

08/08/22

# Data Science

The Great Hack

In God we trust, all others bring data

- W. E. Deming.

→ Data acts as evidence, which helps us to trust others.

\* 3V's of Data → Velocity, Volume & Variety.

\*\* Data Collection → Data Storage → Processing.

\* ETL vs OLAP

Storing data according to use

10/08/22

Case Study of Netflix → organised competition in 2006, challenging to make recommendation system.

In 2010, demonization problem; two UG students  
find which person watch which movie.  
→ fined.

• Throttling Issue → ISP caps the network to

took place in  
US in 2010

found using  
Data Analysis

Slow the speed of network.

{ In case of higher usage, it ~~was~~ done }  
of Netflix.

# Data Science in the Business Context {Business Analytics}

- change is constant in business Environment.
- So, companies have to agile & take quick strategic, tactical, operational decision.

22/08/22

Mode → not affected by outliers.

- only talk abt. value occurring max<sup>m</sup> no. of times.

\* Using in Election or in situation where you want to cater max no. of people.

Median → Value that divides the sorted data in equal half.

- Not affected by outliers. ( $50^{\text{th}}$  percentile)

↳ To get Median from data (Even → avg does not belong to data)  $(\frac{n}{2})$

↳ The  $50^{\text{th}}$  percentile of a list of no. is the smallest no. that is at least as large as  $50\%$  of list.

\* Less Informative as not effected by extremes; so if after median = 3

average = 100  
present  
Notable  
to give info about this

Eg → list → 0, 2, 4, 7, 7, 12

Find  $50^{\text{th}}$  &  $25^{\text{th}}$  percentile.

$\frac{1}{2} \times 6 = 3 \Rightarrow$  At least 3 values should be less than the median.

$$\boxed{\text{Median} = 7}$$

$$\underline{\underline{25^{\text{th}} \text{ percentile} = 4}}$$

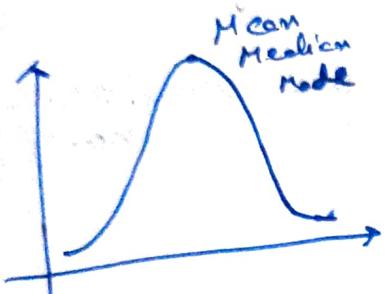
Eg → To decide salary of a Role → Median of salary for the particular role in different competitive companies.

Mean → Affected by the outlier

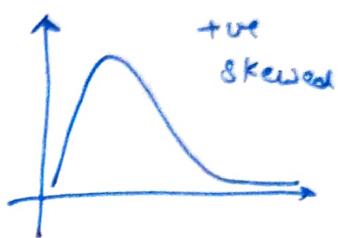
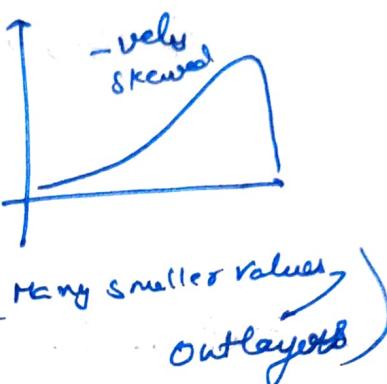
- Holistic picture of set of values.

# Some Questions to be Answered ↴

① Skewness → A long tail at a particular end of graph/distribution.



Normal Curve  
Perfect Distribution  
not skewed



② Cross Sectional vs. Longitudinal.

Can't compare

~~Eg → Diff set of people  
of diff age grp~~

# Measures of Location → If know the avg.

Can I answer  
questions related to  
distribution of  
Data.

Eg → If avg of class is 30%.

what portion of students scored

more than 80%.

Ans →  $\frac{1}{3}$  at max

Can be answered by Markov's Inequality.

→ If the list consists of only non negative numbers then the proportion of entries that are at least as large as  $k$  times the average is at most  $\frac{1}{k}$ .

Q) How much of the portion of people are of greater than half of average?

Ans →  $\frac{1}{1/2} = 200\%$  of people

It is correct. But is quite obvious  
 $\frac{1}{k}$  gives loose bound sometimes.

\* So, To get the tighter bound only the mean is not enough.

Measure of Spread → How much data is deviated from mean

### Standard Deviation.

$$SD = \sqrt{\frac{1}{n} \sum_{i=1}^m (x_i - \bar{x})^2}$$

$\downarrow$   
K/a      sum of deviation from average  
root mean  
Square

\*\* Why Square Instead Of Modulus? ?

### # Chebychev's Inequality

The proportion of entries that are  $k$  or more SDs away from the average is at most  $\frac{1}{k^2}$ .

Q) Suppose the class avg is 30% &  $SD = 10$   
what portion of students scored more than 30%

$$\left\{ \frac{1}{36} \right\}$$

24/08/22

Given that a value is greater than  
Can we predict about mean & mode

## Correlation Coefficient

$$\gamma = \frac{1}{n} \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{\sigma_x \sigma_y}$$

Steps → ① Given  $x \& y$  &  $|x| = |y| = n$

② Convert  $x \& y$  to standard unit.

### Why Std. Unit??

- also known as Z-Score
- have no unit. need this to compare values of different units.
- So call it Std. Unit.
- $\gamma \in [-1, 1]$ 
  - Strong +ve relation
  - No Relation

• find mean & SD

• To convert to std. unit.  $\left\{ \frac{x_i - \bar{x}}{\sigma_x} \right\}$

How far it is from mean in terms of std. deviation

③ Mult. Corresponding pairs of std. units.

④  $\gamma$  is avg of above products.

