
Prepared and Written By
RAZEEN AHMED

CSE-437
[DATA SCIENCE]

HANDWRITTEN NOTE

Computer Science and Engineering
BRAC University

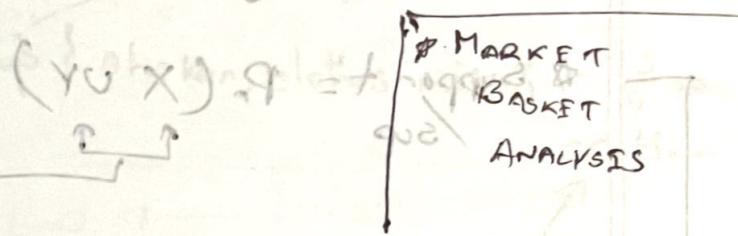
Github: github.com/razeen
LinkedIn: linkedin.com/in/razeenahmed

MINING ASSOCIATION RULES

* Degree of association can be different.

RULES is

→ Assume all data are categorical



The model:

$$(i) I = \{ i_1, i_2, i_3, \dots, i_m \}$$

$$I = \{ A, B, X, Y \}$$

$$(ii) t \in I$$

is a subset of

$$t_1 = \{ A, X, Y \}$$

$$(iii) T = \{ t_1, t_2, \dots, t_n \}$$

$$t_2 = \{ B, X, Y \}$$

$$t_3 = \{ X, Y \}$$

The model rules:

$$I_1 = \{ A, B, X \}$$

$$T = \{ t_1, t_2, t_3 \}$$

$$A, B \rightarrow X$$

[If someone buys A & B he will purchase X]

$$A, B \rightarrow X, B$$

$\boxed{\text{Violation of rule!}}$
not possible.

Application of Association rule

Example → efficient keyboard

⇒ Put the keys

which have higher association together

→ Items → keys

transaction ⇒ most frequent words

transaction ⇒ All the database dictionary

2nd year COMP2022A ALGORITHM

Measure strength of association: To decide if

$$\text{Support} = \Pr(X \cup Y)$$

[There will be some amount of threshold]

Probability of finding X & Y together.

Confidence

Conf =

$$\Pr(Y|X)$$

of transaction

If X is there then what is the probability that Y will also be there.

(ii)

$$\text{Support} = \frac{(X \cup Y). \text{Count}}{n} = \frac{\text{no. of times they were together}}{\text{no. of total transaction}}$$

$$\{X, Y, A\} = ?$$

Confidence =

$$\frac{(X \cup Y). \text{Count}}{X. \text{Count}}$$

$$X \rightarrow Y, A$$

or $Y \wedge A$ and so on

X preceding Y, A

Goal → Find all association rules that support

"support threshold" and "confidence threshold."

EXAMPLE: Transport is forwarded via Optimizing routes.

$$\text{minsup} \text{ minsup} = 30\%$$

$$\text{min conf} = 80\%$$

$$\text{sup} \{ \text{chicken, clothes, milk} \} = \frac{3}{7} = 0.42857$$

= 0.42857

= 0.42857

$$\text{Conf} = \frac{(X \cup Y) \cdot \text{Count}}{X \cdot \text{Count}}$$

$$= \frac{3}{3}$$

$$= 1$$

$$= 100\% \leftarrow \text{This can be considered as association rule as it surpasses minsup threshold}$$

EXAMPLE: This can be considered as association rule as it surpasses minconf threshold

(S, X), (Y, X) \leftarrow threshold

$$\text{Confidence} = \frac{\text{Sup}(X \cup Y)}{\text{Sup}(X)}$$

X \rightarrow Y

Yielding rule

$$\left\{ \{2, N, S, P\}, \{N, S, P\} \right\} = P$$

Stable and durable

$$\{P, S, P\}, \{P, S, P\}, \{N, S, P\}, \{S, P, P\}$$

APRIORI ALGORITHM.

apriori property :- Any subsets of a frequent itemset is also a frequent itemset.

\rightarrow 08 = quantum frequent

\rightarrow 108 = frequent item

Feasible itemset \rightarrow candidate itemset.

$$f_{m+1}(Y \cup X) = f_m$$

$$f_{m+1}(X)$$

Problem

It needs to scan the big dataset again & again.

EXAMPLE

ALGORITHM'S 1st step:-

2. Scan T \rightarrow

frequent itemsets,

C_3 :

(Join operation

$$\rightarrow (X, Y), (X, Z)$$

if this condition holds true then

$(Y \cup X)$ can apply join operation

$$Y \in F_3 \rightarrow \{1, 2, 3\}, \{1, 3, 5\}, \{2, 3, 5\}$$

After join:

$$\{1, 3, 4\}, \{1, 3, 5\}, \{1, 4, 5\}, \{3, 4, 5\}$$

$$C_4 = \left\{ \left\{ \underbrace{\{1, 2, 3, 4\}}_{\text{Subsets from candidate set}}, \underbrace{\{1, 2, 3, 5\}}_{\text{possible subsets}} \right\} \right\}$$

$$\{1, 2, 3\}, \{1, 2, 4\}, \{1, 3, 4\}, \{2, 3, 4\}$$

Now search this subset in F_3 .

$$\{1, 4, 5\}$$

\downarrow no present

If one of the subset is missing in F_3 then we can discard the

Candidate set is or prone to candidate set.
because it violates the apriori algorithm.

ALGORITHM'S 2ND STEP:- GENERATING RULES

for F_3
Association rules for $\{2, 3, 5\}$

conf times = 80.00000000000001% (i)

$$\text{(i) } 2 \rightarrow 3, 5 \quad \text{conf} = \frac{2}{8} = 25\% \quad \text{(ii)}$$

$$\text{(ii) } 3 \rightarrow 2, 5 \quad \text{conf} = \frac{3}{8} = 37.5\% \quad \text{(iii)}$$

$$\text{(iii) } 5 \rightarrow 2, 3 \quad \text{conf} = \frac{5}{8} = 62.5\%$$

$$\text{(iv) } 2, 3 \rightarrow 5$$

$$\text{(v) } 2, 5 \rightarrow 3$$

$$\text{(vi) } 3, 5 \rightarrow 2$$

(conf values are wrong)

Best option after
for F_2 :-

(i) $\{1, 3\}$ (ii) $\{2, 3\}$ (iii) $\{2, 5\}$, (iv) $\{3, 5\}$ (v) existing mining

$$\begin{array}{c|c|c|c|c} 1 \rightarrow 3 & 2 \rightarrow 3 & 2 \rightarrow 5 & 3 \rightarrow 5 & 5 \rightarrow 3 \\ 3 \rightarrow 1 & 3 \rightarrow 2 & 5 \rightarrow 2 & 5 \rightarrow 3 & \end{array} \quad \text{find conf.}$$

LIFT THRESHOLD :-

$$\text{* Lift} = \frac{\text{sup}(X \cup Y)}{\text{sup}(X) * \text{sup}(Y)} = \frac{n * (X \cup Y). \text{count}}{X. \text{count} * Y. \text{count}}$$

$$\text{* Lift} = \frac{P(X \cup Y)}{P(X) P(Y)} = \frac{P(Y|X) P(X)}{P(X) P(Y)} = A$$

if $X \& Y$ has some dependency

$$\frac{P(Y) \& P(X)}{P(X) P(Y)} = 1 \quad \text{if } X \& Y \text{ are independent.}$$

PRINCIPAL COMPONENT ANALYSIS.

20102125

STATISTICS

Basics

(i) MEAN $= \frac{108}{29} = 3.7$

Population $\rightarrow \eta$

(iv) STANDARDS FOR USE

(iii) VARIANCE

(iv) CO-VARIANCE

$$\hookrightarrow \beta = \frac{\sum_{i=0}^n (\cancel{x_i - \bar{x}})(\gamma_i - \bar{y})}{n-2}$$

$\Rightarrow \leftarrow$ $E_8(i)$

$\Sigma \rightarrow \Sigma, \zeta$ (v)

~~Q2 feature~~ (iv)

(now one feature v

with another feature

Covariance matrix properties

$\{ \varepsilon, \delta \}_{\text{B}} \} \{ \varepsilon, \delta \}_{\text{d}}$

$\text{Cov}(x, y) \rightarrow$ Covariance of x & y | $\text{Cov}(x, y) = \text{Cov}(y, x)$

$\text{cod}(y, y) \rightarrow "$ " y

MATHEMATICS Basics

(i) Eigen vector & Eigen value.

$$A = \begin{bmatrix} 7 & 3 \\ 3 & -1 \end{bmatrix} \xrightarrow{(x)^9 (y)^9} \frac{(x)^9 (y)^9}{(y)^9 (x)^9} = \begin{bmatrix} 1 & (y/x)^9 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

$$\text{tanhgabn} \rightarrow D = \frac{(x)^9 (y)^9}{(y)^9 (x)^9}$$

$$\text{Step 3: } \begin{bmatrix} d & a & b \\ e & f & c \\ g & h & i \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \lambda = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = \begin{bmatrix} s\sqrt{2} & 0 \\ 0 & s\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}$$

$$\text{Step 2} \quad A - \lambda I = \begin{bmatrix} 7-\lambda & 3 \\ 3 & -1-\lambda \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 7-\lambda-\lambda^2 & 3 \\ 3 & -1-\lambda-\lambda^2 \end{bmatrix} = \begin{bmatrix} 7-\lambda-\lambda^2 & 3 \\ 3 & -1-\lambda-\lambda^2 \end{bmatrix}$$

Step-3

$$|A - \lambda I| = \begin{vmatrix} 7-\lambda & 3 \\ 3 & -1-\lambda \end{vmatrix}$$

$$= (7-\lambda)(-1-\lambda) - 9$$

$$= -7 - 7\lambda + \lambda^2 + 3\lambda - 9$$

$$= -7 - 4\lambda + \lambda^2 - 9$$

$$\Sigma = 5\lambda^2 - 26\lambda - 56$$

$$\Rightarrow \lambda^2 - 6\lambda - 16 = 0$$

$$\therefore \lambda = 8 \quad \lambda = 2$$

Step-4

Substituting $\lambda = 8$ in Σ in Σ (i)

$$B = \begin{bmatrix} -1 & 3 \\ 3 & -9 \end{bmatrix}$$

Property of eigenvector

$$B \vec{v} = \vec{0}$$

$$\vec{v} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\Sigma = \vec{v}^T B \vec{v}$$

$$\Sigma = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & 3 \\ 3 & -9 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\Sigma = 1(-1) + 1(3) + 1(3) + 1(-9) = -1 + 3 + 3 - 9 = -6$$

$$\Sigma = 6$$

E collar form transformation

$$B \vec{v} = \vec{0}$$

$$\begin{bmatrix} -1 & 3 \\ 3 & -9 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} a & b & e \\ d & e & f \\ g & h & i \end{bmatrix}$$

Eqpt 2

B transform into E collar form \rightarrow

$$R_2 = 3 * R_1 + R_2$$

$$B \Rightarrow \begin{bmatrix} -1 & 3 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Eqpt 2

Now, N-95f3

\Rightarrow primitive

$$-1v_1 + 3v_2 = 0$$

$$\begin{vmatrix} 2 & R-p \\ R-1 & 2 \end{vmatrix} = (2R - A)$$

$$R = 2$$

$$R = 0 * v_1 + 0 * v_2 + 0 * p - p = 0$$

$$0 = 0 + 0 + 0 - p = -p$$

From (1) we get,

$$\Rightarrow -1v_1 + 3v_2 = 0$$

$$\Rightarrow -v_1 = -3v_2$$

$$\therefore v_1 = 3v_2$$

$$\begin{bmatrix} 6 \\ 2 \end{bmatrix} = 2 \begin{bmatrix} 3 \\ 1 \end{bmatrix} \quad \text{if } v_2 = 2 \\ \quad \quad \quad v_1 = 6$$

Why?

$$\begin{cases} 2x + R(3-p) = 0 \\ v_2 = 2 \end{cases}$$

Eigen vector \rightarrow

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

Covariance matrix will be a square matrix

Test \rightarrow

$$A\vec{v} = \lambda \vec{v}$$

Optidimur $\lambda \in \mathbb{R}$ Random \vec{v} eigen vector
vector \Rightarrow eigen vector \times eigen value.

$$\Rightarrow \begin{bmatrix} 7 & 3 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \end{bmatrix} = \begin{bmatrix} 8 \begin{bmatrix} 3 \\ 1 \end{bmatrix} \\ 8 \begin{bmatrix} 3 \\ 1 \end{bmatrix} \end{bmatrix} = \text{solution?}$$

$$\Rightarrow \begin{bmatrix} 7 \times 3 & 3 \times 1 \\ 3 \times 1 & 1 \times 1 \end{bmatrix} = \begin{bmatrix} 24 \\ 8 \end{bmatrix} \text{ (valid)}$$

$$[V \quad \vec{v}] = VF$$

PCA \rightarrow (unsupervised method)

$$\text{Var. } \approx \text{Var. } 0 = \frac{N \bar{x}^2}{N-1} = \frac{\bar{x}^2}{\bar{x}^2 + \bar{x}^2} = 0.25$$

Feature vector, $FV^T = (\vec{v}_1, \vec{v}_2)$

In a way that \vec{v}_1 has largest eigenvalue and \vec{v}_2 has least eigenvalue.

eigen vector = $\begin{bmatrix} 0.99 & 0.1 \\ 0 & 0.99 \end{bmatrix}$

Final data, FD = $FV^T * \text{AdjData}^T$

$$PC = FD^T$$

$$\begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}$$

PC₁ PC₂

shape is set by the variance matrix

How can we select 99% variability?

$$\sqrt{\lambda} = \sqrt{A}$$

eigenvalues = $\begin{bmatrix} 0.049083 & [1] \\ 1.28403 & [8] \end{bmatrix} = \begin{bmatrix} \lambda_1 = 0.049 & [\lambda_1] \\ \lambda_2 = 1.284 & [\lambda_2] \end{bmatrix}$

(b1, v1) $\begin{bmatrix} 0.049 \\ 1.284 \end{bmatrix} = \begin{bmatrix} \Sigma \lambda_1 \\ \Sigma \lambda_1 + \lambda_2 \end{bmatrix} = \begin{bmatrix} 0.049 \\ 1.333 \end{bmatrix}$

$$FV = \begin{bmatrix} V_2 & V_1 \end{bmatrix}$$

↑ PC₁ ↑ PC₂ (bottom being required) $\rightarrow 99\%$

variability holds by $PC_1 = \frac{\lambda_2}{\lambda_1 + \lambda_2} = \frac{1.284}{1.333} = 0.963 \approx 96.3\%$

96.1.

$$PC_2 = \frac{3.7\%}{\text{approx}} \quad (\bar{v}_1, \bar{v}_2) = \text{top } PC_1 \text{ entries} \rightarrow 96.3\%$$

