} (x_k - \theta) = 0$$

$$\text{Since } \Sigma^{-1} \neq \theta , \quad \sum_{k=1}^N (x_k - \theta) = 0$$

$$\left( \sum_{k=1}^N x_k \right) - N\theta = 0$$

(7)

$$\theta = \frac{1}{N} \sum_{k=1}^N x_k$$

= mean when all data are  
Equally likely.

e.g. Consider a Univariate Gaussian with unknown  $\mu + \sigma^2$ .

Find values of  $\mu, \sigma^2$  for MLE :-

Univariate specified gaussian by 1 variable.

→ If 2 variable = Take partial derivative.

$\theta$  = for equally likely

$$P\left(\frac{x_k}{\theta_1, \theta_2}\right) = \frac{1}{(2\pi\theta_2)^{\frac{1}{2}}} \exp\left[-\frac{(x_k - \theta_1)^2}{2\theta_2}\right]$$

Assuming  $\theta_1 = \mu$ ,  $\theta_2 = \sigma^2$

$$P\left(\frac{x_k}{\theta}\right) = \frac{1}{(2\pi\theta_2)^{\frac{1}{2}}} \exp\left[-\frac{(x_k - \theta_1)^2}{2\theta_2}\right]$$

$$l_k(\theta) = \ln P\left(\frac{x_k}{\theta}\right) = -\frac{1}{2} \ln(2\pi\theta_1) - \frac{1}{2\theta_1} (x_k - \theta_1)^2 \quad (3)$$

$$\nabla_{\theta} l_k(\theta) = \begin{bmatrix} \frac{\partial l_k(\theta)}{\partial \theta_1} \\ \frac{\partial l_k(\theta)}{\partial \theta_2} \end{bmatrix} = \begin{bmatrix} \frac{1}{\theta_2} (x_k - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{1}{2\theta_2^2} (x_k - \theta_1)^2 \end{bmatrix}$$

while equating to zero

$$\sum_{k=1}^N \nabla_{\theta} l_k(\theta) = 0$$

$$\sum_{k=1}^N \frac{1}{\theta_2} (x_k - \theta_1) = 0$$

$$\sum_{k=1}^N \left( -\frac{1}{2\theta_2} + \frac{1}{2\theta_2^2} (x_k - \theta_1)^2 \right) = 0$$

Assuming  $\theta_2 \neq 0$

$$\theta_1 = \sum_{k=1}^N \frac{x_k}{N} \quad \theta_2 = \sum_{k=1}^N \frac{(x_k - \theta_1)^2}{N}$$

This is estimate of  $\mu$

This is estimate of  $\sigma^2$

Q:-  $P\left(\frac{x_k}{\theta}\right) = \begin{cases} \theta \exp(-\theta x_k) & \text{if } x_k \geq 0 \\ 0 & \text{otherwise.} \end{cases}$  ①

Sol<sup>n</sup>r goal :- To find parameter  $\theta$  that maximizes likelihood of observing given data.

given probability distribution :—

$$P\left(\frac{x_k}{\theta}\right) = \begin{cases} \theta \exp(-\theta x_k) & \text{if } x_k \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Step 1 For  $n$  independent data points  $x_1, x_2, \dots, x_n$ , function  $L(\theta)$  is product of individual probability, —

$$L(\theta) = \prod_{k=1}^n P\left(\frac{x_k}{\theta}\right)$$

for given pdf, likelihood becomes →

$$L(\theta) = \prod_{k=1}^n \theta \exp(-\theta x_k)$$

Since  $\theta$  and  $x_k$  are independent :— Simplifying n

$$L(\theta) = \theta^n \exp\left(-\theta \sum_{k=1}^n x_k\right)$$

Step 2 To simplify maximization problem, take log of likelihood ↗

$$\log L(\theta) = \log \left( \theta^n \exp \left( -\theta \sum_{k=1}^n x_k \right) \right)$$

Simplifying ↗

$$\log L(\theta) = n \log \theta - \theta \sum_{k=1}^n x_k$$

Step 3 Differentiate log-likelihood w.r.t  $\theta$ , set it equal to 0,  
find maximum ↗

$$\frac{d}{d\theta} \log L(\theta) = \frac{n}{\theta} - \sum_{k=1}^n x_k = 0$$

Step 4 Solve for  $\theta$  ↗

$$\theta = \frac{n}{\sum_{k=1}^n x_k}$$

It means that  $\theta$  is reciprocal of average of observed data points  $x_k$ .