Correlation coefficient

\* Association is not causation → Just Happening,

{ Eg → Rain & Exam }

Not like { Not causal Relation }  
{ Due to inc of x, y is inc } skip

\*  $\gamma$  measures linear association

\* Even one outlier might affect  $\gamma$

# Given mean =  $m$  & std. deviation =  $\sigma$

if n added to all values in list → mean =  $m+n$  & Std. Dev =  $\sigma$

if all values mult by  $gn$  → mean =  $mn$  & std Dev =  $\sigma gn$

# # Random Variables →

$\{68.95 \rightarrow 99\% \text{ rule}\} \rightarrow \text{Normal Distribution}$

$$\begin{array}{c} + \\ 68.95 \\ + \\ \sigma = 20 \\ + \\ 30 \end{array}$$

Weight of people

Small deviations from mean

Throughout of 99% of nodes → Power Large Dictionary Supporting.

## 26/08/22 Types of Prob. Distributions

- Uniform →  $f(x) = \frac{1}{b-a} \quad x \in [a, b]$  Eq. → Use least significant digit to sort the 3 digits.  $b-a$  → Code: 143 → 3
- Bernoulli → toss coin once.  $f(0) = 1-p$  &  $f(1) = p$ .  $B(n, p)$  →  $f(k) = {}^n C_k p^k (1-p)^{n-k}$ ;  $\mu = np$  &  $\sigma^2 = np(1-p)$
- Binomial Repeat several times independently.  $f(k) = {}^n C_k p^k (1-p)^{n-k}$ ;  $\mu = np$  &  $\sigma^2 = np(1-p)$
- Poisson for very large  $n$  in  $B(n, p)$   $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$ ;  $\mu = \lambda$  &  $\sigma^2 = \lambda$  \* Here, no of Success/people can be any no.
- Geometric no. of attempts till we get success.  $f(k) = (1-p)^{k-1} p$ ;  $F(k) = 1 - (1-p)^k$ ;  $\mu = \frac{1}{p}$ ;  $\sigma^2 = \frac{1}{p^2}$
- Exponential →  $f(x) = \lambda e^{-\lambda x}$ ;  $F(x) = 1 - e^{-\lambda x}$ ;  $\mu = \frac{1}{\lambda}$ ;  $\sigma^2 = \frac{1}{\lambda^2}$
- Normal Distribution →  $f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$  bell shaped curve  $= \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{x^2}{2}}$  \* Z-table

31/08/22

\* A pattern is a pt. on D-dimensional space  
where  $D \rightarrow$  no. of features present in  
Dataset

- To find the Closeness/Similarity among the two patterns
  - ↳ measure distance b/w two Patterns  $\underline{d(P_1, P_2)}$
  - ↳ Larger distance  $\Rightarrow$  Much dissimilar.

\* Many types of distances can be measured  $\rightarrow$  i) Euclidean Distance

Eg →

$$d(P_1, P_2) = \sqrt{(120-100)^2 + (5-3)^2} \\ = 10.$$

	$f_1$	$f_2$
$P_1$	160	3
$P_2$	170	5
$P_3$	160	5

Here first attribute,  $f_1$ , is dominating on the dist. value  
 $\Rightarrow f_2$  is insignificant over distance  
 So,  $f_2$  has no ~~not~~ larger effect on the outcome

So, To avoid the dominance Effect, we need to normalize all the data  $\rightarrow$  so that they come on same scale.

Normalization  $\rightarrow \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^D$

$$\text{Max}_j = \text{Max} \{x_{1j}, x_{2j}, \dots, x_{nj}\} \quad \text{where } j \rightarrow \text{attribute}$$

so  $x_{ij}$  = Value of  $P_i$  at attribute  $j$

$$\text{Min}_j = \text{Min} \{x_{1j}, x_{2j}, \dots, x_{nj}\}$$

$$y_{ij} = \frac{x_{ij} - \text{Min}_j}{\text{Max}_j - \text{Min}_j}; i = 1 \text{ ton. } \text{ hence } y_{ij} \in [0, 1] \quad \text{for } j = 1 \text{ to } D$$

Remove/Reduce the effect of dominance of a attribute due to diff. of scale in values

$$\text{if } y_{ij} \in [0, 1] \rightarrow y_{ij} = 2 \left[ \frac{x_{ij} - \text{Min}_j}{\text{Max}_j - \text{Min}_j} \right] - 1$$

Metric → Let  $x \in \mathbb{R}^D$  be a set  
 gives measurement

$$d: X \times X \rightarrow [0, \infty)$$

$d$  be a metric if

$$1) d(x, y) = d(y, z); \forall (x, y) \in X$$

$$2) d(x, y) = 0; \text{ if } x=y.$$

$$3) d(x, y) + d(y, z) \geq d(x, z); \forall x, y, z \in X$$

# Minkowski Metric →

$$d_p(x, y) = \left( \sum_{i=1}^D |x_i - y_i|^p \right)^{\frac{1}{p}} \text{ where } p \geq 1$$

$$\textcircled{1} \quad p=1; d_1(x, y) = \sum_{i=1}^D |x_i - y_i| \rightarrow \text{called as}$$

City Block Distance

OR Manhattan Distance

$$\textcircled{2} \quad p=2; d_2(x, y) = \left( \sum_{i=1}^D |x_i - y_i|^2 \right)^{\frac{1}{2}} \rightarrow \text{Euclidean Distance}$$

$$\textcircled{3} \quad p=\infty; d_\infty(x, y) = \max \{|x_i - y_i|\}$$

$$\text{Eg} \rightarrow d_1(P_1, P_2) = |170-160| + |5-5| = 12$$

$$d_2(P_1, P_2) = \sqrt{10^2 + 2^2} = \sqrt{104}$$

$$\star \underline{d_1 > d_2 > d_3 \dots > d_\infty}$$

$$d_\infty(P_1, P_2) = \max(10, 2) = \underline{10}$$

## Weighted distance

# If we know that a attribute (given) is very discriminant for a particular task

Then we can give the higher weight to the higher priority attribute.