If it is guided we can drop one then take 96% of the variability then we can take PC₁ and drop PC₂

$$T_{\text{stab}}(bA) * T_{\text{Vf}} = 0.7, \text{stab limit}$$

$$T_{\text{DF}} = 0.9$$

$$\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

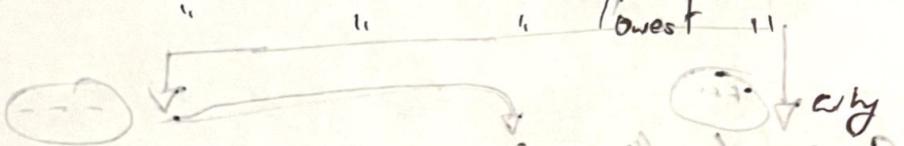
↑ 0.9 ↑ 0.9

Random Forest

GROSS IMPURITY :- which has the highest imp

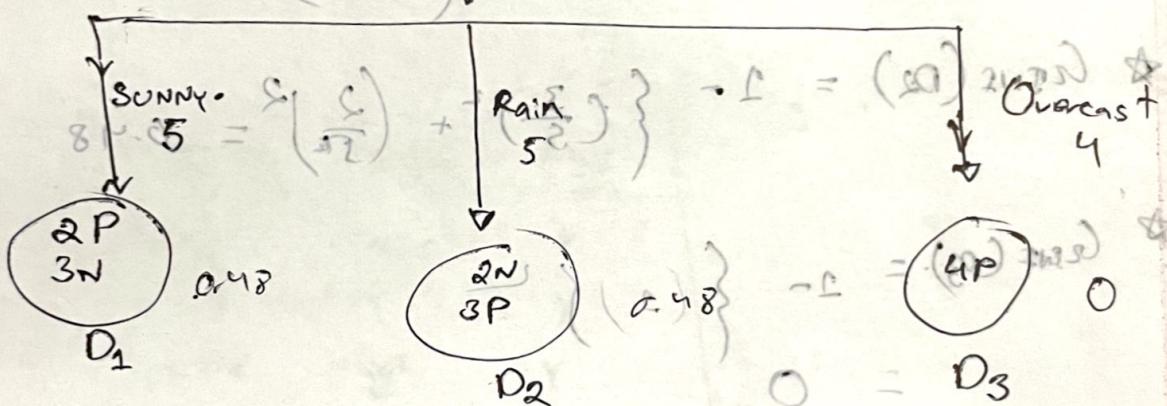
→ ↑ ↑ ↑ ↑

" " " lowest "



Impurity?

$$G.P.O = \left\{ \left(\frac{2}{2} \right) + \left(\frac{2}{2} \right) \right\} = (0.5)^2 = 0.25$$



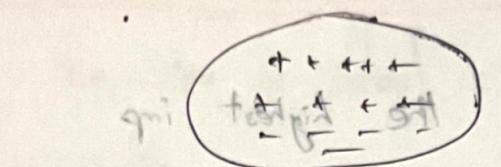
$$GGINI(D) = \frac{4}{6} \times \frac{2}{6} + \frac{2}{6} \times \frac{4}{6} = \left(\frac{2}{3} \right)^2 + \left(\frac{2}{3} \right)^2 = \frac{4}{9}$$

Genre Another any:-

$$GGINI(D) = \frac{1}{6} \times \left(1 - \frac{4}{6} \right) + \frac{2}{6} \times \left(1 - \frac{2}{6} \right)$$

$$\Sigma a(P_i) = 1 - \sum_{i=1}^n P_i^2$$

$$GGINI(D) = 1 - \left\{ \left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right\}$$



Date : 4/29/2023

$$G_{\text{SNS}}(D_1) = 1 - \left\{ \left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right\} = 0.459$$

* $G_{\text{SNS}}(D_2) = 1 - \left\{ \left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right\} = 0.48$

* $G_{\text{SNS}}(D_2) = 1 - \left\{ \left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right\} = 0.48$

* $G_{\text{SNS}}(D_3) = 1 - \left\{ \left(\frac{1}{2} \right)^2 \right\} = 0.5$

$$\rightarrow G_{\text{SNS}}(D_{\text{outlook}}) = \frac{5}{14} * G_{\text{SNS}}(D_1) + \frac{5}{14} * G_{\text{SNS}}(D_2) + \frac{4}{14} * G_{\text{SNS}}(D_3)$$

$$\left(\frac{5}{14} - \frac{1}{2} \right) * \frac{5}{14} + \left(\frac{5}{14} - \frac{1}{2} \right) * \frac{5}{14} + \frac{4}{14} * 0 = \frac{5}{14} * 0.48 + \frac{5}{14} * 0.48 + \frac{4}{14} * 0 = 0.48$$

$$\left\{ \left(\frac{5}{14} - \frac{1}{2} \right) + \left(\frac{5}{14} - \frac{1}{2} \right) \right\} - 0 = 0.48$$

$$\frac{5}{14} * \frac{1}{2} - \frac{1}{2} = 0.48$$

* Most Reduction will be the disseminator,
not the end user

Reduction of Impurity by Outlook = $0.459 - 0.343$

$$= 0.116$$

$$2.0 \quad (-+) \quad \leftarrow (S) \quad \text{Hot}$$

Decision

tree construction $\xrightarrow{\text{use}} \text{Cart algorithm}$

01/03/25

Abs to find out the most important feature?

humidity
A
temperature

Day	Temperature	Humidity	Play
10	Mild	Normal	Y
11	Mild	N	Y
12	Mild	High	Y
13	Hot	N	Y
14	Mild	High	N
15	Hot	High	N

$$D = \frac{++}{2} + \frac{+-}{2}$$

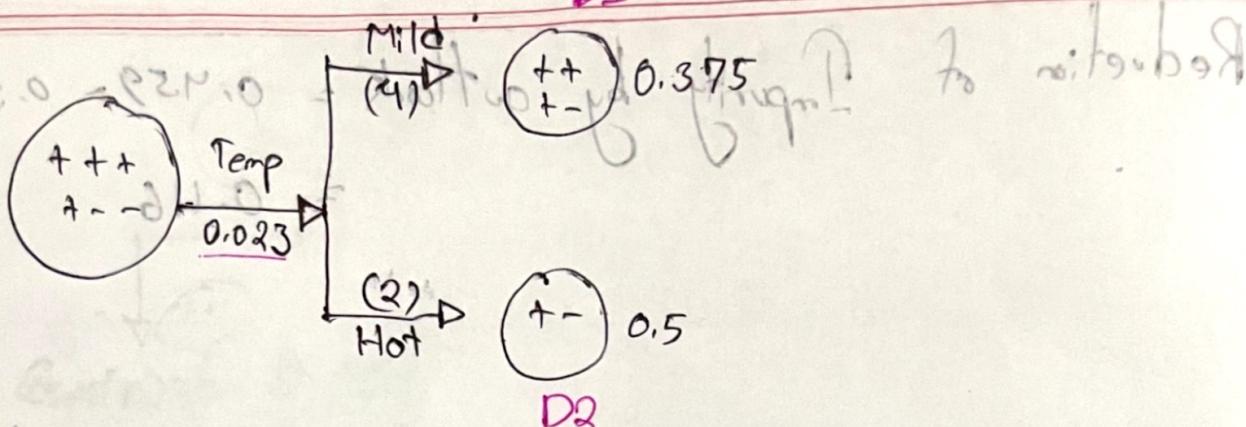
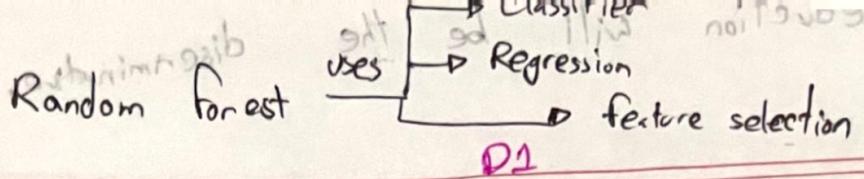
$$GINI(D) = \frac{1}{10} \left\{ P_+^2 + P_-^2 \right\}$$

$$= 0.44 \times \frac{5}{10} + 0.56 \times \frac{5}{10} = \frac{1}{2} \left\{ \left(\frac{5}{10}\right)^2 + \left(\frac{5}{10}\right)^2 \right\}$$

$$= 1 - \frac{16}{36} = \frac{4}{9}$$

$$= 1 - \frac{20}{36}$$

$$= \frac{16}{36} = \frac{4}{9} = 0.44$$



$$GINI(D_1) = 1 - \left(P_+^2 + P_-^2 \right)$$

$$\text{Count} = \frac{1}{16} \left\{ \left(\frac{3}{8}\right)^2 + \left(\frac{5}{8}\right)^2 \right\} \rightarrow \text{Grid of left}$$

$$= 1 - \left(\frac{9}{16} + \frac{1}{16} \right) \rightarrow \text{Left side} = 0.5$$

$$= 1 - \frac{10}{16} \rightarrow \text{Left side} = 0.5$$

$$= \frac{6}{16} \rightarrow \text{Left side} = 0.375$$

$$= \frac{3}{8} = 0.375 \rightarrow \text{Left side} = 0.375$$

$$GINI(D_2) = 0.5$$

$$\Rightarrow GINI_{Temp}(D) = \frac{|D_1|}{|D|} GINI(D_1) + \frac{|D_2|}{|D|} GINI(D_2)$$

$$= \left(\frac{9}{16} \right) + \left(\frac{7}{16} \right) = \frac{9}{16} * 0.375 + \frac{7}{16} * 0.5$$

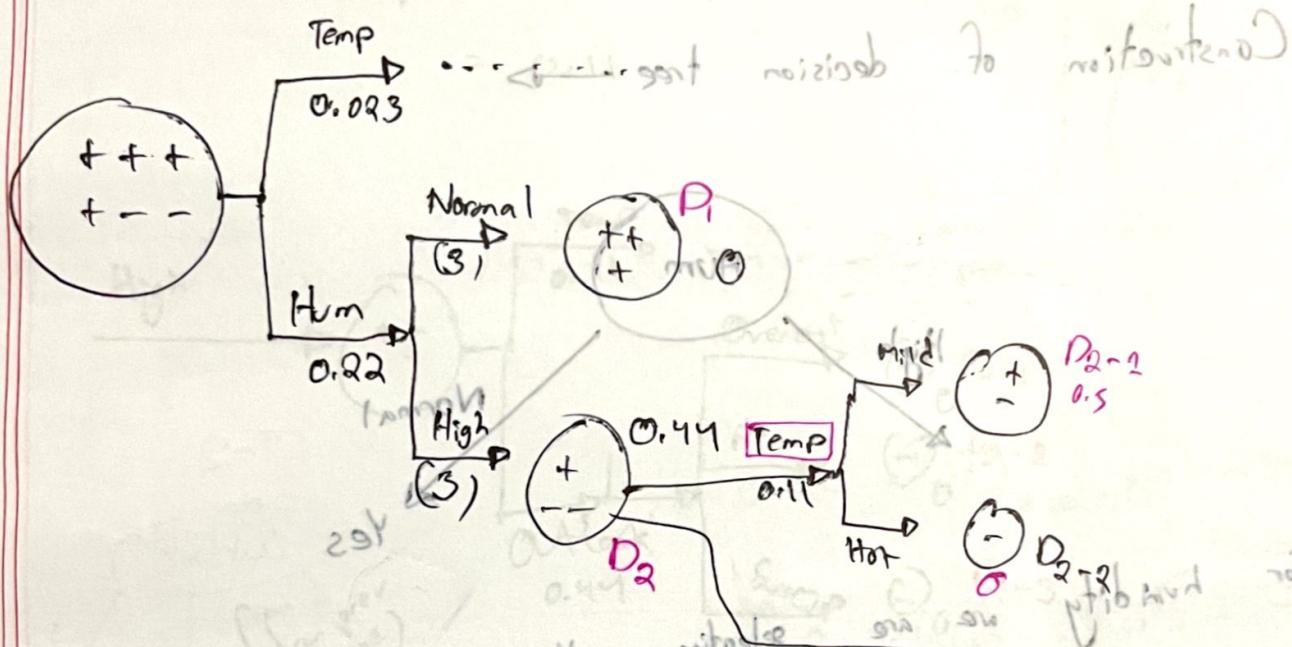
$$\frac{9}{16} * 0.375 = 0.417$$

$$\frac{7}{16} * 0.5 = 0.25$$

$$0.417 + 0.25 = 0.667$$

Reduction of G_{INSI} by Temp = $0.44 - 0.41$ $\Rightarrow 0.03$

Reduction Temp = $0.44 - 0.33$



$$G_{\text{INSI}}(D) = 1 - \left\{ \left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right\} \cdot 0.44$$

$$= 1 - \frac{1}{9} - \frac{4}{9} = 0.44$$

Reduction at look = $0.44 - 0.33 = 0.11$

$$G_{\text{INSI}}_{\text{HUM}}(D) = \frac{|D_1|}{|D|} G_{\text{INSI}}(O_1) + \frac{|D_2|}{|D|} G_{\text{INSI}}(O_2)$$

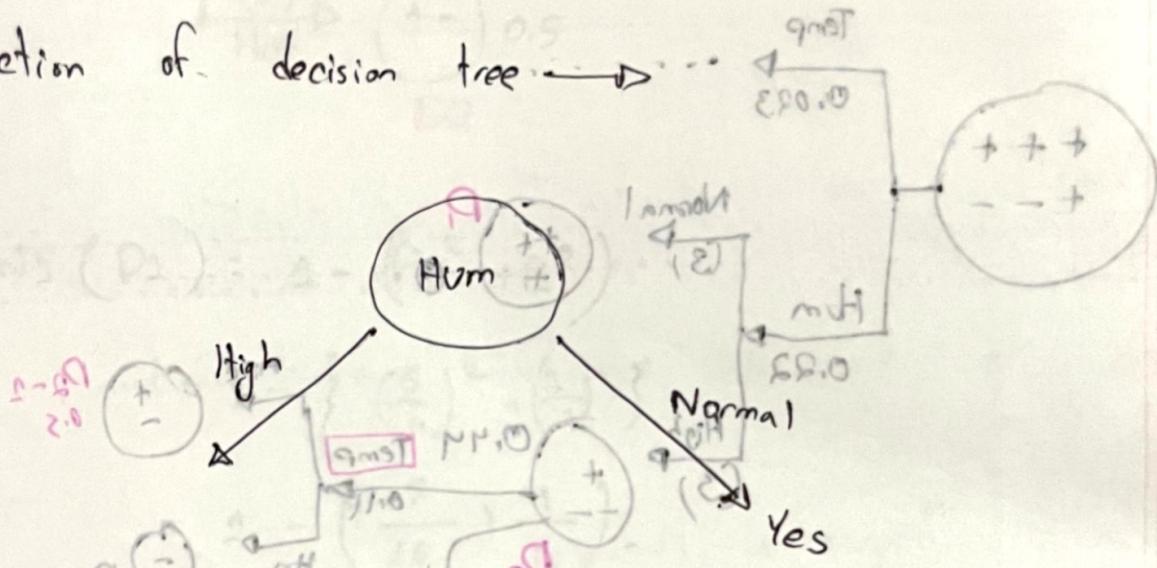
$$= \frac{s}{6} * 0 + \frac{s}{6} * 0.44$$

$$= 0.22$$

have to take at least two col.