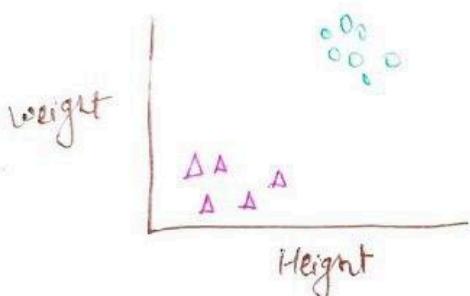
## Lecture - 18

14-80pt-2024

①



Some sample data  
of a few animals



In clustering, we create group of samples in such a way that sample in group are more similar to each other than sample in other groups

Let's 1st see Notion of Distance →

There are multiple types of distances : —

Given : 2 points  $(x_1, y_1)$   $(x_2, y_2)$  in 2D plane.

Goal : Is to find distance

Approaches are many : —

euclidian distance : —

$$d_E = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

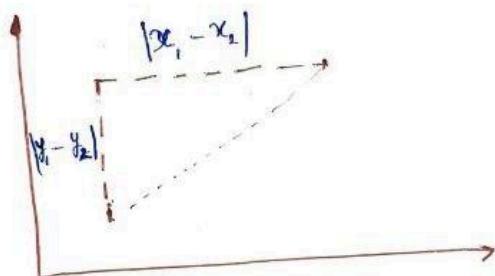
data point in 2 dimension feature.

(2)

When we move to higher dimension :-

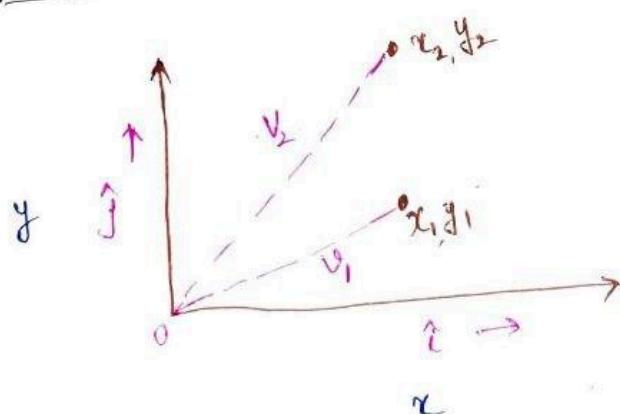
Eucleadian distance tend to loose its significance.

(ii) Absolute distance / Manhattan distance :-



$$d_m = |x_1 - x_2| + |y_1 - y_2|$$

(iii) Mostly for higher dimensions ← we prefer cosine distance.



- Distance of  $x_1, y_1$  and origin is  $\theta_1$
- $x_2, y_2$  and origin is  $\theta_2$

It indicates Angular distance.

(3)

$$v_1 = x_1 \hat{i} + y_1 \hat{j}$$

$$v_2 = x_2 \hat{i} + y_2 \hat{j}$$

$$\vec{v}_1 \cdot \vec{v}_2 = |\vec{v}_1| \cdot |\vec{v}_2| \cos\theta$$

$$\cos\theta = \frac{\vec{v}_1 \cdot \vec{v}_2}{|\vec{v}_1| \cdot |\vec{v}_2|} = \frac{x_1 x_2 + y_1 y_2}{\sqrt{x_1^2 + y_1^2} \sqrt{x_2^2 + y_2^2}}$$

Cosine distance tells how much I need to rotate  $v_1$  to reach  $v_2$

→ It is indicator of direction mostly.

## Clustering

Some examples of it

• Targeted Advertising to group people of similar interest.