$$d_w(P_1, P_2) = \sum_{i=1}^D w_i |x_i - y_i|$$

$$= w_1(10) + w_2(2)$$

$$= 1 \times 10 + 10 \times 2 = \underline{\underline{30}}$$

$$\text{Let } w_1 = 1$$

$$\text{Let } w_2 = \underline{\underline{10}}$$

## Equation of Distance in Matrix form

$$x = \begin{bmatrix} 160 \\ 3 \end{bmatrix}_{D \times 1}$$

$$y = \begin{bmatrix} 170 \\ 5 \end{bmatrix}_{D \times 1}$$

for considered  
eg.

$$\boxed{D=2}$$

$$* x' = [160 \ 3]_{1 \times 2}$$

$$\boxed{* d_2^2(x, y) = (x-y)'(x-y)}$$

$$* x-y = \begin{bmatrix} -10 \\ -2 \end{bmatrix}$$

$$* (x-y)'(x-y) = \begin{bmatrix} -10 & -2 \end{bmatrix}_{1 \times 2} \begin{bmatrix} -10 \\ -2 \end{bmatrix}_{2 \times 1} = [100 + 4] = \boxed{104} = d_2^2(x, y)$$

$$* (x-y)' I_{D \times D} (x-y) = \begin{bmatrix} -10 & -2 \end{bmatrix}_{1 \times 2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}_{2 \times 2} \begin{bmatrix} -10 \\ -2 \end{bmatrix}_{2 \times 1} = [-10 \ -2] \begin{bmatrix} -10 \\ -2 \end{bmatrix} = \boxed{104}$$

In case of any other matrix  
(not I)  $\rightarrow$  weighted euclidean  
distance

02/09/22

$$Q) \text{ Avg} = 67 \text{ in} \quad SD = 3 \text{ in}$$

63 in. to 67 in.

$$\frac{63-67}{3} = -\frac{4}{3} \\ = -1.33 \text{ to } 0$$

$\sigma_{\text{Sample}} = \frac{\sigma}{\sqrt{n}}$   
 ↳ Std. Error

$$0.5000 - 0.0918 \Rightarrow 50 - 9.18 = 40.82\%$$

~~Central Limit Theorem~~

Central Limit Theorem → ~~Standard Deviation~~

Eg → Dice → Avg. of 3 dice  
 More the no. of dice, more close to  
Normal Distribution

Eg → Student

Let there be any sample / distribution.

Mean of the sample Mean will be Mean of Population

Distribution of sample Mean will be ~~mean of the Normal~~  
Distribution

05/08/22  
# Inferential Statistics

\* Not the Sample Size that matters but the quality of Sample is what matters.

Simple Random Sample → all need to be equally likely.  
Eg → Tossing a coin to select.

Random Sample → all need not be equally likely {But same proportion}  
Eg → Population divided in equal classes

Estimation & Std. Error

Q) Given avg age = 37 for population.  
std dev = 15

find avg age of sample of size 200 from population

Sol:

Expected age for sample = 37 {as avg of population}

$$\text{Std. error} = \frac{15}{\sqrt{200}} = \underline{\underline{1.06}}$$

07/09/22

## Testing of Hypothesis

\* One Sample Z-test  $\rightarrow$  Size of Sample  $> 30$

\* Signification Level / Power

09/09/22 t-test (for sample of size  $< 30$ )  $\Rightarrow$  Use t statistics

Eg = population of weights  $\rightarrow$  app. normal.

$\mu$  &  $\sigma$  of population  $\rightarrow$  Unknown

Data  $\rightarrow 31.8, 30.9, 34.2, 32.1, 28.8 \text{ gm}$   $\leftarrow$  small sample

Sample mean = 31.56 gm

\* H<sub>0</sub>:  $\mu = 30$

$\rightarrow H_0: \mu = 30 \text{ & } H_A: \mu > 30$

If H<sub>0</sub> true  $\rightarrow$  sample mean = 30

$$SE \text{ of sample mean} = \frac{\sigma}{\sqrt{n}}$$

$$SE = \frac{\sigma}{\sqrt{n}} \Rightarrow \frac{1.96}{\sqrt{5}} = 0.87$$

$$t \text{ score} = \frac{(31.56 - 30)}{0.87} = 1.79 \Rightarrow P = 7.2\%$$

\* Degree of freedom =  $n - 1$

$$\sigma = 1.96 \frac{\text{sample size}}{n}$$

Cal from  $\left( \frac{1}{3} \right)_{37}$  data

but divide by  $(n-1)$   
as sample is too small

Two Sample Test → Eg - Apple & Egg  
 or test + norm

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad \text{for Z test}$$

Unpaired Test

Paired Test →

diff sets of 100 eggs taken before after  
temp

If  $\text{Var}(s^2)$  of both  
 same  
 same

\* for t-test

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

$$df = n_1 + n_2 - 2$$

If var are diff

What if More than Two sample? {Anova}

### Method-Between

$$\text{Sum of S. B/w (SSB)} = n \sum_{i=1}^a (\bar{y}_i - \bar{y})^2$$

$$D.o.F = a-1$$

$$\text{Mean Sq. Between (MSB)} = SSB / D.o.F$$

### Method - Within

$$SSE_{\text{Error}} = \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2$$

$$D.o.F = N-a$$

$$MSE = \frac{SSE}{D.o.F}$$

$$F \text{ score} = MSB / MSE$$

↓  
every individual value

$\bar{y}$  (Global Mean)

#

	$x_{11}$	$x_{12}$	...	$x_{1n}$	mean $\bar{y}_1$
P	$x_{21}$	$x_{22}$	...	$x_{2n}$	$\bar{y}_2$
B	$x_{31}$	$x_{32}$	...	$x_{3n}$	$\bar{y}_3$
C	$x_{41}$	$x_{42}$	...	$x_{4n}$	$\bar{y}_4$
D					

If  $MSB$ ,  
 $MSE$  &  $MST$   
are same  
↓  
Null. Hypo.