Reduction of GINI by humidity = $0.44 - 0.22$ ~~0.44 - 0.22~~
 $E_{\text{Gini}} = 0.22$

Construction of decision tree \rightarrow



For humidity we are selecting another two columns from the bootstrapped dataset \rightarrow

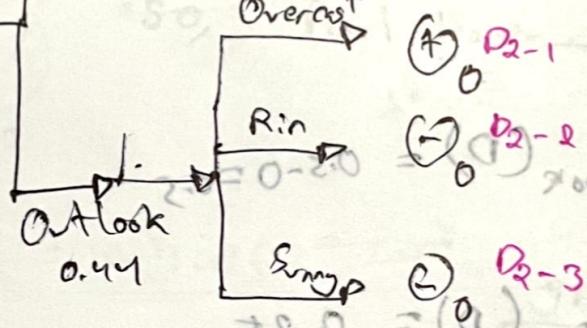
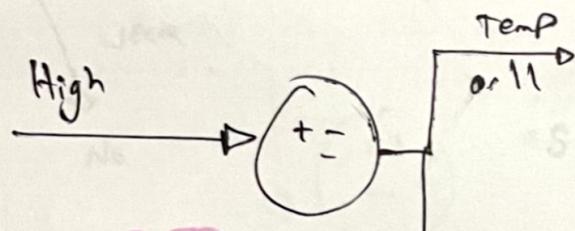
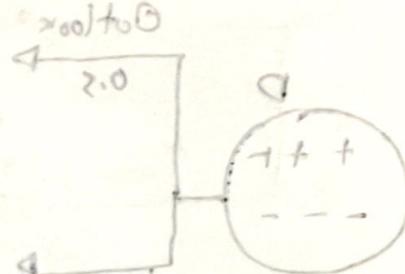
Day	Outlook	Temperature	Play	$=$
12	Overcast	Mild	Y	$MM.O =$
14	Rain	HOT	N	$\frac{1}{10} = (1)_{\text{Mild}}$
2	Sunny	Mild	N	$\frac{3}{8} =$

How many bootstrap do we need?
 or No. of trees we need

$$G_{\text{INST}}(DQ) = \frac{2}{3} \times 0.3 + 2 + \frac{5}{3} \times 0 = -0.83$$

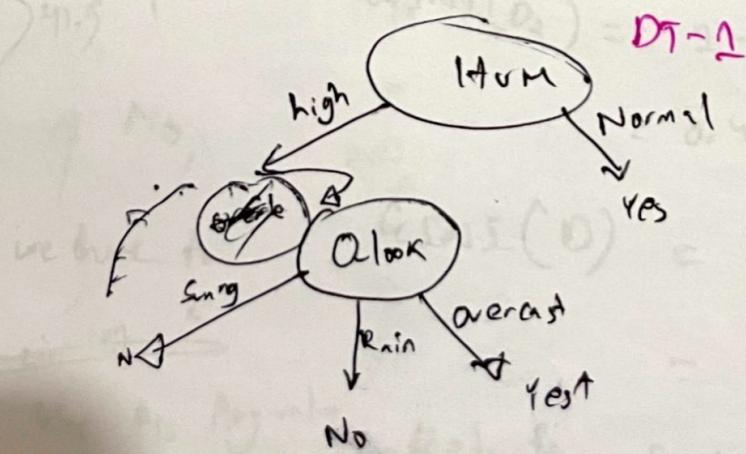
Reduction Temp = $0.44 - 0.33$

$$= 0.11$$



$$G_{\text{INST}}(\text{outlook}) = \frac{1}{3} * 0 + \frac{1}{3} * 0 + \frac{1}{3} * 0 = 0$$

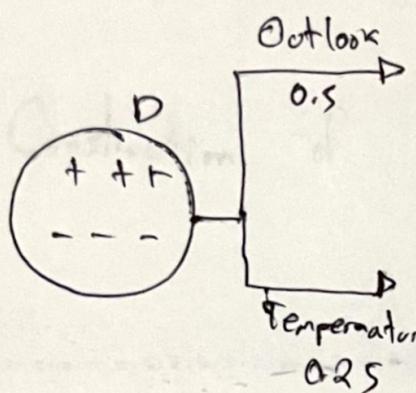
Reduction outlook = $0.44 - 0 = 0.44$



Bootsrap DATASET -2

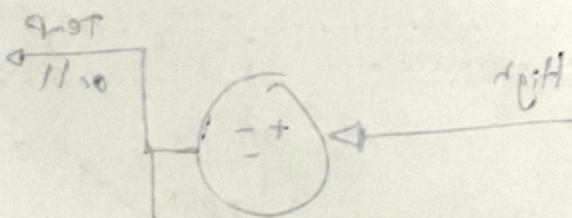
Bootstrapped - 3
(A, U, J)

$$E2.0 = \text{Decision} + \text{Error} = (SD)_{\text{min}} 25.25$$



$E2.0 - NN.0 = \text{Error}$ reduction

$$NN.0 =$$

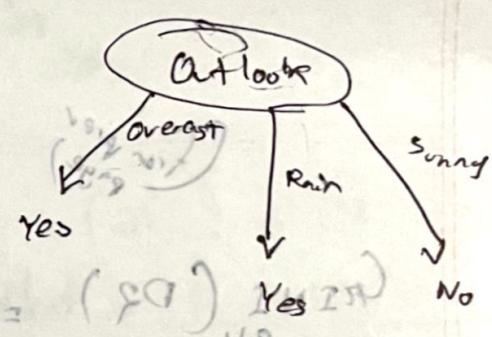


Reduction $E2.0 - NN.0$

$$\text{Reduction}_{\text{Outlook}}(D) = 0.5 - 0 = 0.5$$

$$\text{Reduction}_{\text{temp}}(D) = 0.25$$

DT-2



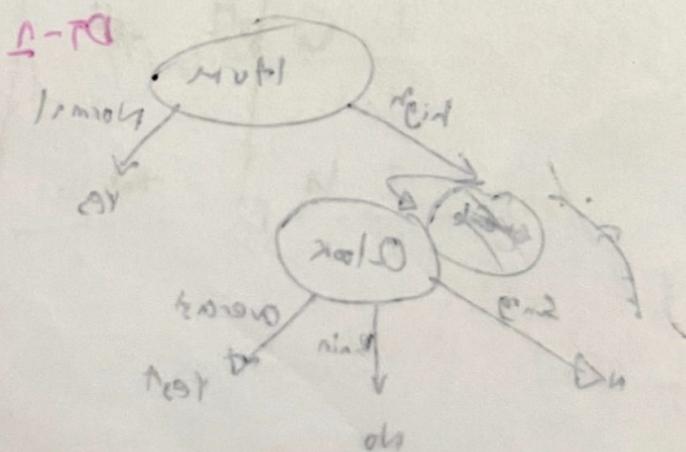
$$0 \times \frac{1}{2} + 0 \times \frac{1}{2} + 0 \times \frac{1}{2} = (SD)_{\text{min}} 25.25$$

$$O =$$

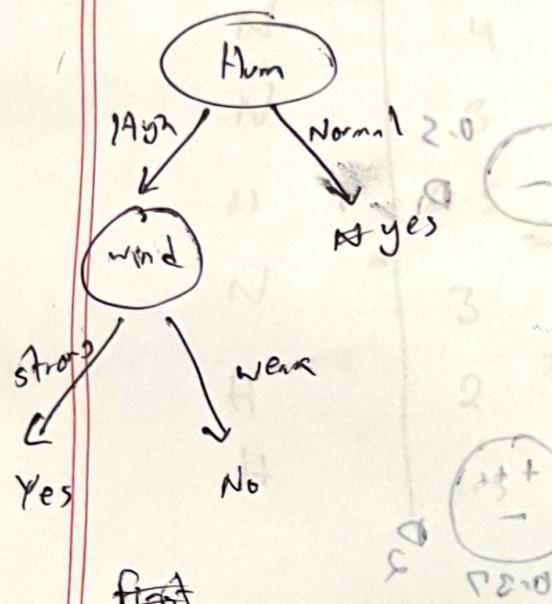
overcast

$$NN.0 = 0 - NN.0 = x_{001} + x_{010} \quad \text{reduction}$$

DT-3



DT-3



Bagging approach:-

DT-1 → 1
Yes 2.0
No 2.0

DT-2 → 1
Yes 2.0
No 2.0

DT-3 → 1
Yes 3.0
No 3.0

yes with wins

so final decision → yes

first

2st → sort

$$\text{Temp} = \frac{P}{R} = \frac{20 \times \frac{2}{3}}{\frac{1}{2}} = 20 \times \frac{4}{3} = 26.67$$

$$20 > 20.5 \quad Y$$

$$21 > 20.5 \quad N$$

$$25 > 20.5 \quad Y$$

$$27 > 20.5 \quad Y$$

$$31 > 20.5 \quad Y$$

$$41 > 20.5 \quad Y$$

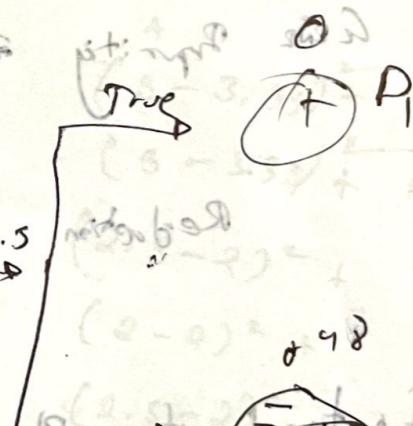
$$42 > 20.5 \quad N$$

we have to

use Avg

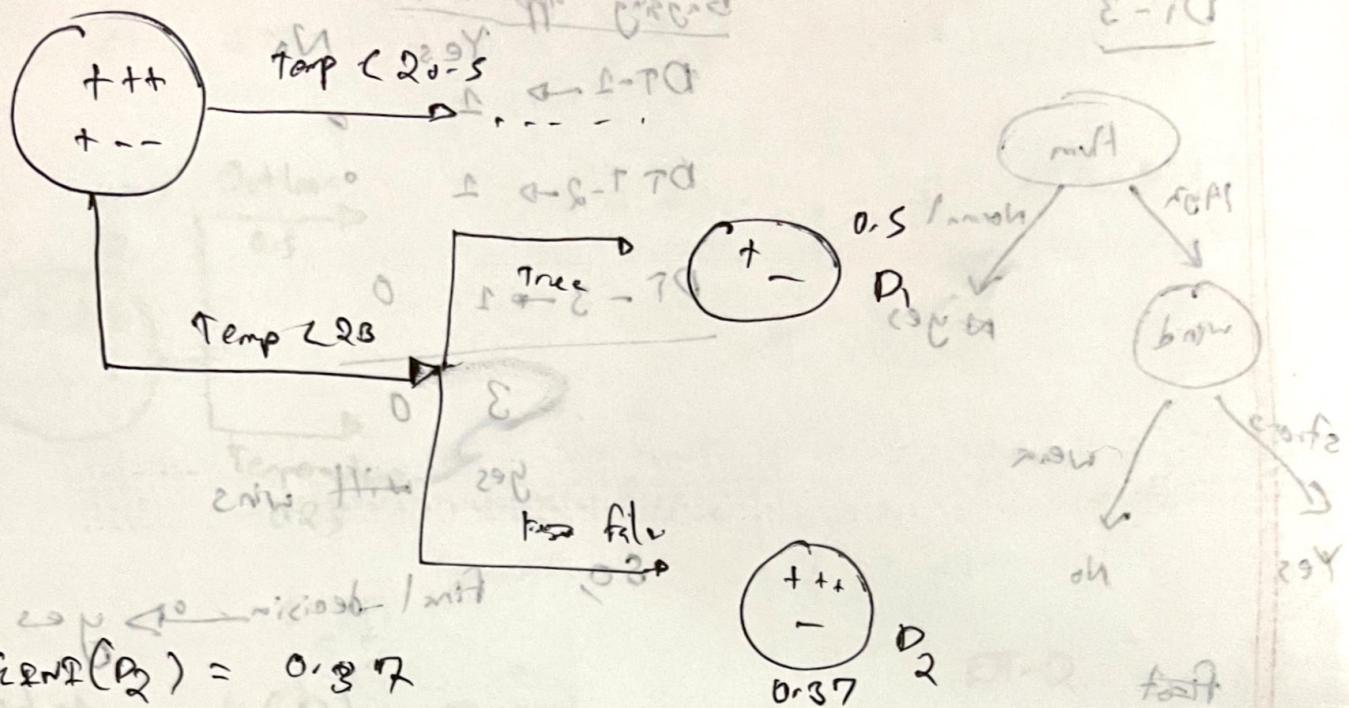
use this Avg value
as splitter

Temp > 20.5 →



cutting & cutting

utilizing



$$G_{ENT}(D_2) = 0.37$$

Leave Priority ~~$G_{ENT}(D)$~~ = $G_{ENT}(D) = \frac{2}{6} \times 0.5 + \frac{4}{6} \times 0.37$

$$= 0.167 + 0.2467 = 0.41287$$

Reduction ~~$G_{ENT}(D) = 0.45$~~ $2.087 - G_{ENT}(D) = 2.087 - 0.41287 = 1.674$

Instead of Play tennis \rightarrow play tennis

Likability score Likability \rightarrow regression proba

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

$$8^{0.8} = \frac{3}{2.5}$$

$$8^{0.8} \times \frac{2}{3} + 0 \times \frac{1}{3} = (0)2.623$$

$$P_{-B} = 1.3$$

$$- P_{-B} = -1.3 \rightarrow \text{not good}$$

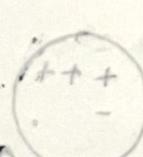
win 0.2 2nd 50%
lose 0.8 50%

Hvn	Lav
N	4
N	3
H	2.5
N	3
H	2
H	1.5

if ($Hvn = N$)

True

False



2nd row
Avg $\rightarrow 3.33$

4, 3, 3

2.5, 2, 2.5

Avg $\rightarrow 3.33$

Avg $\rightarrow 2.5$

$$SSR_{(Hvn = N)} = (4 - 3.33)^2 +$$

$$\left\{ \left(\frac{3}{n} \right) + \left(\frac{3}{n} \right) \right\} - 1 = \text{Actual} - \text{Pred}$$

$$\frac{s^2 + s^2}{2} - 1 = (3 - 3.33)^2 +$$

$$(3 - 3.33)^2 +$$

$$(2.5 - 2)^2 +$$

$$(2 - 2)^2 +$$

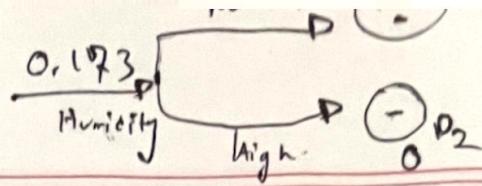
$$(2.5 - 2)^2 +$$

$$\left\{ \left(\frac{1}{n} \right) + \left(\frac{1}{n} \right) \right\} - 1 = 2.1667$$