→ Recommendation system

Clustering Problem: don't assign label

Classification " : It can assign label

## Clustering Techniques

• Centroid model : K mean, K mediod

• Graph based model : Spectral clustering

(4)

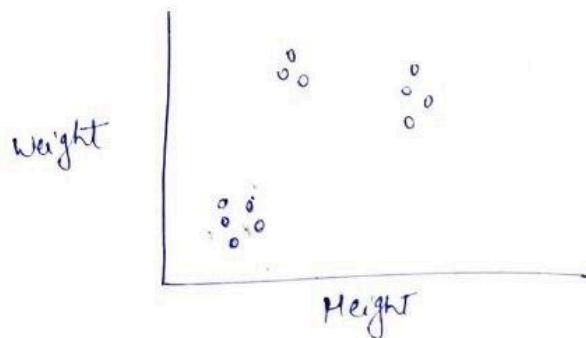
• Distribution model : Gaussian Mixture Models

• Density models : DBSCAN

• Connectivity Based : Hierarchical clustering

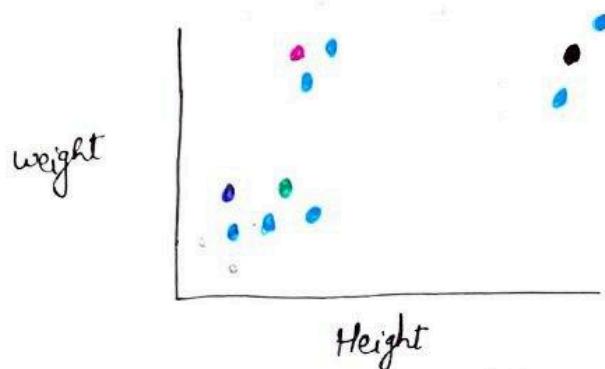
• Neural models : Self-organizing map

⑩ K-means :- Take data points :-



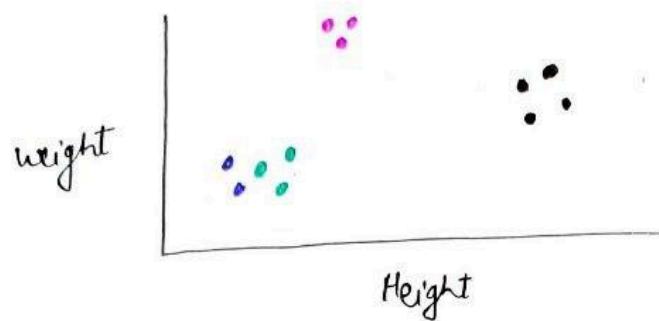
Let's say I want  $m$  clusters , ie =  $m=4$

• Randomly initialise  $m=4$  cluster centers.

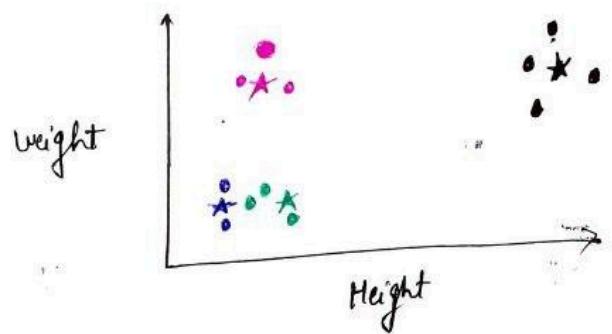


- Assign each of other points to nearest cluster center.

(5)

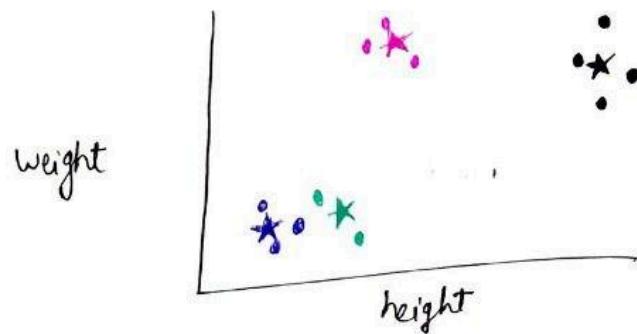


→ Recalculate to cluster mean.



→ Reassign points based on updated mean.

→ Keep in loop to mean calculation, reassign point until there is no change in assignment.