Else  
is not Null.  
hypo  
 $MSB > MSE \Rightarrow$   
MSE

$$\text{Method-Total} \quad SST = \sum_{i=1}^a \sum_{j=1}^n (y_{ij} - \bar{y})^2$$

$$D.o.F : N-1$$

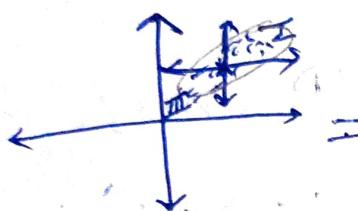
$$MS1 = SST / D.o.F$$

12/09/22

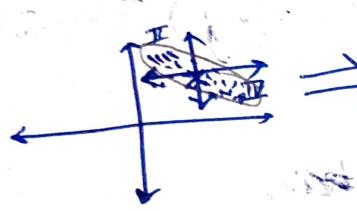
$$\text{Mean} \rightarrow \mu_x = \frac{1}{n} \sum_{i=1}^n x_i = \sum_{i=1}^n p_i x_i$$

$$\text{Variance} \rightarrow \sigma^2(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(x_i - \mu_x)$$

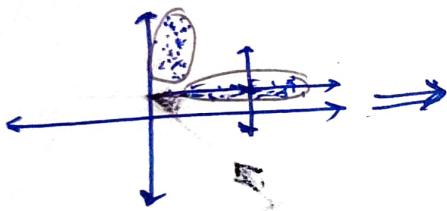
\* Covariance  $\rightarrow \text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$



$\Rightarrow \boxed{\text{Cov}(x, y) > 0}$  if most of Data pt lies in I & III quad of shifted origin



$\Rightarrow \boxed{\text{Cov}(x, y) < 0}$  if most of Data pt lies in II & IV quad of shifted origin



$\Rightarrow \boxed{\text{Cov}(x, y) = 0}$   $x \& y$  are unrelated

# Correlation  $\rightarrow$  Correlation of two variables ( $x, y$ )

$$f(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}}$$

$$-1 \leq f(x, y) \leq 1$$

# Measuring dist b/w a data point & a cluster/class of  $k$

i) Measuring dist b/w data pt. & the mean of class.



$d_k(x, \mu_k) \rightarrow$  if dict more - not in class else Belong to class

Problem  $\rightarrow$  Distribution/Spreadness of classes is ignored

ii) So, we need another measure to consider spread



$$\Sigma = \begin{matrix} f_1 & f_2 & \dots & f_D \end{matrix}$$

Dispersion Matrix

Variance-Covariance Covariance Matrix

$f_1$	$\text{Cov}(f_1, f_1)$	$\text{Cov}(f_1, f_2)$	$\dots$	$\text{Cov}(f_1, f_D)$
$f_2$	$\text{Cov}(f_2, f_1)$	$\text{Cov}(f_2, f_2)$	$\dots$	$\text{Cov}(f_2, f_D)$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$f_D$	$\text{Cov}(f_D, f_1)$	$\text{Cov}(f_D, f_2)$	$\dots$	$\text{Cov}(f_D, f_D)$

Each cell  $\sigma_{ij}$  contains covariance of  $f_i, f_j$

$D \times D$  → no. of features

$$* \text{Cov}(f_i, f_i) = \underline{\text{Var}}(f_i) * \text{Cov}(f_i, f_j) = \text{Cov}(f_j, f_i)$$

\*  $\Sigma$  is real, symmetric, positive definite Matrix.

all cells contains  
real values only

symm wrt diagonal  
as  $\text{cov}(f_i, f_j) = \text{cov}(f_j, f_i)$

A matrix is positive definite  
iff  $x'Ax > 0$  for  
any vector  $x$

$$\rightarrow |\Sigma| > 0$$

$\rightarrow \Sigma^{-1}$  also be ~~also~~ positive definite

$\rightarrow$  All the eigenvalues are positive

$Ax = \lambda x$  Eigen equation

$$\boxed{|A - \lambda I| = 0}$$

$$\text{let } A = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix}$$

$$x = \begin{bmatrix} a \\ b \end{bmatrix}$$

$$x'Ax = [a \ b] \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

$$= [ab] \begin{bmatrix} 2a - b \\ -a + 3b \end{bmatrix}$$

$$= 2a^2 - 2ab + 3b^2$$

$$= 2 \left[ \frac{a-b}{2} \right]^2 + \frac{5}{2} b^2 > 0$$

\* Mahalanobis Dist b/w class & pt.

$$d_M^2(c_1, x) = (x - \mu_1)' \Sigma^{-1} (x - \mu_1)$$

\* Mahalanobis dist b/w two classes

$$d_M^2(c_1, c_2) = (\mu_2 - \mu_1)' \Sigma^{-1} (\mu_2 - \mu_1)$$

$$\boxed{\Sigma = \frac{\Sigma_1 + \Sigma_2}{2}}$$

14/09/22

Find  $\cos \theta$  b/w two vectors.