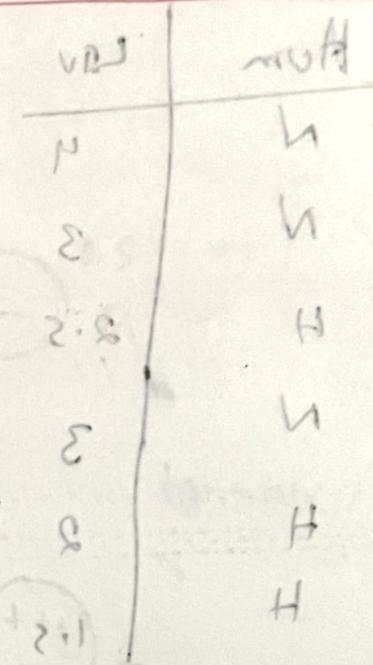
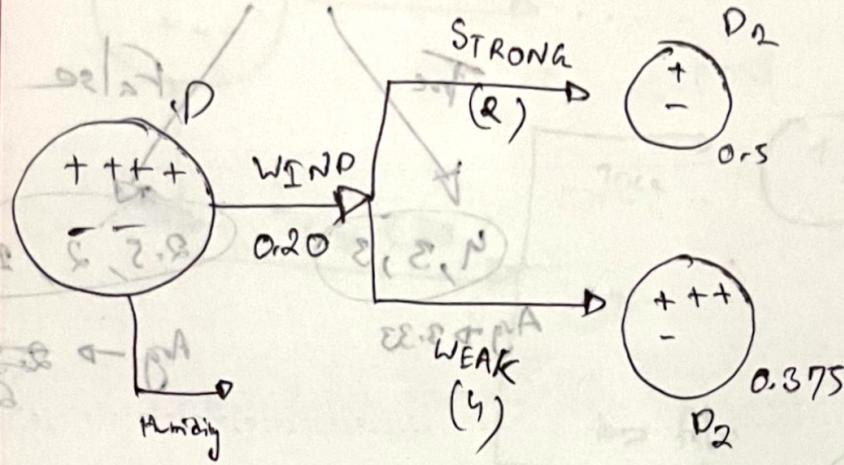
$$\frac{0.5}{2} - 1 =$$

$$250.0 - 1 =$$

$$250.0 =$$



BOOTSTRAPPED DAT - 3



$$+ (\varepsilon \cdot \varepsilon - N) = M = MVH$$

$$C_{INS}(D) = 1 - \left\{ \left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right\}$$

$$+ (\varepsilon \cdot \varepsilon - S) = 1 - \frac{4^2 + 2^2}{36}$$

$$+ (\varepsilon \cdot \varepsilon - S) = 1 - \frac{20}{36}$$

$$+ (\varepsilon \cdot \varepsilon - S) = 0.444$$

$$C_{INS}(D_1) = 0.5$$

$$C_{INS}(D_2) = 1 - \left\{ \left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right\}$$

$$= 1 - \frac{20}{36}$$

$$= 1 - 0.625$$

$$= 0.375$$

$$G_{\text{IND WIND}}(D) = \frac{2}{6} \times 0.5 + \frac{u}{6} \times 0.375$$

~~2~~
= ~~0.25~~
Or $0.417 R$

Reduction of G_{IND} by Temp = ~~0.445 - 0.25~~ $0.417 R$

~~0.445 - 0.25~~ of ~~0.445~~ To balance
= ~~0.20~~ 0.23

$\left[\begin{array}{l} \text{of } A \\ \text{principle} \end{array} \right]$

$$G_{\text{IND hum}}(D_1) = 1 - \left\{ \left(\frac{u}{5}\right)^2 + \left(\frac{1}{5}\right)^2 \right\}$$

~~of Planing~~ $= 0.32$

$$G_{\text{IND hum}}(D_2) = 0$$

$$G_{\text{IND hum}}(D) = \frac{5}{6} \times 0.32 + 0$$

~~(E-E)~~
 $= 0.267 R$

Reduction of G_{IND} by $\text{num} = 0.44 - 0.267$

$$= 0.173$$

~~0.445 - 0.267~~

~~0.178~~

$E = \text{DNA}$

$$\therefore G_{\text{IND hum}}(E-E) + G_{\text{IND hum}}(E-A) + G_{\text{IND hum}}(S-E) + G_{\text{IND hum}}(S-A) = 0.267$$

$S =$

We need to visualize the data  If follows linear curve

\rightarrow Non linear algorithm come
If doesn't follow

If doesn't follow

- ✓ SSR instead of Gini for Regression tree.
- ✓ loss function

Considering Age]

$$\text{Age} < 9 \left\{ \begin{array}{l} \xrightarrow{\text{using}} \text{Her. } \left(\frac{n}{z} \right) \\ \end{array} \right\} - 1 = \left(\text{f}^n \right)_{\text{mod } z}$$

$$A \cdot g = 2$$

$$\begin{array}{c} \cancel{3} \\ 3, 2, 4, 5 \\ \hline 8 \end{array} \quad \text{Avg} = 3$$

Residual = ~~Original~~ Original - Predicted
= ~~Pred~~ PQ

$$= \begin{pmatrix} 0 \\ 1-2 \end{pmatrix}$$

$$+ \frac{1}{(x-3)^2}$$

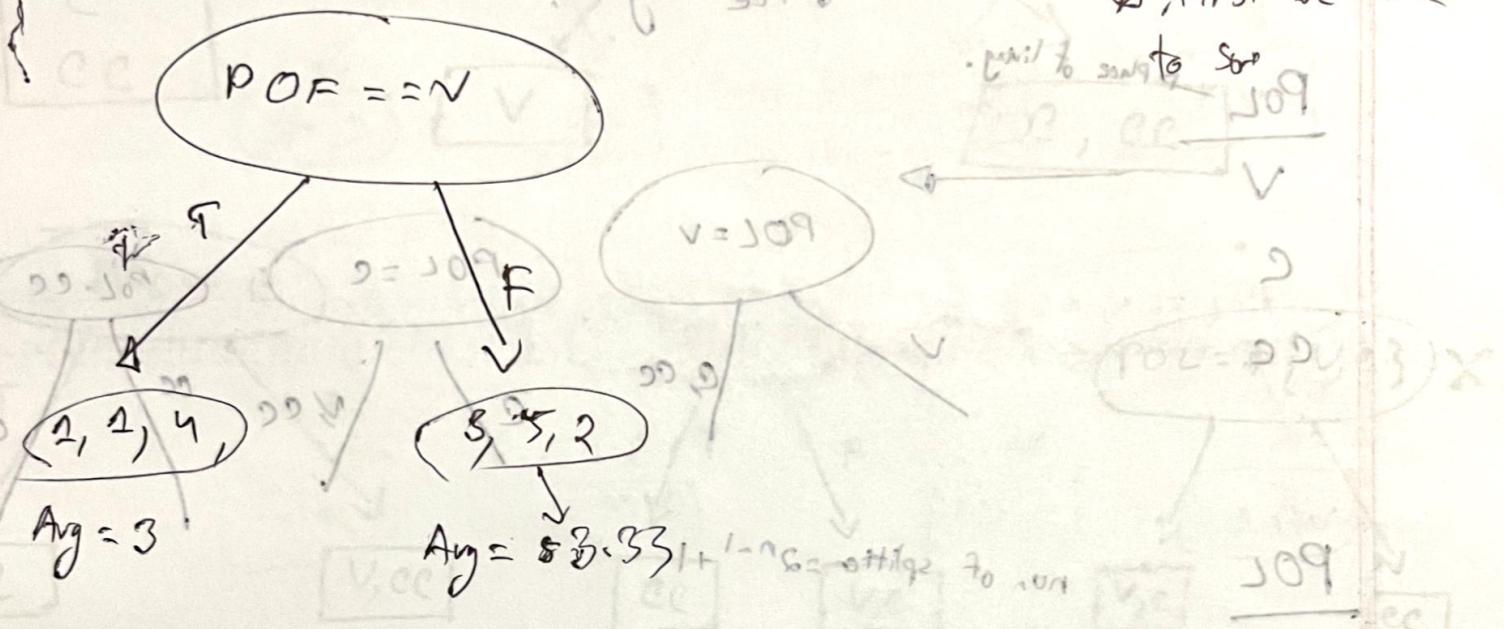
$$0 + 58.0 \times \frac{3}{2} = 1((3-3)^2_{\text{min}} + (2-3)^2_{\text{max}})$$

$$\begin{array}{c} \nearrow (3) \\ \times (y-3)^2 + (x-3)^2 \\ (x-3)^2 \end{array}$$

Age < 22
 T
 F
 2, 4, 5 + 2,
 $\text{AVG} = 3$

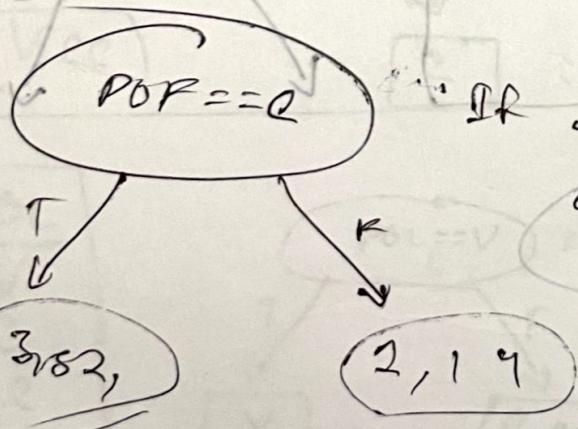
$$SSR = (2-2)^2 + (3-2)^2 + (4-3)^2 + (7-3) \cdot (5-8)^2 + (2-3)^2 \\ = 12$$

Aga-21 provides the lowest SSR for 21 out of 27 species.



$$SSR_{\text{village}} = (1-3)^2 + (1-2)^2 + (1-2)^2 + (5-3)^2$$

$(2-3.33)^2$

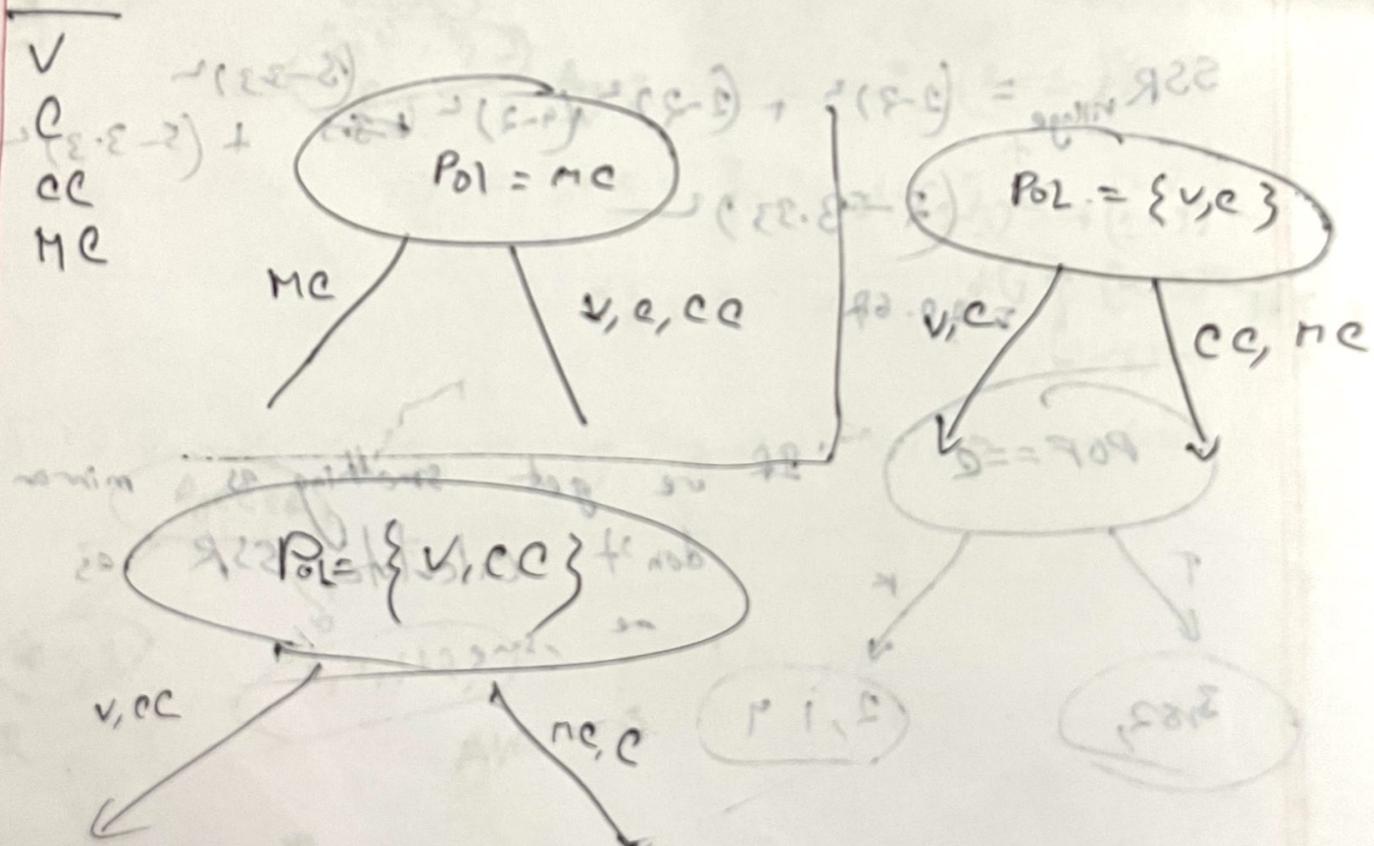
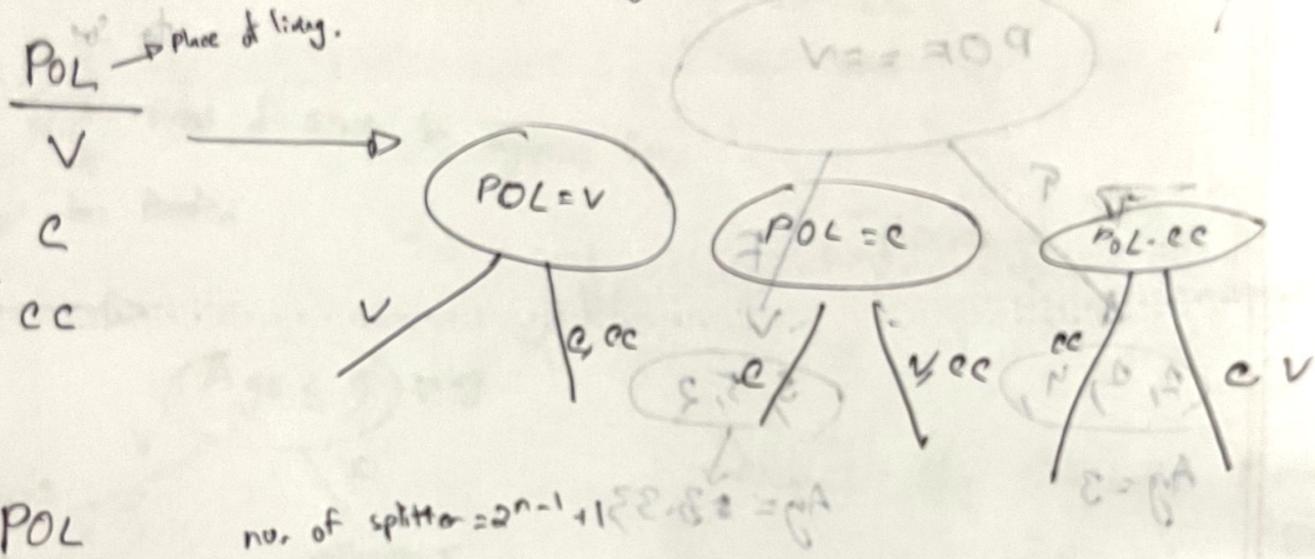