(6)

- Randomly assign points to 4 points centers.
- Each data point is d-dimensional data.

$$\begin{aligned} x^1 &\rightarrow x_1^1, x_2^1, \dots, x_N^1 \\ x^2 &\rightarrow x_1^2, x_2^2, \dots, x_N^2 \\ &\vdots \\ x^N &\rightarrow x_1^N, x_2^N, \dots, x_d^N \end{aligned}$$

Coordinate of centroid for cluster 1, Cl<sub>1</sub>

$$Cl_1 = \{c_1^1, c_2^1, \dots, c_d^1\}$$

$$c_1(1) = \frac{x_1^1 + x_2^1 + \dots + x_p^1}{p}$$

If p points are there

$x_i^1 \rightarrow$  no of data point = p data point

$x_i^1 \rightarrow$  dimension e.g. d dimension

- No guarantee if there is convergence or not, infinite loop.
- We set some threshold when less than i points changes their cluster assignment.

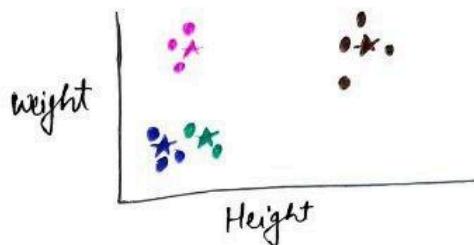
(7)

• We can tune it to find  $K$ .

↳ There is no full-proof method to find  $K$  (# of clusters)

Steps for clustering :-

(i) Let's say, I want  
 $m$  clusters



(ii) Let's take  $m=4$

(iii) Randomly initialize  $m=4$  cluster centers

(iv) Assign each of other points to nearest cluster center

(v) Recalculate cluster mean

(vi) Reassign points based on updated mean

(vii) Go to step 5, if there is any change in the assignment, Otherwise, stop.

Limitations :-

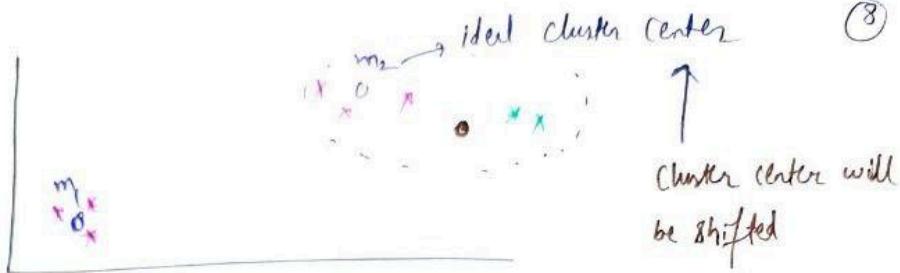
(i) Not guaranteed to converge

(ii) Difficult to find  $K$

(iii) Sensitive to outliers

(iv) may be significantly affected by initialization

eg ↴

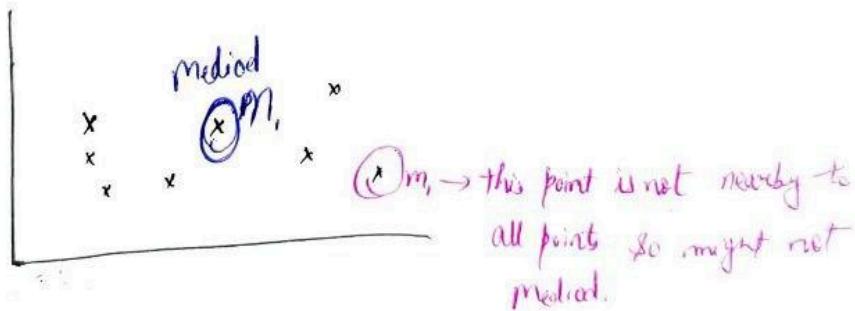


(i) If data mean is far away from  $m_1, m_2$  so mean shifted  
one of the way is to discard these points.

→ K mean significantly affected by initialization.

- Mean may not lie on actual data point.
- Medoid lies on actual data point.