$$X_1 = \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1D} \end{bmatrix} \quad X_2 = \begin{bmatrix} x_{21} \\ x_{22} \\ \vdots \\ x_{2D} \end{bmatrix}$$

$$X_1 \cdot X_2 = |X_1| \cdot |X_2| \cos \theta = X_1^T X_2$$

⇒

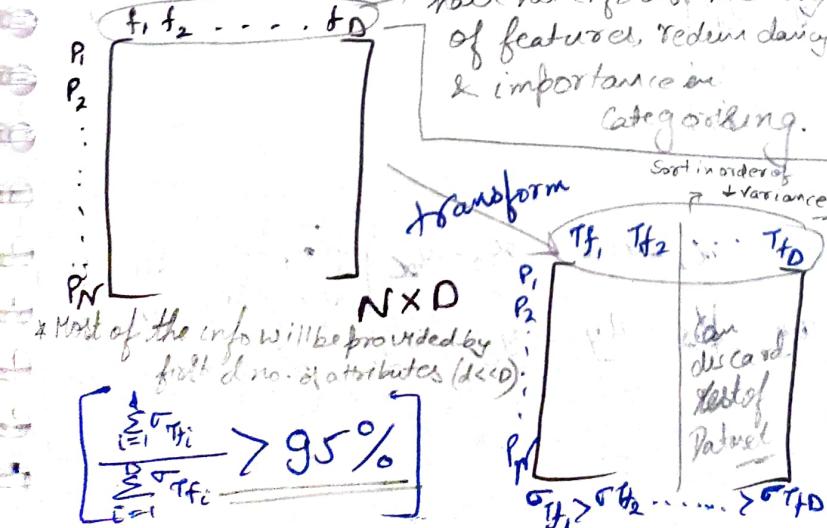
$$\cos \theta = \frac{X_1^T X_2}{\sqrt{\sum_{i=1}^D x_{1i}^2} \sqrt{\sum_{i=1}^D x_{2i}^2}}$$

$$\begin{aligned} A = & a_i + b_j = \begin{bmatrix} a \\ b \end{bmatrix} \rightarrow x_1 \\ B = & c_i + d_j = \begin{bmatrix} c \\ d \end{bmatrix} \rightarrow x_2 \\ A \cdot B = & ac + bd. \\ X_1^T X_2 = & \begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} c \\ d \end{bmatrix} \\ = & [ac + bd] \end{aligned}$$

\* Attribute with high variance is better to be considered for categorisation.

\* Choosing all the features ~~for~~ for categorization will just increase the complexity.

So, need to select some very imp. features. But How??



M-1 → Transform the dataset (of which we don't know the feature meaning) into the dataset whose ~~features~~ features are ordered for order of decreasing variance of values of that features.

1. Let  $X$  be the mean subtracted data i.e. matrix of  $N \times D$ .
- $N \rightarrow$  no. of pattern.
- $D \rightarrow$  no. of attributes.

## 2. Calculate Dispersion matrix

$$\Sigma = \frac{1}{N} (X^T X)$$

Eg  $\rightarrow$   $P_1 \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ x_3 & y_3 \end{bmatrix}$

$P_2$

$P_3$

$x_1 - \bar{x}_1 \quad y_1 - \bar{y}_1$

$x_2 - \bar{x}_1 \quad y_2 - \bar{y}_1$

$x_3 - \bar{x}_1 \quad y_3 - \bar{y}_1$

$\underbrace{\quad}_{N \times D}$   
 $[3 \times 2]$

$$\bar{x}_x = \frac{x_1 + x_2 + x_3}{3} \quad \& \quad \bar{y}_y = \frac{y_1 + y_2 + y_3}{3}$$

~~Context~~

$$\frac{1}{3} \left[ (x_1 - \bar{x}_x)(y_1 - \bar{y}_y) + (x_2 - \bar{x}_x)(y_2 - \bar{y}_y) + (x_3 - \bar{x}_x)(y_3 - \bar{y}_y) \right]$$

$$\star \Sigma = \begin{bmatrix} \text{Var}(x) & \text{Cov}(x,y) \\ \text{Cov}(x,y) & \text{Var}(y) \end{bmatrix}$$

$$\text{Var}(x) = \frac{1}{3} \left[ (x_1 - \bar{x}_x)^2 + (x_2 - \bar{x}_x)^2 + (x_3 - \bar{x}_x)^2 \right]$$

$$\star X = \begin{bmatrix} x_1 - \bar{x}_x & y_1 - \bar{y}_y \\ x_2 - \bar{x}_x & y_2 - \bar{y}_y \\ x_3 - \bar{x}_x & y_3 - \bar{y}_y \end{bmatrix}$$

Product =  $N \times \text{Var}(x)$

$$X^T X = \begin{bmatrix} x_1 - \bar{x}_x & x_2 - \bar{x}_x & x_3 - \bar{x}_x \\ y_1 - \bar{y}_y & y_2 - \bar{y}_y & y_3 - \bar{y}_y \end{bmatrix}$$

$$\begin{bmatrix} x_1 - \bar{x}_x & y_1 - \bar{y}_y \\ x_2 - \bar{x}_x & y_2 - \bar{y}_y \\ x_3 - \bar{x}_x & y_3 - \bar{y}_y \end{bmatrix}$$

\*  $\Sigma$  is a  $D \times D$  matrix.

We need to find eigenvalues & eigen vectors.