If we get something as minor we
don't calculate SSR as they
are same

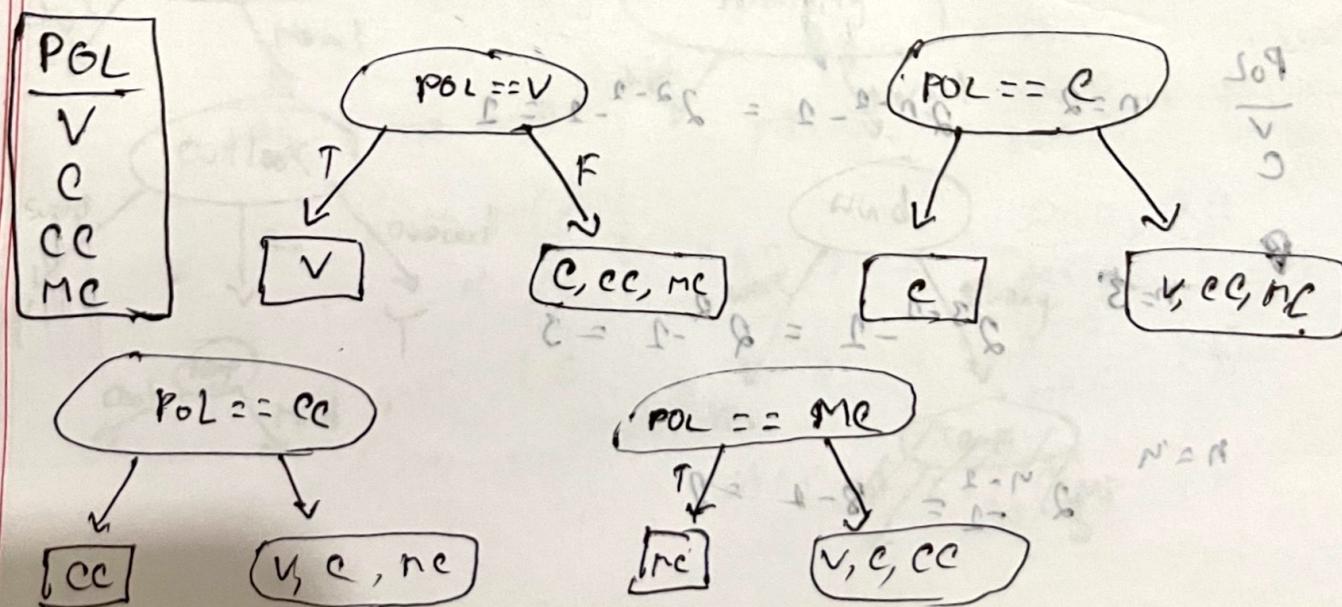
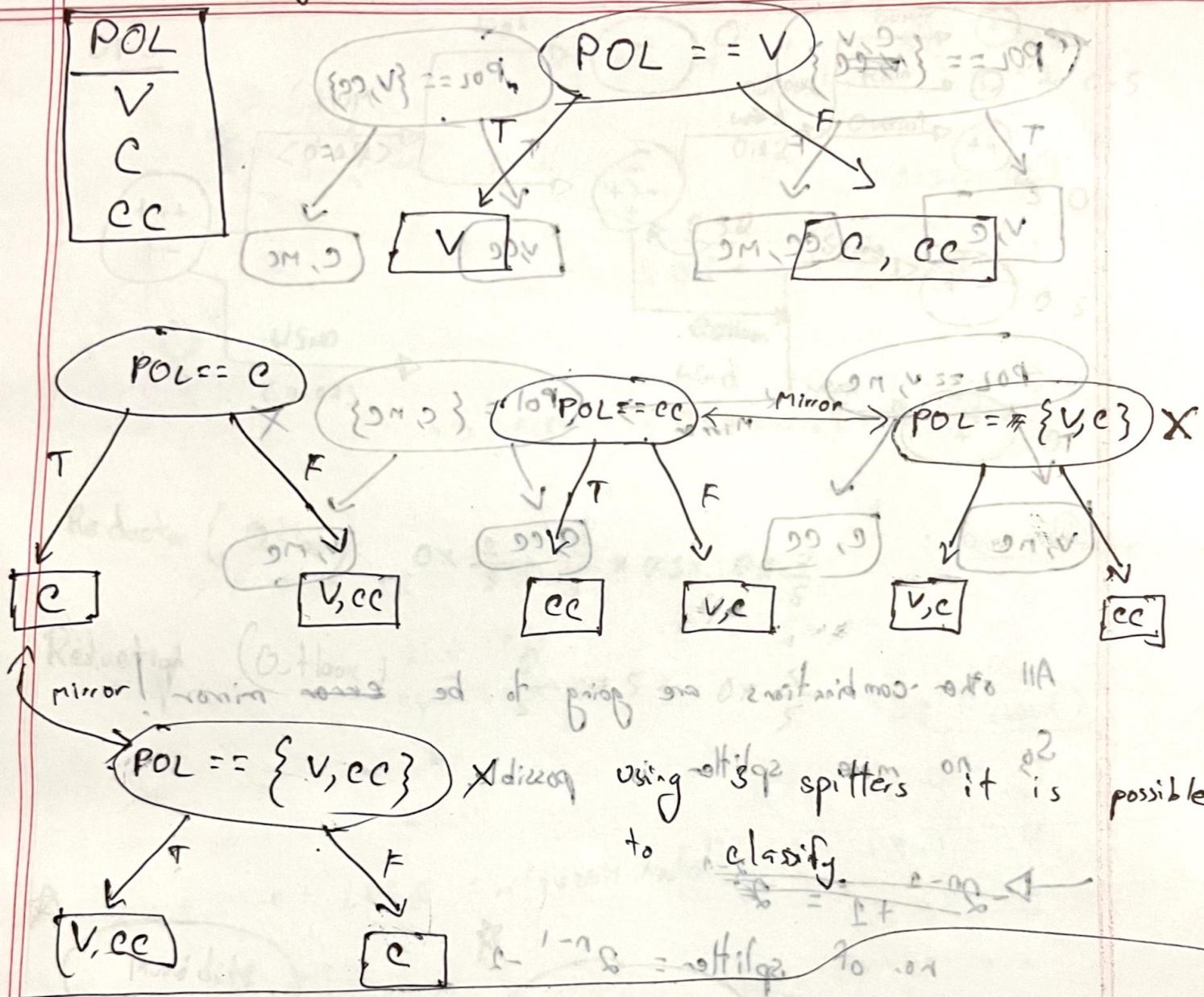
$$\text{SSR}_{age} = 7.155$$

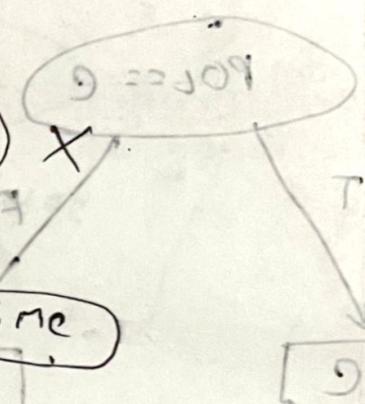
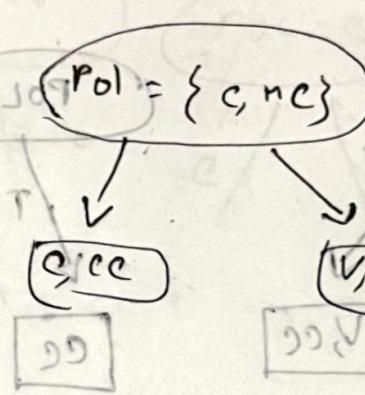
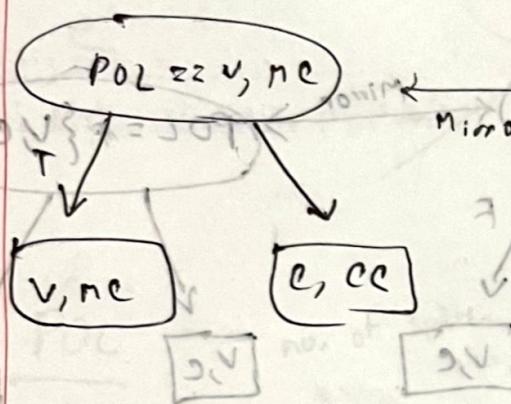
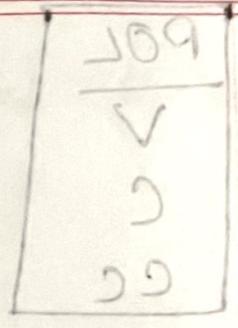
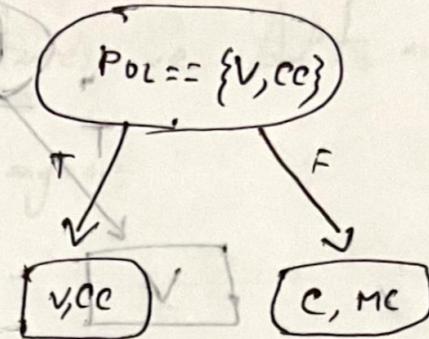
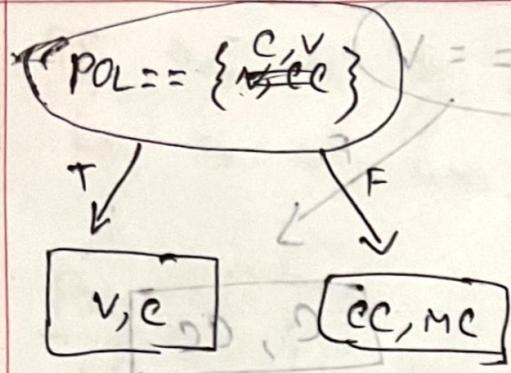
$$SBR_{VII, R} = 20.6\%$$

If there is only 1 instances we don't need to expand them. The 1st instances is tree anymore.



Place of Living



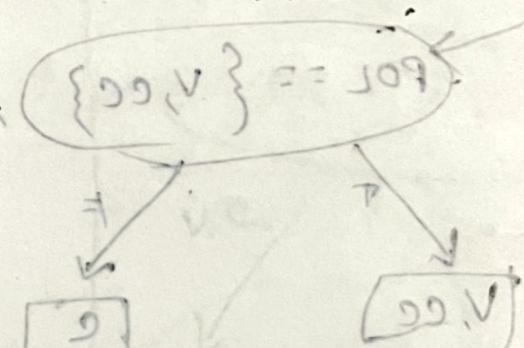


All other combinations are going to be error minor!

So, no more splitter is possible.

$$\rightarrow 2^{n-1} + 1 = 2^{n-1}$$

$$\text{no. of splitter} = 2^{n-1} - 1$$



$\text{POL} = \{V, C\}$

$n=2$

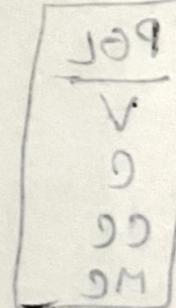
$2^{2-1} - 1 = 2^1 - 1 = 1$

$\therefore \text{no. of splitters} = 1$

$\text{POL} = \{V, CC\}$

$n=3$

$2^{3-1} - 1 = 2^2 - 1 = 3$



$n=4$

$2^{4-1} - 1 = 8 - 1 = 7$

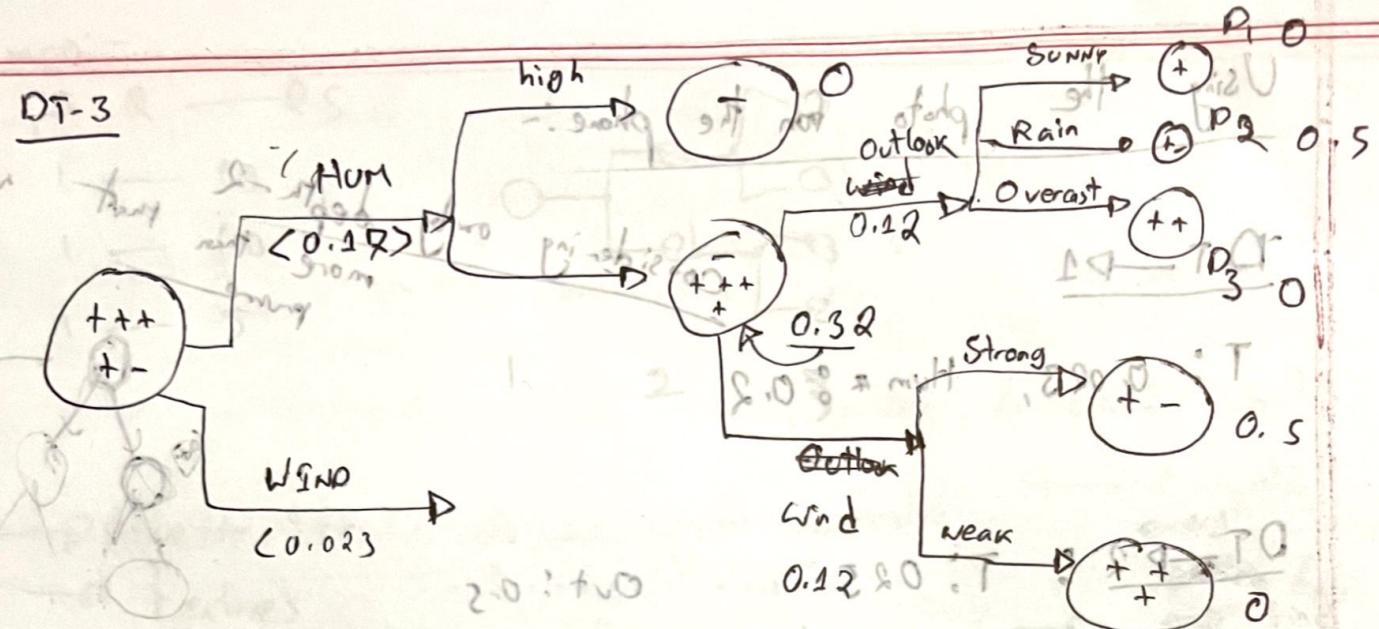
$\text{POL} = \{V, CC, NC\}$

$\text{POL} = \{V, NC\}$

$\text{POL} = \{CC, NC\}$

$\text{POL} = \{CC\}$

DT-3

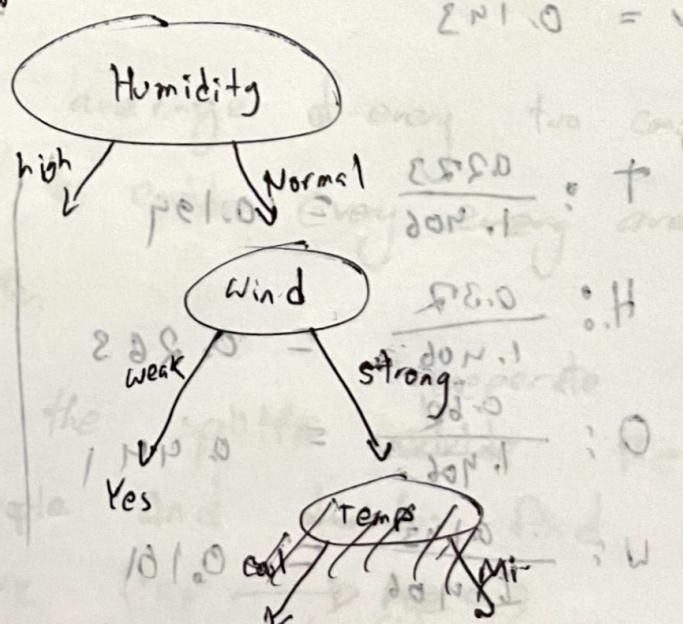
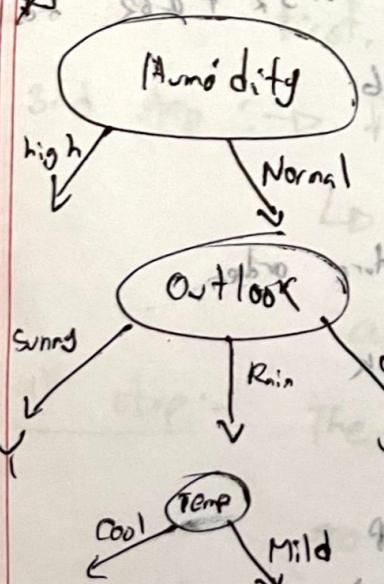