Mediod



$m_i$  = mediod → it is nearby to all points (centre)

Mediod = minimum distance

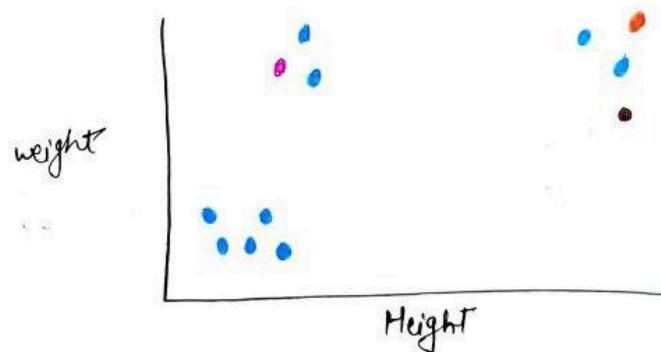
Mediod + we have to options +

(9)

- (a) To assign points randomly to any point
- (b) To assign point to same cluster.

### Clustering Technique : K-mediod :-

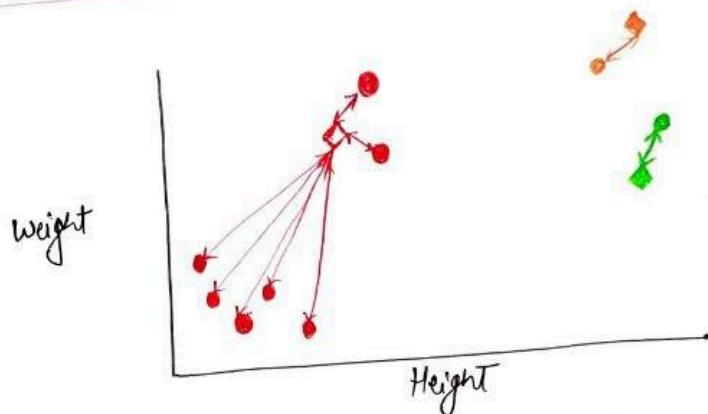
- A mediod is a point in the cluster from which sum of distances to other data points is minimal.
- A Mediod is a point in cluster from which dissimilarities with all other points in cluster are minimal.
- Instead of centroid as reference points in k-means algorithm, k-mediod algo takes a Mediod as reference point.



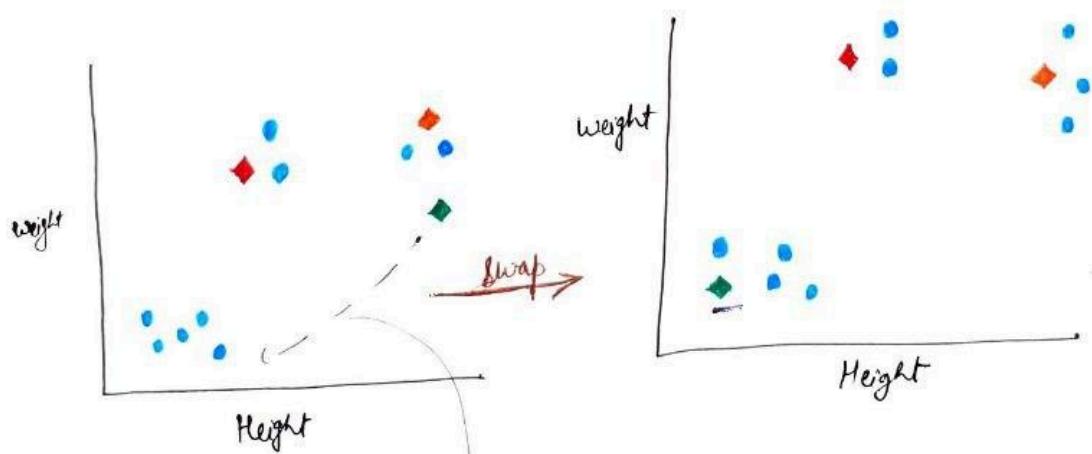
- (i) Take the data points
- (ii) Randomly assign mediods
- (iii) For each non-mediod data points, find nearest mediod and assign data point to the corresponding cluster.