D' eigen values  
will be there

$$(\lambda_1, \lambda_2, \dots, \lambda_D) > 0$$

$\rightarrow \lambda_1 > \lambda_2 > \dots > \lambda_D$  be the ordered eigen values.  
 $v_1, v_2, v_D \rightarrow$  eigen vectors

$$\lambda_0, \lambda_i > 0$$

$$\& v_i = \begin{bmatrix} v_{i1} \\ v_{i2} \\ \vdots \\ v_{iD} \end{bmatrix}_{D \times 1}$$

$\rightarrow \text{Trace}(\Sigma) \rightarrow$  sum of variances w.r.t. all features

$$= \sum_{i=1}^D \text{var}(x_i) = \sum_{i=1}^D \lambda_i$$

$$W = \begin{bmatrix} v_1, v_2, \dots, v_D \\ v_1, v_2, \dots, v_D \\ \vdots & \vdots \\ v_D, v_2, \dots, v_D \end{bmatrix}_{d \times D}$$

Transformation Matrix

\* Can arrange in form of rows also.

Select largest d eigen values  
 $(\lambda_1, \dots, \lambda_d)$

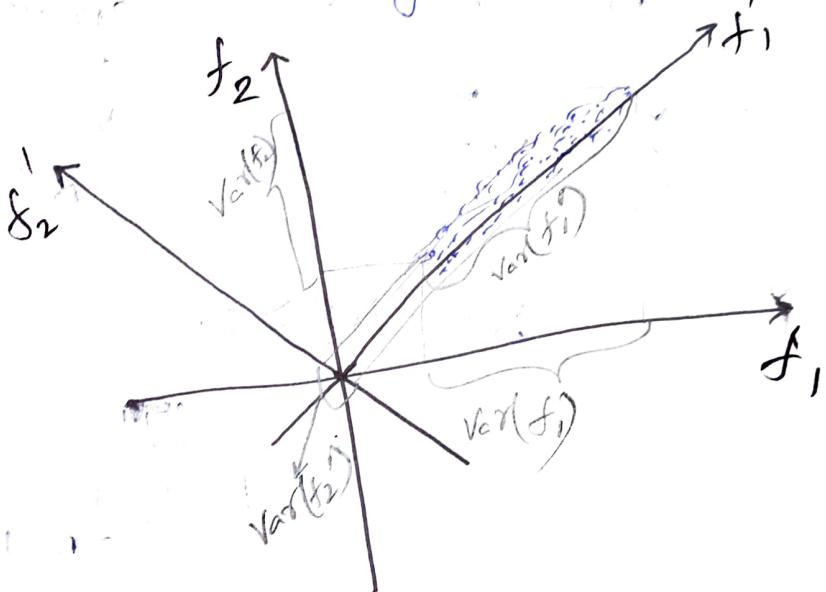
$\sqrt{\lambda_i}$   
 $(v_1, \dots, v_d)$   
→ first d eigen vectors from  $W$ .

19/09/22

\*  $\frac{\sum_{i=1}^D \lambda_i}{\sum_{i=1}^D \lambda_i} > 98\% \rightarrow$  choose d & remove all others from  $f_d$  to  $f_D$

→ After transformation → variance of a feature is much larger than the others.

Eg →



$$\rightarrow Y = W X$$

$D \times N$        $D \times D$        $D \times N$   
 ↓                  ↓                  ↓  
 initial      transformed      Dataset  
 Transfor-      matrix      Matrix  
 mation      matrix

This Concept is  
Principal Component  
Analysis

$$\text{Eg} \rightarrow X = \begin{bmatrix} 2 & 1 \\ 3 & 5 \\ 4 & 3 \\ 5 & 6 \\ 6 & 7 \\ 7 & 8 \end{bmatrix} \rightarrow X' = \begin{bmatrix} -2.5 & -4 \\ -1.5 & 0 \\ -0.5 & -2 \\ 0.5 & 1 \\ 1.5 & 2 \\ 2.5 & 3 \end{bmatrix}_{6 \times 2}$$

$$M_1 = 4.5 \quad M_2 = 5$$

$$\Sigma = \frac{1}{N} (X'^T X')$$

here

$$N = 6$$

~~$$10 + 0 + 1 + 0.5 + 3 + 7.5 = 24$$~~

$$\Sigma = \frac{1}{N} \begin{bmatrix} 6.25 + 2.25 + 0.25 + 0.25 + 2.25 + 6.25 & 22 \\ 22 & 16 + 0 + 4 + 1 + 4 + 9 \end{bmatrix}$$

$$\Sigma = \frac{1}{N} \begin{bmatrix} 17.5 & 22 \\ 22 & 34 \end{bmatrix} = \frac{1}{6} \begin{bmatrix} 175 & 22 \\ 22 & 34 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix}$$

$$|A - \lambda I| = 0$$

~~Is it the final DS~~  
How can we be  
sure of no  
redundancy.

$$\Rightarrow \begin{vmatrix} 2.92-\lambda & 3.67 \\ 3.67 & 5.67-\lambda \end{vmatrix} = 0$$

$$\Rightarrow (2.92-\lambda)(5.67-\lambda) - (3.67)^2 = 0$$

$$\Rightarrow 2.92 \times 5.67 - 5.67\lambda - 2.92\lambda + \lambda^2 - (3.67)^2 = 0$$

$\Rightarrow$

$$\Rightarrow \boxed{\lambda_1 = 8.22} \quad \& \quad \boxed{\lambda_2 = 0.37}$$

$$v_1 = \begin{bmatrix} 2.55 \\ 3.67 \end{bmatrix}$$

$$\boxed{* \frac{\lambda_1}{\lambda_1 + \lambda_2}}$$

As  $\lambda_1$  is higher  $\Rightarrow W = \begin{bmatrix} v_1 \\ 2.55 \\ 3.67 \end{bmatrix}_{2 \times 1}$

'final DS'  
but lost the  
characteristics/meaning  
of all the  
attributes  
of Initial DS

$$Y = \begin{bmatrix} 2.55 & 3.67 \end{bmatrix} \begin{bmatrix} 2 & 3 & 4 & 5 & 6 & 7 \\ 1 & 5 & 3 & 6 & 7 & 8 \end{bmatrix}_{2 \times 6}$$

$$Y = \begin{bmatrix} 8.77 & 26 & 21.21 & 34.47 & 40.92 & 41.21 \end{bmatrix}$$

\* If we know the meaning of the features  $\Rightarrow$  we know their importance, so we do not need to perform PCA.