Reduction (Outlook) =

$$0 \times \frac{1}{3} + \frac{2}{3} \times 0.5 + 0 \times \frac{2}{3}$$

Reduction (Outlook) =

$$\frac{2}{3} \times 0.5 + 0 \times \frac{3}{5}$$



Using the photo from the phone:-

E-70

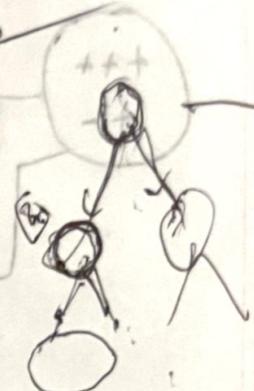
DT -> D1

$$T: 0.023, \text{ Hum} = 0.2$$

$$\underline{DT -> D2}: T: 0.25$$

$$\text{Out: } 0.5$$

only depth-2
more than
parent



DT -> D3:

$$\frac{\sum}{2} \text{Hum: } 0.17 \times \frac{5}{2} \text{ Wind: } 0.023$$

(root) not used
~~0.12~~

$$\text{Out: } 0.12$$

(root) not used

Overall %:-

$$T \approx 0.273$$

$$H = 0.87$$

$$O = 0.62$$

$$W = 0.143$$

$$\text{Total Reduction} = 0.273 + 0.32 + 0.62 = 1.206$$

Now,

$$T: \frac{0.273}{1.206} = 0.194$$

$$H: \frac{0.32}{1.206} = 0.263$$

$$O: \frac{0.62}{1.206} = 0.491$$

$$W: \frac{0.143}{1.206} = 0.101$$

Best feature order

1) Outlook

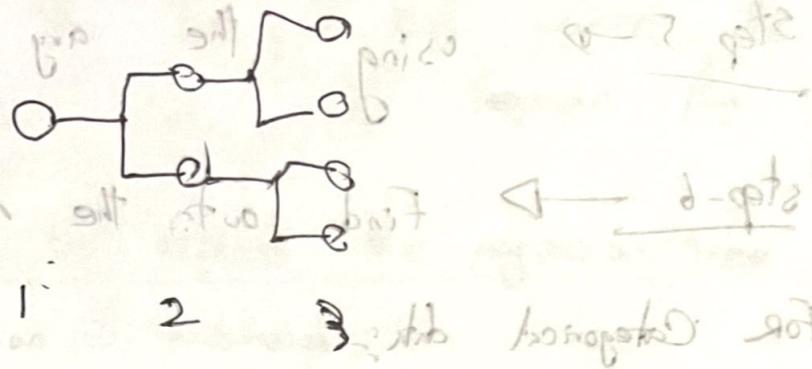
2) Hum

3) Temp

4) Wind

Ans:-

95 Q → 25
1 → 10
1 (node) → 50
1 → 50



→ Regression trees: To find pathologies

→ features

$$n = 100$$

type of features

SDP

categorical

numerical

find SSR

every splitter

find SSR

REGRESSION TREE:-

For 1st step → identify the column type feature. To

2nd step → for continuous values (numerical) sort the numbers first.

3rd step :- find the average of every two consecutive data points. → you have to consider every average as a splitter.

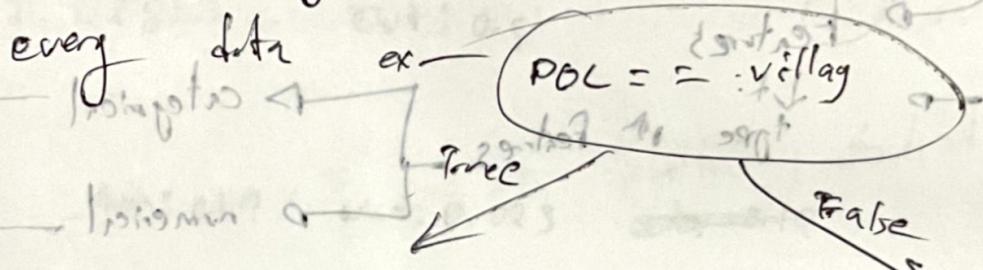
4th step:- Then using the splitter divide the IQ scores people and find the average of the value like $\rightarrow \text{Age} < 9$.
1. $\leftarrow \text{avg}$ 2. $\leftarrow \text{left}$ $\rightarrow \text{right}$
 $3, 4, 5, 2, 1$

Step 5 → using the arg find the SSR value.

Step 6 → Find out the minimum SSR value (which has low error)

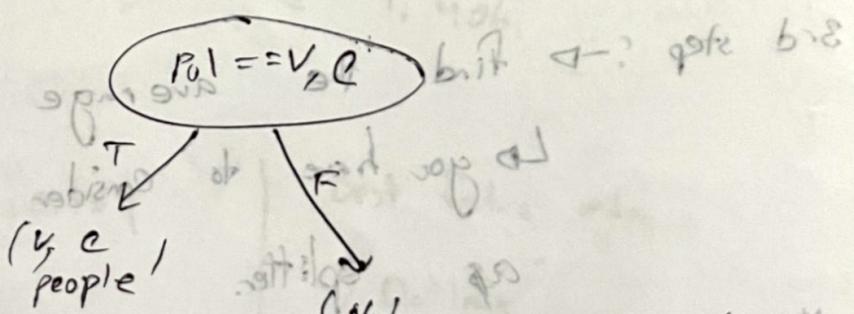
For Categorical data:

Step 4 → selecting any one of sector apply True or False



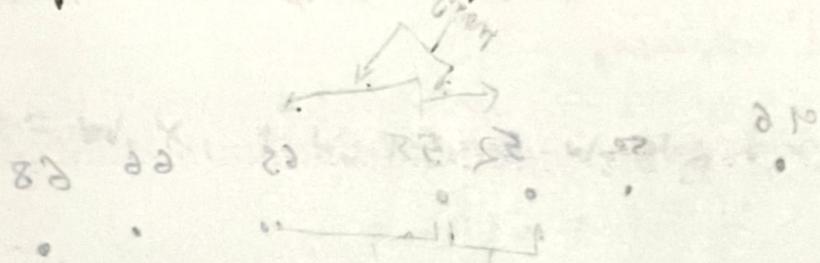
Repeat & find arg of SSR, classify & find the best value

exception: If there are more than 4 categorical values



Step 7 Comparing the two SSR values, the one which has less will become the splitter.

Step 8 :- You have to expand the regression tree up to expansion threshold.



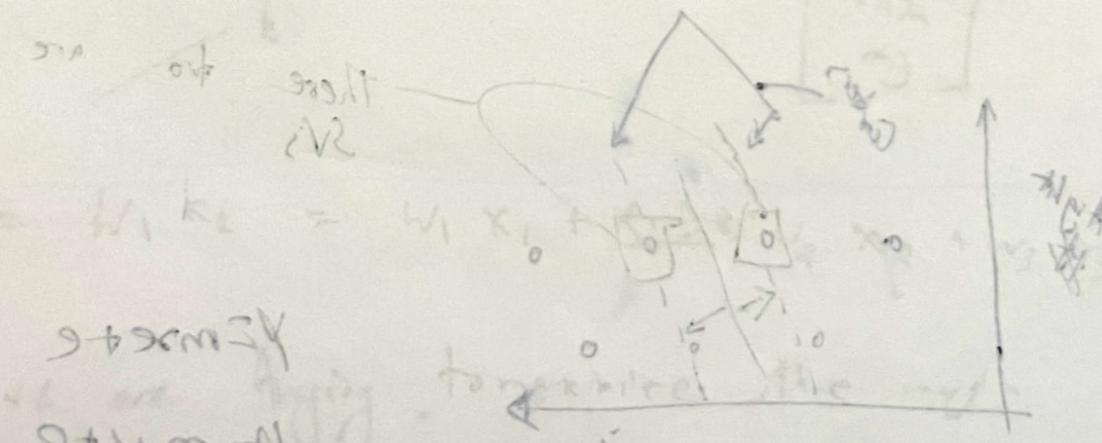
$$E = H = \sum_{i=1}^n w_i$$

$$E = \sum_{i=1}^n w_i$$

...
pseudo

$$\frac{\text{pseudo-total} + \text{extreme-right}}{2} = H$$

$$H = \frac{E_1 + E_2}{2}$$



Step 4
if we trying to minimize the error of the tree
it will be going to the tree

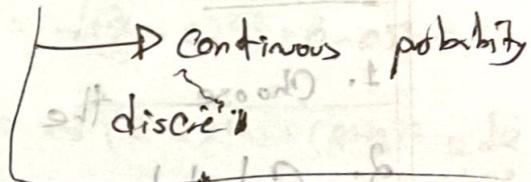
INFERENTIAL STATISTICS

We want to infer something from the data.

We "prove a hypothesis."

* It is random sample of data taken from population to describe and make inferences about the population.

Discrete



HYPOTHESES TESTING

→ It can really be anything as long as you pt. it along the test

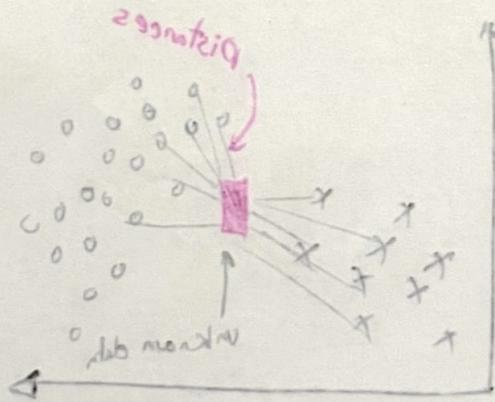
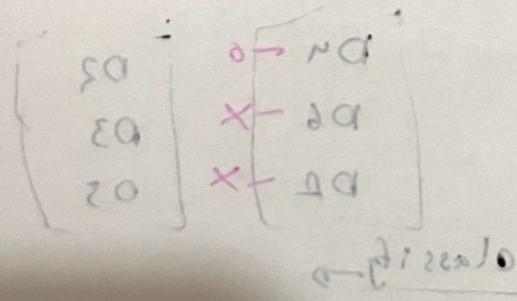
standard score / Z score
(i) subtract the mean

Hypothesis testing

to calculate the standard D.

(i) One tailed test

(ii) two tailed test



KNN Algorithm

for large dataset

$K \leftarrow$ hyperparameter, since

High space complexity

doesn't need a training set
require

Rules →

1. Odd number

\sqrt{n} where $n =$ total no. of data.
Ex. $n = 100$

$$K = \sqrt{100} \approx 10.0$$

STEPS:-

↳ 3, 5, 7, 9

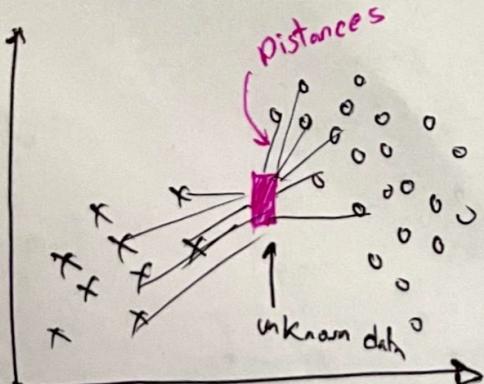
1. Choose odd no. of K .

(odd no.) except 1.

2. Calculate the distance between the new points in the and all points in the data.

3. Sort the distance values.

4. Assign the class with most common among the K neighbours.



sort →

D_4	D_2	$(K=3)$
D_6	D_3	
D_2	D_5	

classify →

majority is 'X' so, classify as

$$\text{Example B: } \left(1, 2\right) = (0, 0) \Delta$$

	X	$S_{2,1} = 1 - 2 \cdot 2$	label $D(P, Q)$
--	-----	---------------------------	--------------------

- A. 1 2 Red $n=4$ $\sqrt{n} = \sqrt{4} = 2$
 $2 \cdot 0 = (0, 0) \Delta$
- B. 2 3 Red $2 \cdot 1 = 2 \cdot 1$
- C. 3 3 Blue $K=3$ $2 \cdot 0 = (0, 0) \Delta$
- D. 6 5 Blue $D(P, A) = \sqrt{(3-1)^2 + (2-2)^2} = 2$
 $2 \cdot 0 = (0, 0) \Delta$
- P. 3 2 $D(P, B) = \sqrt{(3-2)^2 + (2-2)^2} = \sqrt{2}$
 $2 \cdot 0 = (0, 0) \Delta$

$$D(P, C) = \sqrt{(3-3)^2 + (2-3)^2} = \sqrt{2}$$

$$D(P, D) = \sqrt{(3-6)^2 + (2-5)^2} = \sqrt{18} = 4.24$$

Sorting $\rightarrow (1, 1) = (8, 9) \Delta$

$$D(P, C)$$

$K=3$

$$D(P, B)$$

\hookleftarrow

$$D(P, C) = R$$

$\epsilon = 2$

$$D(P, A)$$

$$D(P, B) \in R$$

$$D(P, D)$$

$$D(P, A) = R$$

Majority is

RED.

$$22.8 \Leftarrow$$

Regression,

EXAMPLE-2

X	Y	$D(P, A) = 3.5 - 1 = 2.5$
A 1	1	$D(P, B) = 3.5 - 2 = 1.5$
B 2	2	$D(P, C) = 3.5 - 3 = 0.5$
C 3	1.5	$D(P, D) = 3.5 - 4 = 0.5$
D 4	3.75	$D(P, E) = 3.5 - 5 = 1.25$
$\Sigma = 15 - 9 + 5(2 - 3) = (A, 9) \Delta$		
$(P_B) + 3(5 - 3) h = (0.9) \Delta$ Sort \rightarrow		
$K = \sqrt{5} \quad n=5$		$D(P, C) = 0.5$
$= 2.25$		$D(P, D) = 0.5$
$K = 3$		$D(P, B) = 1.15$ <u>error</u>
Mean $\Rightarrow \bar{x} = (0.9) \Delta$		$(0.9) \Delta$
$0.5 + 0.3 + 1.5 + 2.25 + 3 = 7.5$		$D(P, E) = 1.25$ <u>values</u>
$\bar{x} = (0.9) \Delta$		$(0.9) \Delta$
$= 2.25$		
$\Rightarrow 2.25$		