(v) Calculate distance of every data point from corresponding mediod. ⑩  
Sum of all these distances is called Cost.

$$\text{Cost} = \text{total red distance} + \text{total orange distance} + \text{total green distance}$$

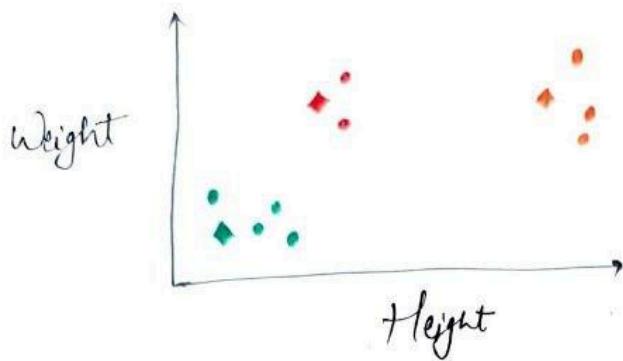


Rohmku representants = mediods



if restricted to same cluster  
green couldn't come at Bottom.

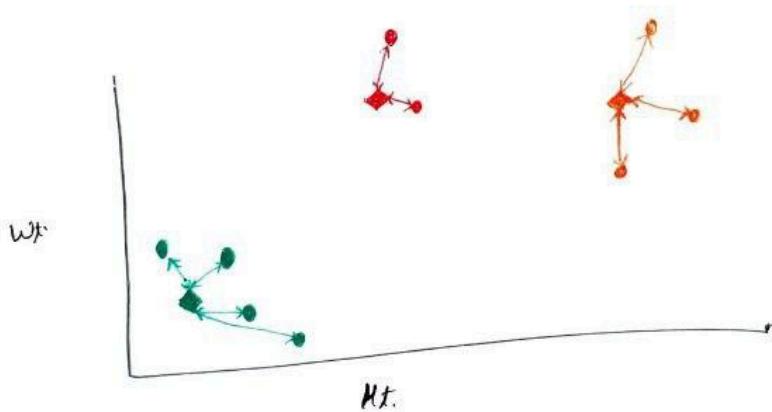
- With new Medoids, create the clusters using nearest medoids. (11)



→ with the new Medoids, recalculate the cost.

- Calculate distance of every data point from corresponding Medoid.  
Sum of all distances is called Cost.

$$\text{Cost} = \text{total red distance} + \text{total yellow distance} + \text{total black distance}$$



If

new cost > old cost → discard new medoids

↓  
go back to old medoids.

Algorithm converges.

else , keep new Mediod and redo random swapping.

(12)

Steps ↗

- (i) Select K random points from the dataset.
- (ii) Assign data points to the cluster of nearest mediod.
- (iii) Calculate total sum of distances of data points from their assigned mediod for each mediod.
- (iv) Swap a non-mediod with mediod point and recalculate cost.
- (v) Undo swap if recalculated cost with new mediod exceeds previous cost.

Until new cost  $>$  old } Terminate condition

This applies when already tried for 3-4 times if always cost is greater than old one then stop.

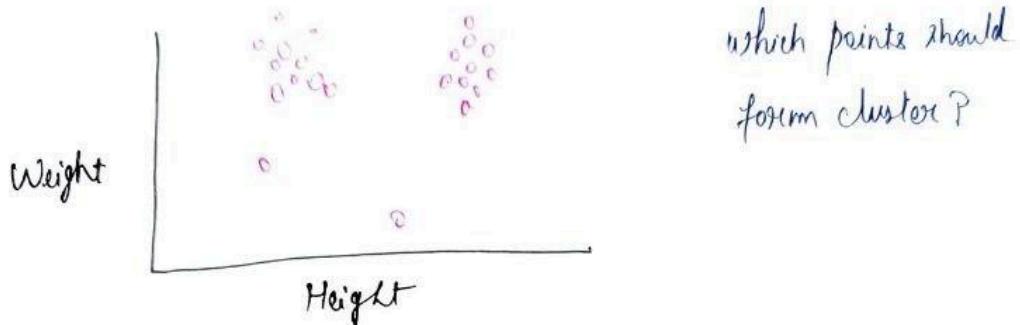
To get rid of K mean, K mediod initialization based method; —

Density Based clustering.