So, needs to perform PCA, to know the importance which we do based on Variance

Use only when no. info abt meaning/ characteristics of features

## # Measurement of compactness of classes.

$\rightarrow$  Within classes Distance

Intra class Distance

\* Should be small



C<sub>1</sub>



C<sub>2</sub>

better  
as within class Distance (dist of every pt. from mean) is less

## # Measurement of Better Dataset

$\rightarrow$  Between class Distance

Inter class Distance

\* Should be large

$\rightarrow$  proper distributed  
categorised

$\rightarrow$  Less Error

Dataset-2

Dataset-1



C<sub>1</sub> C<sub>2</sub>

21/09/22 Learning → Why? → to get success & make less mistakes  
↳ getting some knowledge & use it to increase your performance

- Supervised → prior info. given → Parents teach their child.  
↳ training info / data given.
- Unsupervised → learn from own experience → no prior info.  
↳ need to learn directly.

\* Training Data → Along with patterns (of diff features), class labels are also given.

\* Classification Task → Supervised learning using training data

\* Clustering Task → unsupervised learning based on similarities & dissimilarities of patterns of features.

→ need to find Natural grouping

look on within & b/w class dist

→ generate and train a classification model by using  
finding the  
Training Data (Pattern + Classes)

→ for an unknown pattern, assign it to a predefined  
classes based on the classifier.

# Various types of Classification Models

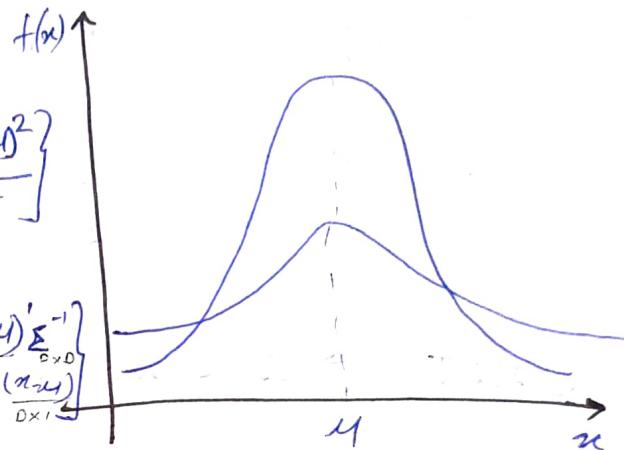
## ① Bayes' Decision Rule / Bayes' Classifier.

prior prob.,  $P_i / P(w_i)$  → Prob. of any new data to be a part of  $C_i$

\* In Gaussian Dist./Normal Dist

Exponential Rule  
68% - 95% - 99.7%

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}\right\}$$



→ Let there be  $M$  classes ( $M \geq 2$ ). The prior probability be  $P_1, P_2, \dots, P_M$ , where  $P_i \in [0, 1]$  &  $\sum_{i=1}^M P_i = 1$ .

Let  $p_1(x), p_2(x), \dots, p_M(x)$  be the cond<sup>nal</sup> prob. density func. of  $M$  classes.

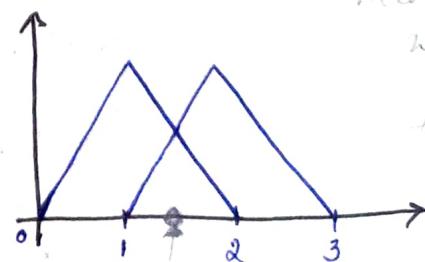
Now, any unknown pattern  $x_0$  to be in class  $i$  if

$$P_i P_i(x_0) \geq P_j P_j(x_0); \forall i, j \text{ & } i \neq j.$$

Example →  $M=2$ ; Prior Prob:  $P, 1-P$ .

$$P_1(x) = \begin{cases} x & ; 0 \leq x \leq 1 \\ 2-x & ; 1 \leq x \leq 2 \\ 0 & ; \text{otherwise} \end{cases}$$

$$P_2(x) = \begin{cases} x-1 & ; 1 \leq x \leq 2 \\ 3-x & ; 2 \leq x \leq 3 \\ 0 & ; \text{otherwise} \end{cases}$$



need to find x which divides the upcoming data into two diff. Classes

Case-I  $\rightarrow 0 \leq x \leq 1$

$$P_1 = P \quad P_2 = 1 - P$$

$$P_1(x) = x \quad P_2(x) = 0$$

$$\rightarrow P_1 P_1(x) \geq P_2 P_2(x) \stackrel{>0}{\Rightarrow} x \text{ be in class } \underline{1}$$

$$\Rightarrow \boxed{Px \geq 0} \xrightarrow{\text{as } P \text{ is const.} (>0)} \Rightarrow \boxed{x \geq 0}$$