Example - 3

$N = 5 / \max$

	length	width	weight	label	envelope
A	5	3	150	Apple	8
B	6	3.5	160	Apple	8
C	7	4	180	Apple	9
D	10	6	300	Mango	10
E	11	6.5	320	Mango	10
P	6.5	3.8	170	?	9

$$n = 5 \Rightarrow k = \sqrt{5} \approx 2.23 \approx 3$$

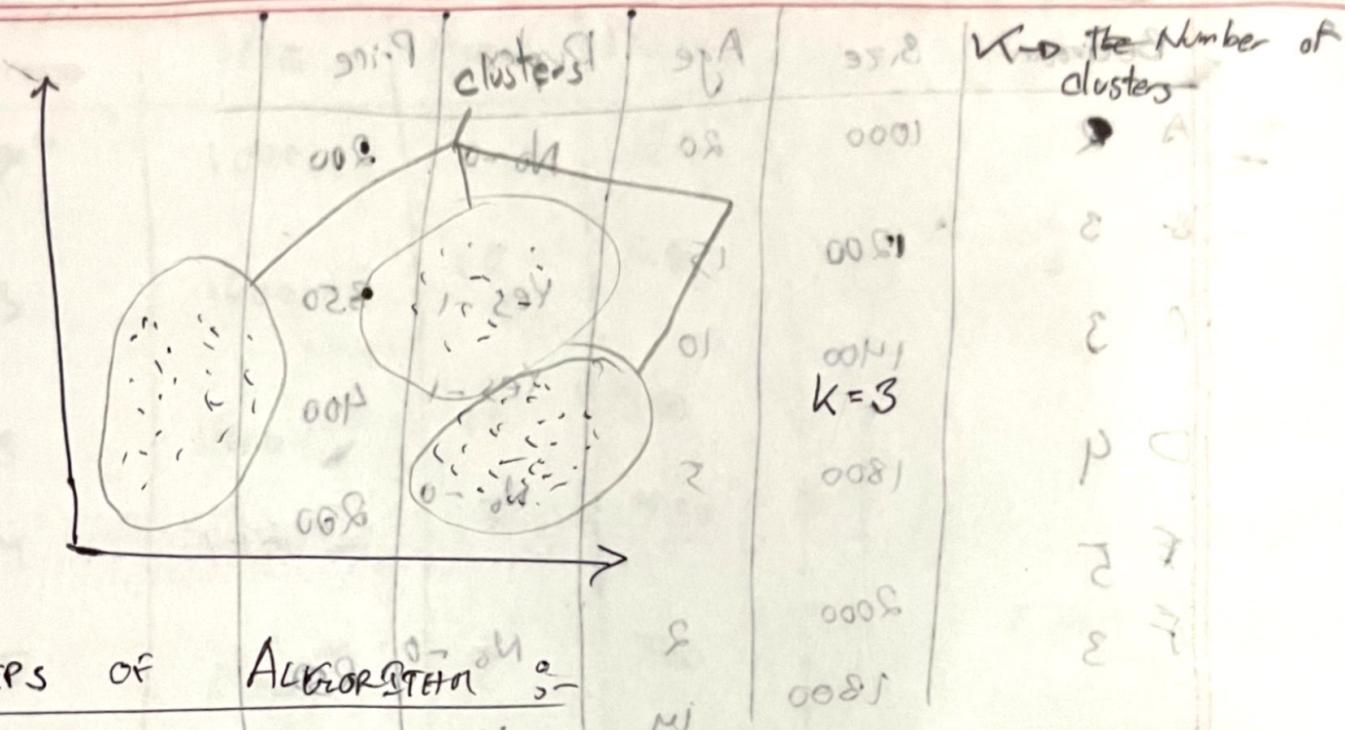
$$D(P, A) = \sqrt{(6.5 - 5)^2 + (3.8 - 3)^2 + (170 - 150)^2}$$

$$D(P, B) = \sqrt{(6.5 - 6)^2 + (3.8 - 3.5)^2 + (170 - 160)^2}$$

$$D(P, C) = \sqrt{(6.5 - 7)^2 + (3.8 - 4)^2 + (170 - 180)^2}$$

$$D(P, D) = \sqrt{(6.5 - 10)^2 + (3.8 - 6)^2 + (170 - 300)^2}$$

$$D(P, E) = \sqrt{(6.5 - 11)^2 + (3.8 - 6.5)^2 + (170 - 300)^2}$$

K-Means Clustering :-STEPS OF ALGORITHM :-

1. Initialize the ~~points~~ ^{centroids} randomly. (for each cluster)

2. Assign the ~~data points~~ to the nearest centroid.

~~Repeat~~ 3. Recalculating centroid.

When to stop :-

(i) The centroids are not changing

(ii) The clusters ...

(iii) If the centroids & clusters are changing continuously

Then run upto fixed number of iteration.

Example:-

Q-3 Solved

		x_1	x_2	$\underbrace{K=2}_{c_1}$	$c_1 \rightarrow (1, 2)$	$c_2 \rightarrow (10, 0)$	
$c_1 \rightarrow (1, 0)$		2	0	0	0	0	
A	1	1	0	2	2	2	
B	1	5	0	0	5	5	← calculate all the distances
C	10	1	2	10	8	8	← find out the shortest distances,
D	10	9	2	80.9	2	2	
$c_2 \rightarrow (10, 0)$		0	0	0	0	0	

Random centroid

$d_c(c_1, c_2) = \sqrt{(1-10)^2 + (2-0)^2} = 9.22$

Q. 2 (A) \therefore Centroid

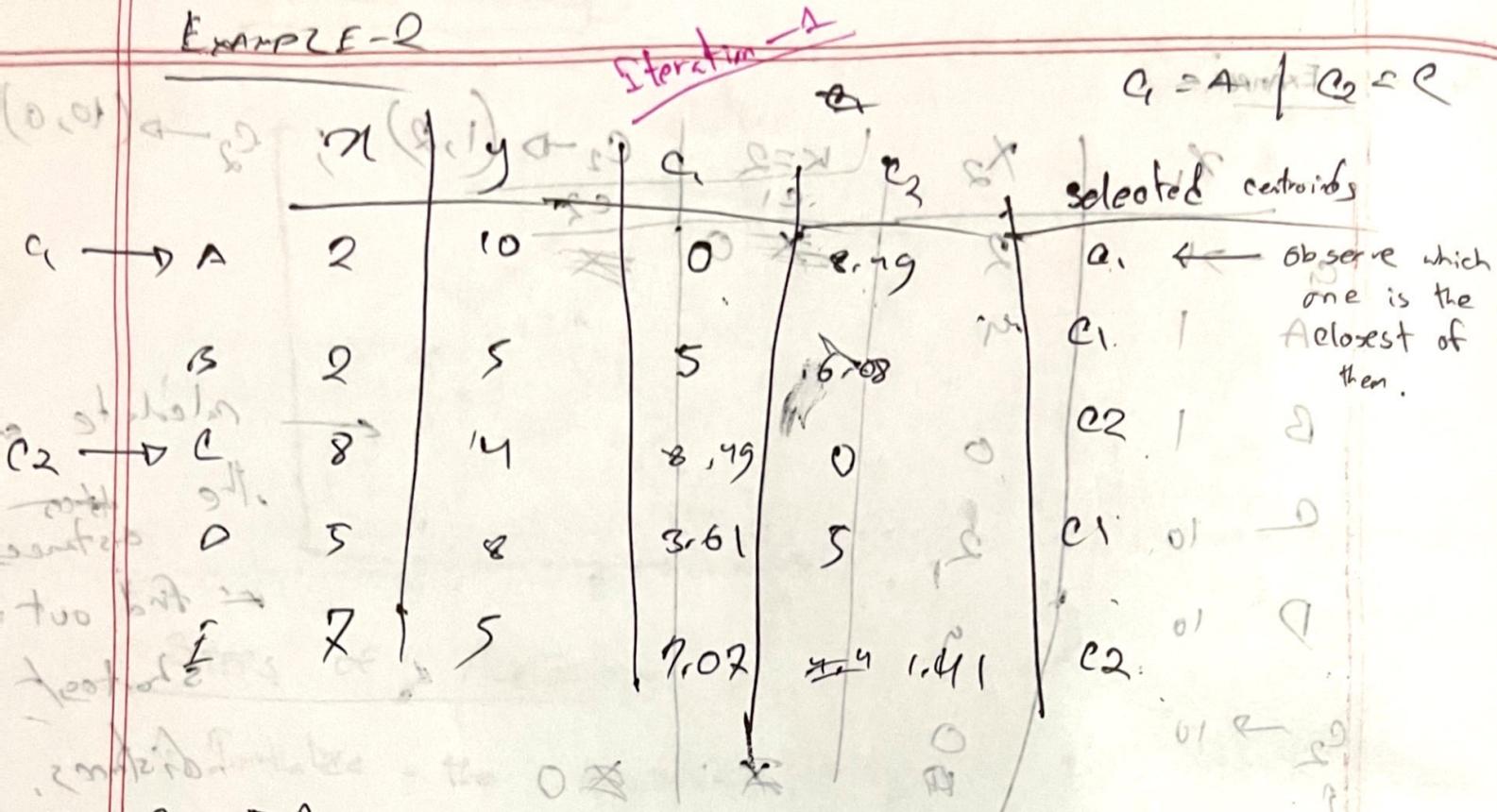
$$c_1(A, B) \mid c_2 = (C, D)$$

$$A = (1, 4) \mid B = (1, 0) \mid c_1 = (1, 2)$$

$$\frac{1+1+1}{3} = 1 \mid \frac{4+0+2}{3} = 2$$

\therefore Centroid $c_1 = (1, 2)$

Example-2



$$C_1 = A \quad C_2 = B$$

selected centroids

$C_1 \leftarrow$ observe which one is the closest of them.

$$C_2 \rightarrow B$$

$$C_1 \rightarrow A$$

$$C_2 \rightarrow B$$

$$C_1 \rightarrow A$$

Both are equal

$$C_2 \rightarrow B \Rightarrow \sqrt{(2-8)^2 + (0-4)^2} = \sqrt{32} = \sqrt{8+16} = \sqrt{24} = 4.9 \quad \leftarrow \text{distance between } (C_2, B) \text{ & } C_2$$

Cluster-1 :- A, B, D

$$(A, B) = 8 \quad | \quad (A, D) = 10$$

cluster 2 :- C, E

$$(A, C) = 8 \quad | \quad (A, E) = 10$$

From C_2 new cluster \Rightarrow

$$C_2 \Rightarrow \frac{2+8+5}{3} = \frac{15}{3} = 5 \quad | \quad \frac{10+5+8}{3} = \frac{23}{3} = 7.67$$

From C_2 new cluster \Rightarrow

$$C_2 \Rightarrow \frac{8+16}{2} = \frac{24}{2} = 12 \quad | \quad \frac{4+5}{2} = \frac{9}{2} = 4.5$$

then again find the distance

New points

$$\text{Initial center } C_1 \rightarrow (3, 7.67) \quad C_2 \rightarrow (7.5, 4.5)$$

Iteration 2

	x	y	C ₁	C ₂	DFE
A	2	10	2.45	9.98	$C_1 = (3, 7.67)$
B	2	5	2.83	5.5	$C_2 = (7.5, 4.5)$
C	8	9	6.01	0.71	
D	5	8	2.04	4.5	
E	8	5	4.072	0.21	
F	3	3	0.678	5.5	

cluster point
Cluster - 1 → A, D, E, F
Cluster - 2 → B, C

New feature ref 2

$$\text{Analogous to } \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \cdot (z_2 - z_1)$$

For F point →

$$d(F, C_1) \rightarrow \sqrt{(3-3)^2 + (7.67-3)^2}$$

$$d(F, C_2) \rightarrow 5.5$$

$$0 = d + (x)^T w \leq H$$

$$\frac{|d + (x)^T w + (x)^T b|}{\|(w)\|} = b$$

$$\frac{d + (x)^T w}{\|(w)\|} = (x)^T b$$

SUPPORT VECTOR MACHINE (Machine learning Model)

Concept:- SVM aims to find the optimal hyperplane that best separates data into different classes.

This hyperplane is a line [In two dimensions]

This hyperplane is a plane [In 3 D]

COMMON TERMS NEEDED:-

- (i) Margin :- It is the distance between the hyperplane and the closest data points from both the classes.
- These closest points are called support vectors.

Place the hyperplane such that the data points that lie closest to the hyperplane.

MATHEMATICAL TERMS

Equation of hyperplane:-

$$H \Rightarrow w^T(x) + b = 0$$

Distance from the hyperplane:-

$$d = \frac{|w^T(x_0) + b|}{\sqrt{w^2 + b^2}}$$

$$d_H(x_0) = \frac{w^T(x_0) + b}{\|w\|_2}$$

Place the hyperplane in such a position that from the positive samples and negative samples will be in the farthest possible distances.

In other words, we have to maximise the margin.

Why maximise margin is necessary?

→ If the point is on more distant then we will be more confident to say that their class labels.

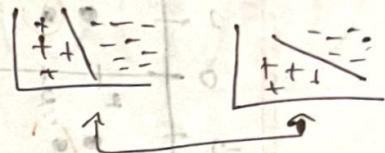
$\|w\|_2$ is the euclidean norm.

$$S = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$$

So we need to find the slope that can maximize the margin

Goal is to maximize the minimum distance

we need to find the slope that can maximize the margin



If we insert a point $p(x, y)$ on a hyperplane equation and we get

positive value then the point is from the positive group

$$w^T(x) + b > 0$$

And, if we get a negative value then the point

from the negative group

$$w^T(x) + b < 0$$

VERIFICATION

$$y_n \underbrace{[w^T \phi(x) + b]}_{\text{predicted}}$$

origin

The no. of support vectors reduced to try -

$$\rightarrow d+1 \quad \left(\begin{array}{c} d \\ 1 \end{array} \right) \quad [d = \text{dimension of data}]$$

+1

$$x \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \cdot +1 = +1$$

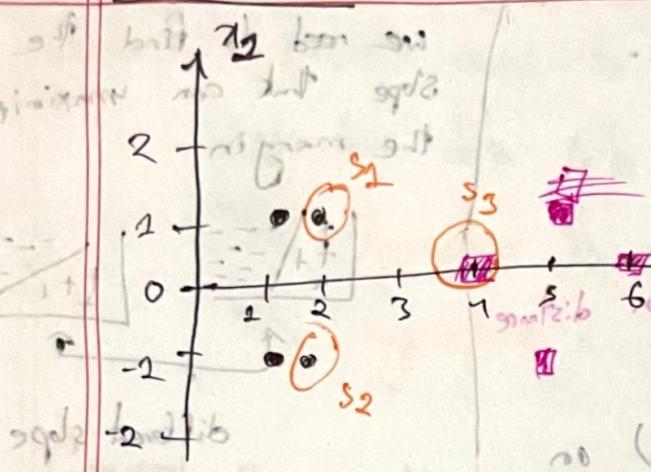
$$-1 \cdot -1 = +1$$

$$+1 \cdot -1 = -1$$

$$-1 \cdot +1 = -1$$

$$+1 \$$

EXAMPLE :-



Step-2

Dimension is $n_2, n_3 = 2$

Required support vectors = 2+1

$$S_1(2, 2); S_2(2, -1)$$

$$S_3(4, 0)$$

Step-3: Augment the support vectors. Given eqn

$$S_1 = (2, 1, 1), S_2 = (2, -1, 1), S_3 = (4, 0, 1)$$

$$0 < d + (x)^T w$$

Step-3: find the 3 parameters $\alpha_1, \alpha_2, \alpha_3$

$$\alpha_1 S_1 \cdot S_1 + \alpha_2 S_2 \cdot S_2 + \alpha_3 S_3 \cdot S_3 = -1 \quad (i)$$

$$\alpha_1 S_1 \cdot S_2 + \alpha_2 S_2 \cdot S_2 + \alpha_3 S_3 \cdot S_2 = -1 \quad (ii)$$

$$\alpha_1 S_1 \cdot S_3 + \alpha_2 S_2 \cdot S_3 + \alpha_3 S_3 \cdot S_3 = -1 \quad (iii)$$

From (i) substituting S_1, S_2, S_3

$$\alpha_1 \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 0 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 \cdot 6 + \alpha_2 \cdot 4 + \alpha_3 \cdot (-7) = -1 \quad (iv)$$