## Lecture - 19

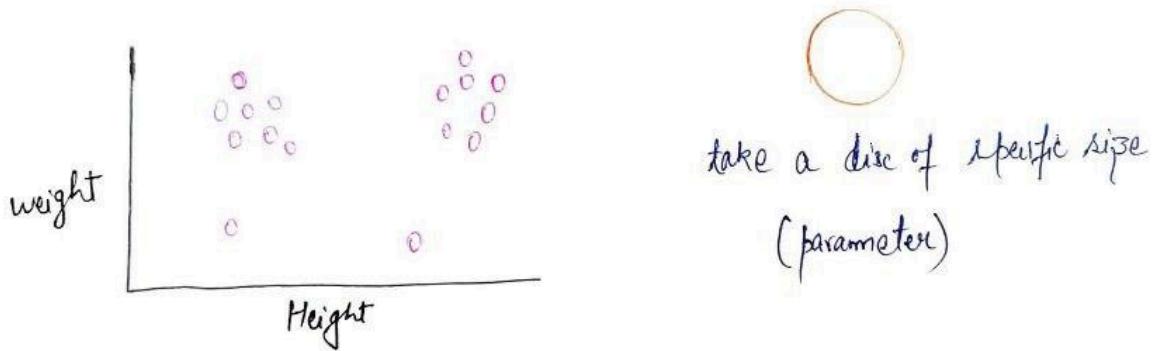
Date : 17/09/2024 ①

### Clustering Techniques : Use of Density +



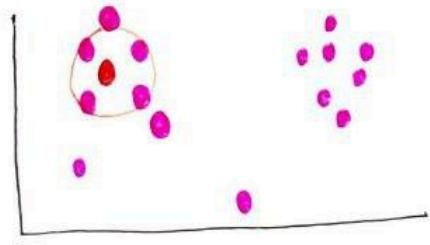
which points should  
form cluster?

- Take a disc of specific size (parameter)



take a disc of specific size  
(parameter)

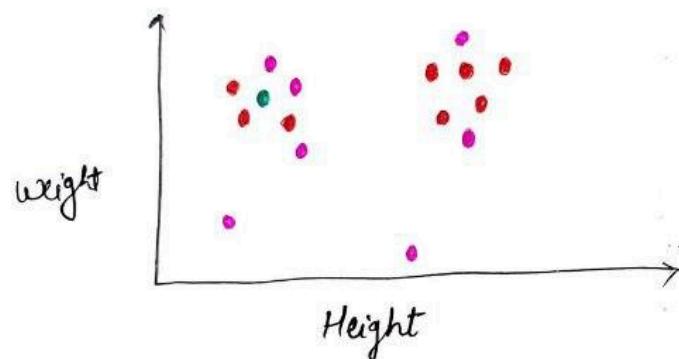
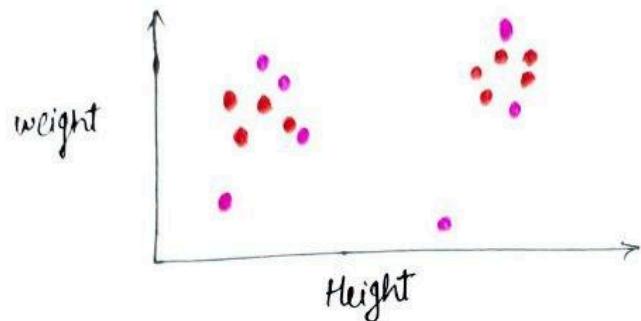
- Put it across every point.
- If there are  $m$  no. of points within disc when positioned on point  $x_i$ ,  
The point  $x_i$  is called as Core Point.  $m$  is parameter



- Put it across every point.  
If there are  $m$  no. of points within the disc when positioned on point  $x_i$ ,

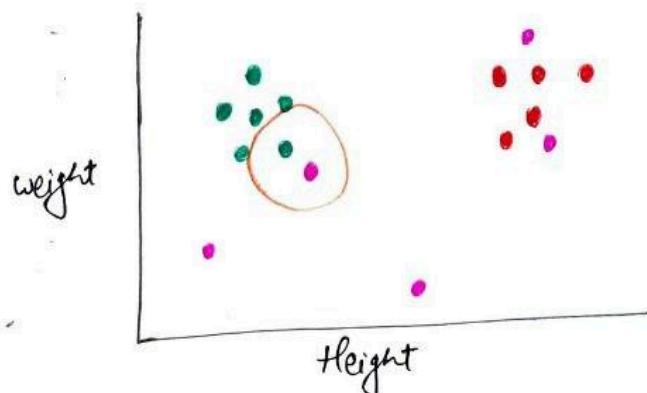
→ Point  $x_i$  will be called a core point.  $m$ - parameter.