Case-II  $\rightarrow 2 \leq x \leq 3$

$$P_1 = P \quad P_2 = 1 - P$$

$$P_1(x) = 0 \quad P_2(x) = 3 - x$$

$$\rightarrow P_2 P_2(x) \geq P_1 P_1(x) \stackrel{>0}{\Rightarrow} x \text{ be in class } \underline{2}$$

$$\Rightarrow (1-P)(3-x) \geq 0 \Rightarrow \boxed{Px - 3P - x + 3 \geq 0} \xrightarrow{\substack{P \in (0,1) \\ \text{as } P \text{ is const}}} \Rightarrow \boxed{x \leq 3}$$

Case-III  $\rightarrow 1 \leq x \leq 2$

$$P_1 = P \quad P_2 = 1 - P$$

$$P_1(x) = 2 - x \quad P_2(x) = x - 1$$

$$\text{for } x \text{ be in class 1} \rightarrow P_1 P_1(x) \geq P_2 P_2(x)$$

$$\Rightarrow P(2-x) \geq (1-P)(x-1)$$

$$\Rightarrow 2P - Px \geq x - 1 - Px + P$$

$$\Rightarrow \boxed{x \leq P+1}$$

so  $\rightarrow$  if  $\boxed{x \leq P+1} \rightarrow$  in class 1

if  $\boxed{x \geq P+1}$

$x$  in class 2

23/09/22

$$\text{Error, } E = P_1 \int_{1+P}^2 p_1(x) dx + P_2 + \int_1^{1+P} p_2(x) dx \text{ for Case } - \underline{\underline{III}}$$

class 2 element  
in the region of class II.

In general,

$$E = \sum_{i=1}^c P_i \int_{\Omega_i^c} p_i(x) dx$$

$\hookrightarrow$  no. of classes

$$E = \sum_{i=1}^c P_i \int_{\Omega_i^c} p_i(x) dx$$

$\hookrightarrow$  Complement of  
class- $i$  region

$$\Omega = [0 \text{ to } 3]$$

$$\Omega_1 = [0 \text{ to } 1+P]$$

$$\Omega_1^c = [1+P \text{ to } 3]$$

$$\Omega_2 = [1+P \text{ to } 3]$$

$$\Omega_2^c = [0 \text{ to } 1+P]$$

So, overall,

$$E = \sum_{i=1}^c P_i \int_{\Omega_i^c} p_i(x) dx$$

$$E = P_1 \int_{1+P}^3 p_1(x) dx + P_2 \int_0^{1+P} p_2(x) dx$$

$$E = P_1 \int_{1+P}^2 p_1(x) dx + P_2 \int_2^3 p_2(x) dx + P_1 \int_1^{1+P} p_2(x) dx + P_2 \int_1^0 p_2(x) dx$$

$$E = P_1 \int_{1+P}^2 p_1(x) dx + P_2 \int_1^3 p_2(x) dx$$

$$E = P \int_{1+P}^2 (2-x) dx + (1-P) \int_1^3 (x-1) dx$$

$$E = P \left[ \frac{1+P}{2} \left( 2x - \frac{x^2}{2} \right) \right]_{1+P}^2 + (1-P) \left[ \frac{1}{2} \left( x^2 - x \right) \right]_1^{1+P}$$

$$= P \left[ \frac{4-2-2(1+P)+1+P^2}{2} \right] - \frac{(1+P)^2}{2} + 1+P$$

$$F = P(1-P) + \frac{(P^2+2P+1)}{2} - P + 1 = \frac{-2P^2+2P+P^2+2P+1-2P-2}{2} + \frac{(1+P)^2}{2} - 1-P$$

$$E = \frac{-P^2+2P-1}{2}$$

$$= -(P-1)^2$$

So, for bayes' Decision Rule

$$E_B = \frac{P(1-P)}{2}$$

$$\Omega_1 = [0 \text{ to } 1+P]$$

$$\Omega_2 = [1+P \text{ to } 3]$$

Let there be another class after

$$\Omega_1 = [0 \text{ to } 1.5]$$

Here:  ~~$E = 0.125$~~

$$\Omega_2 = [1.5 \text{ to } 3]$$

$$E = P \left[ 2x - \frac{x^2}{2} \right]_{1.5}^2 + (1-P) \left[ \frac{x^2}{2} - x \right]^{1.5}$$

$$E = P \left[ 4 - 2 - 3 + \frac{(1.5)^2}{2} \right] + (1-P) \left[ \frac{(1.5)^2}{2} - 1.5 - \frac{1}{2} + 1 \right]$$

$$E = P \left[ -1 + \frac{(1.5)^2}{2} - \frac{(1.5)^2}{2} + 1 \right] + \left[ \frac{(1.5)^2}{2} - 1 \right]$$

$$E = \frac{2.25 - 2}{2} = \frac{25}{200} = \frac{1}{8}$$

$$E_{1.5} = 0.125$$

$\rightarrow$  Is  $E_B \leq E_{1.5}$ ?  $\rightarrow$  Yes

$$E_B \leq E_{1.5}$$

$E$  of any classifier

$E_B$  is less than  $E_{1.5}$  if values of  $P \neq 0.5$

$E_B$  is equal to  $E_{1.5}$  for  $P = 0.5$

NOTE → Minimum Error classifier / Classifier with minimum error prob. is Bayes' classifier

So, Bayes' Classifier is the best classifier.

\* What is the need of any other classifier then

↳ Real Time Data's distribution cannot be found / Real Time Data does not follow any fixed Distribution. { In Most of the Cases }

↳ Estimating  $p_i(x)$ , conditional Probability Density function, is very difficult & computationally very complex  
for the Real Time Datasets. —