②



Red : core points

Assign a cluster label in a core point (green)



→ Assign a cluster label in a core point (green)

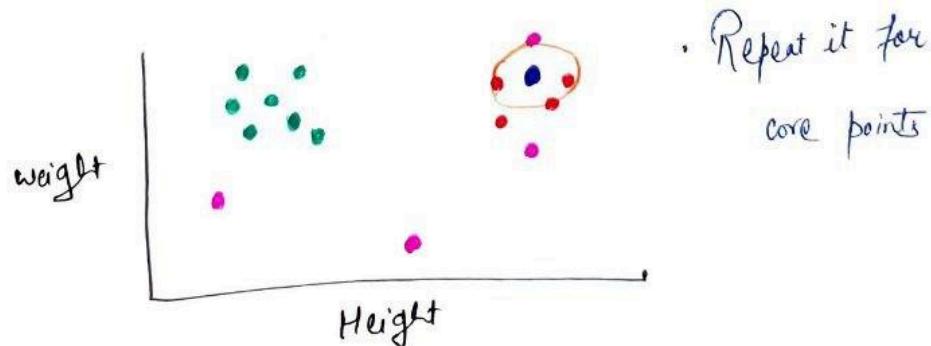
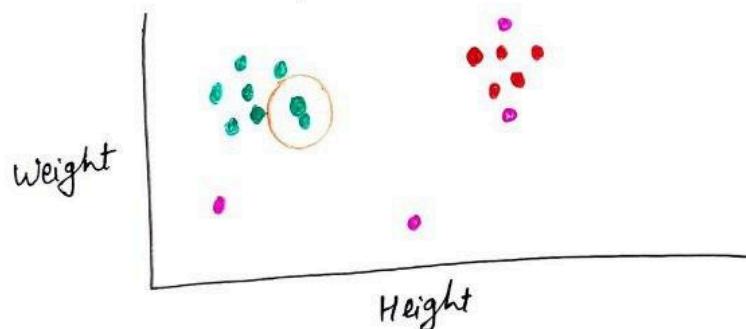
• Put every point within the disk inside the 1st cluster.

- Apply the disk from every point in the cluster. ③
- e.g. Customer has no idea, how many customer for delivery boy.  
So based on density based clustering.
- Based on density, I am making ~~top, bottom, right~~ left Clustering.

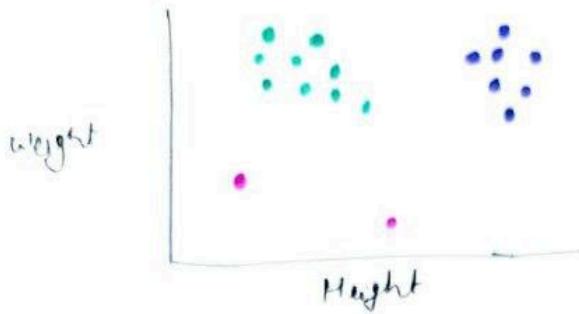
→ If no. of point inside disc is greater than outer, this point is core points.

"Core point = Crowded region", if  $m = \underline{\text{some threshold}}$ .

→ Apply disk from every point in cluster.



(4)



when we are done with all core points and left with only non-core points that can't be added to any clusters, we are done.

These non-core points is called outliers

- It's 2D, so considering circle, hypersphere else.
- Bottom two aren't core point as doesn't contain point greater than equal to  $m$ .

(i) Find core point

↓  
based on putting circle disc  $m$  points, check if falling as core point, randomly choose any core, find its neighbour

↓  
repeat.

(ii) Repeat with different radius values find which ' $r$ ' considers more point.

(iii) Choose intuitive value of  $r$

There are two parameters  $\sigma$ ,  $m$ ,  $r_L$