From (ii) we get, $\varepsilon = d$ $\left(\begin{array}{c} 1 \\ 0 \end{array}\right) = w$

$$d_1 \left(\begin{array}{c} ? \\ 1 \end{array}\right) \left(\begin{array}{c} 2 \\ 0 \end{array}\right) + d_2 \left(\begin{array}{c} 2 \\ 1 \end{array}\right) \left(\begin{array}{c} ? \\ 0 \end{array}\right) + d_3 \left(\begin{array}{c} 2 \\ 0 \end{array}\right) \left(\begin{array}{c} 2 \\ 1 \end{array}\right) = -1$$

$$d_1 \times 4 + d_2 \times 6 + [d_3] \times 9 = -1$$

\vdots

$$\varepsilon = d \times \left[\begin{array}{c} 4 \\ 6 \\ 9 \end{array}\right] = 0 \quad \text{--- (v)}$$

From (iii) we get,

$$(d_1 \times 2, d_2 \times \left(\begin{array}{c} ? \\ 1 \end{array}\right) \left(\begin{array}{c} 9 \\ 1 \end{array}\right) + d_2 \left(\begin{array}{c} 2 \\ 1 \end{array}\right) \left(\begin{array}{c} 9 \\ 0 \end{array}\right) + d_3 \left(\begin{array}{c} 4 \\ 0 \end{array}\right) \left(\begin{array}{c} 9 \\ 1 \end{array}\right) = 1$$

$$d_1 \times 9 + d_2 \times 9 + d_3 \times 17 = 1$$

Step - 4

$$\frac{d_2 \times 9}{d_2 \times 9} = -3.25 \quad / \quad d_3 = 3.5 \quad \varepsilon = \left[\begin{array}{c} 4 \\ 6 \\ 9 \end{array}\right] \quad \text{(vi)}$$

find the equation $d_1 \approx d_2, d_3$ using (iv), (v), (vi)

$$w = \sum_i d_i s_i$$

$$w = d_1 \left(\begin{array}{c} ? \\ 1 \end{array}\right) + d_2 \left(\begin{array}{c} 2 \\ 1 \end{array}\right) + d_3 \left(\begin{array}{c} 4 \\ 0 \end{array}\right)$$

$$w = -3.25 \left(\begin{array}{c} ? \\ 1 \end{array}\right) + \left(\begin{array}{c} 2 \\ 1 \end{array}\right) (-3.25) + (3.5) \left(\begin{array}{c} 4 \\ 0 \end{array}\right)$$

$$w^* = \left(\begin{array}{c} 1 \\ 0 \\ 3 \end{array} \right) \rightarrow w^*$$

$$\text{Given } \omega^* = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, b = -3 \quad (\text{by } y = \omega x + b) \text{ meant}$$

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} = 1 \cdot \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} \left(\begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} b = 3 \right) + \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \in \mathbb{R}$$

The equation is $x_1 \cdot [1 \ 2] + x_2 \cdot [0 \ 1] + b = 3$

$$y = \cancel{\begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}} \times x + -3$$

(equation of line)

Testing Hypothesis $\left(\begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} \in \mathbb{R} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} \in \mathbb{R} \right) \text{ by popline}$

Test point :- $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$

$$y = x_1 \cdot [1 \ 2] + x_2 \cdot [0 \ 1] + b$$

Hence,

$$(iv) \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1 \\ 2 \end{pmatrix} - 3 \xrightarrow{\text{M - opre}} 2.2 = \infty \quad | \quad 2.2 - 3 = \infty = 6 \in \mathbb{R}$$

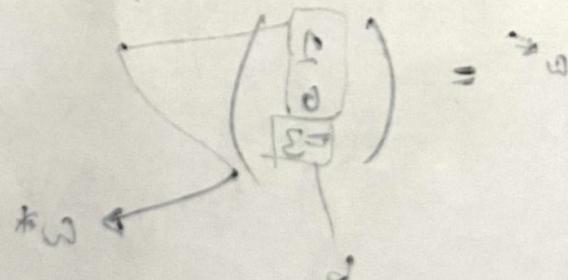
$$\begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1 \\ 2 \end{pmatrix} - 3$$

$$= 2 - 3$$

$$= -2 \quad (2 \neq 0 \neq -2 = 0)$$

$$\left(\begin{pmatrix} 1 \\ 0 \end{pmatrix} \in \mathbb{R} + \begin{pmatrix} 1 \\ -1 \end{pmatrix} \in \mathbb{R} \right) \times \begin{pmatrix} 1 \\ 2 \end{pmatrix} + b = 0$$

$$\left(\begin{pmatrix} 1 \\ 0 \end{pmatrix} (2.2) + (2.2) \begin{pmatrix} 1 \\ 2 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} 2.2 - 3 = 0 \right)$$



"INFERENTIAL STATISTICS"

- Inferential statistics allows you to make predictions/inferences from chart or graph. It allows you take data from samples and make generalizations about a population.
- Inferential statistics are concerned with making inferences based on the relations we found from the sample data.

ONE METHOD: Z score / Standard score
OR TESTS

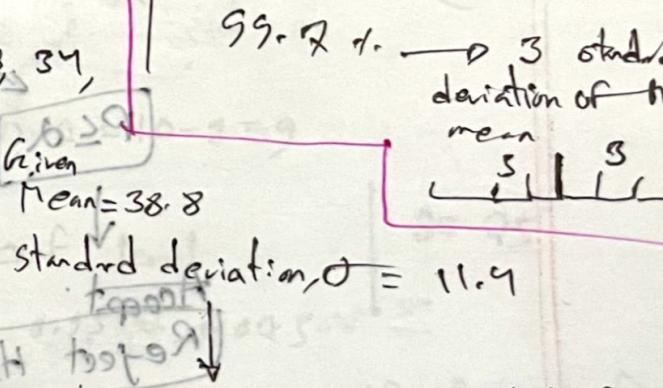
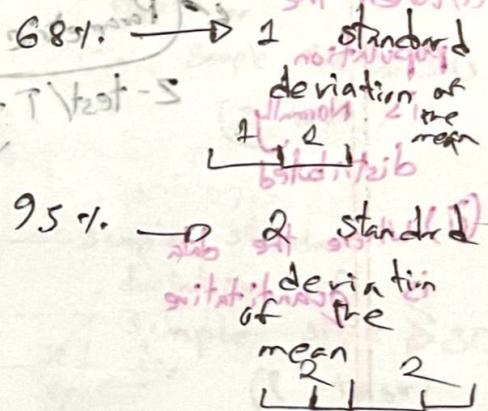
$$Z(\text{Score}) = \frac{\text{value} - \text{mean}(\bar{x})}{\text{Standard deviation } (\sigma)}$$

Sample standard deviation $\sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N-1}}$

EXAMPLE →

sample data: 26, 33, 65, 28, 34,
calculate z score!

Value	calc	Z-score
26	$\frac{26 - 38.8}{11.4}$	-1.12
33	$\frac{33 - 38.8}{11.4}$	-0.51
65	$\frac{65 - 38.8}{11.4}$	2.29
28	$\frac{28 - 38.8}{11.4}$	-0.91
34	$\frac{34 - 38.8}{11.4}$	-0.42



how much deviated from the mean.

HYPOTHESIS TESTING

"TESTING THE HYPOTHESIS"

STEPS IN HYPOTHESIS TESTING

NULL HYPOTHESIS (H_0)

Propose a base line statement

STATEMENT WHICH STATES THAT THERE IS NO RELATIONSHIP BETWEEN THE VARIABLES

ALTERNATE HYPOTHESIS (H_A)

STATEMENT THAT STATES THERE IS A RELATIONSHIP BETWEEN THE VARIABLES

HYPOTHESIS TESTING

(i) When the population is Normally distributed

(ii) Where the data is quantitative

Parametric Z-test/T test

NON PARAMETRIC Chi squared Test

(i) Where the population is not Normally distributed

(ii) Where the data is qualitative.

Compute P value

$$P \leq \alpha$$

Accept

Reject H_0

$$P > \alpha$$

Accept

Reject H_0

CONFIDENCE INTERVAL VS CONFIDENCE LEVEL

(U.S. Census Bureau) use: 90% confidence level

→ If Bureau Census Bureau repeats the survey using the same technique's 90% of the time the result would fall between 35,534 to 37,315 people in poverty.

→ 35,534 to 37,315 is the confidence interval.

$$\text{Confidence Interval} : \bar{x} \pm Z \frac{s}{\sqrt{n}}$$

standard deviation
Population
For Sample size ≥ 30
(Z test)

[z-score]

$$\text{Confidence Interval} : \bar{x} \pm t \frac{s}{\sqrt{n}}$$

sample standard deviation of sample size < 30
(t test)

(t score)

COMPUTATION OF T-Score

$$\text{def degrees of freedom, } df = n-2 = 10-2 = 8$$

$$\text{Significance level, } \alpha = 1 - \text{conf level}$$

$$S = 25$$

COMPUTING CONFIDENCE INTERVAL

$$\text{Standard Error, } SE = \frac{s}{\sqrt{n}}$$

$$\begin{aligned} \text{Upper Confidence interval} &= 240 + 2.62 \times \frac{25}{\sqrt{10}} \\ &= 259.883 \end{aligned}$$

Using df & significance level find the ~~t~~ t-score from T-table

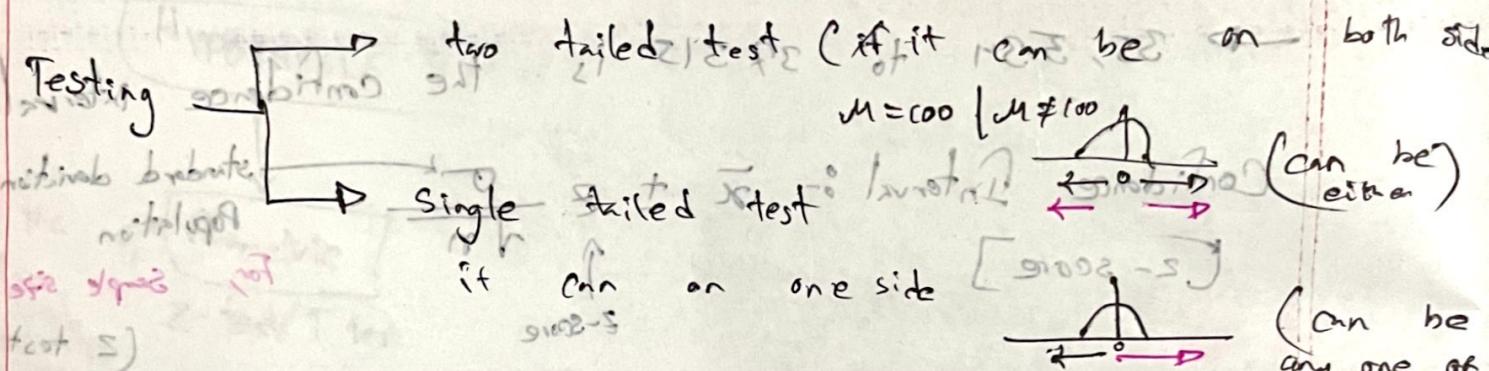
$$t_{0.025} = 2.0$$

$$\text{Lower Confidence Interval} = \bar{x} - t \times \frac{s}{\sqrt{n}}$$

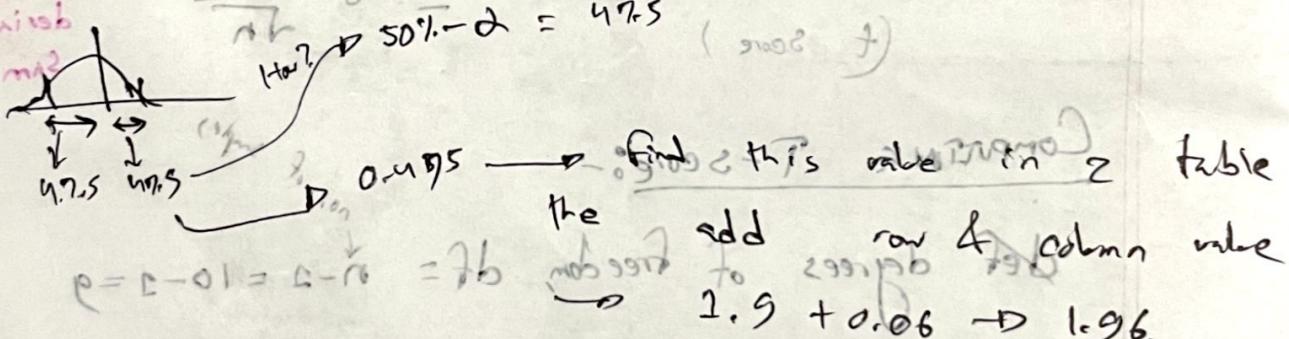
$$= 240 - \frac{2.62}{\sqrt{10}} \times \frac{25}{\sqrt{10}} \rightarrow 214$$

$$= 222.117$$

Confidence Interval: 257.88 + 222.117 between 480.00 & 218.88



Computing Z-score from table:



$Z = 2$ | Level $\alpha = 0.05$ = level committing

Z-test statistics of $\bar{x} - \mu_0$ which will be null hypothesis

$$Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

where $\bar{x} = 240$, $\mu_0 = 257.88$, $\sigma = 25$, $n = 10$

$$Z = \frac{240 - 257.88}{25/\sqrt{10}} = -3.88$$

Chi Square Test

chi squared statistics

FORMULA :-

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

→ use this value (to) find a value from
the table.

Conditions

$\text{SSP} > 10$ not significant

$P \leq .1$ marginally

$P \leq 0.5$ significant

$P \leq .01$ highly significant.

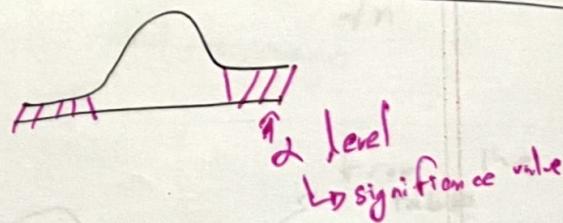
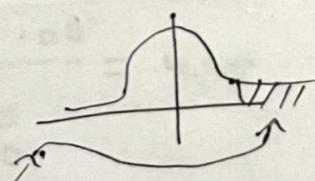
If this value is not present on
the table then find the range!

→ then (are) get P value

[null hypothesis can't be rejected]

$$\begin{aligned} \text{Expected} &= \frac{\text{Total artist}}{\text{no. of category}} \\ (\text{For } 2 \text{ category}) &= \frac{256}{12} \\ &= 21.33 \end{aligned}$$

$$\begin{aligned} \text{For 2 category} \\ \text{Expected} &= \frac{\text{row total} \times \text{column total}}{\text{All total}} \end{aligned}$$



→ If z-score from test statistics is greater than
rejection area then reject the null hypothesis.

$$T\text{-test statistics} \rightarrow t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

$$6 - x_1 + (x_1 - x_c)^2$$

carrying capacity step

test sample id

Non Linear SVM

STEP-2 → Transform

$$\phi(x_1, x_2) = \begin{cases} \left(6 - x_1 + (x_1 - x_c)^2\right) & \text{if } \sqrt{x_1^2 + x_2^2} \geq 2 \\ \left(6 - x_2 + (x_1 - x_c)^2\right) & \text{otherwise} \end{cases}$$

TRANSFORM

blue class → $\left(\frac{1}{2}\right)$ (if $\sqrt{x_1^2 + x_2^2} < 2$)

Red class → $\left(\begin{array}{c} 2 \\ 0 \end{array}\right)$ (if $\sqrt{x_1^2 + x_2^2} \geq 2$)

$\left[\begin{array}{c} 2 \\ 0 \end{array}\right] \rightarrow \phi\left(\begin{array}{c} 2 \\ 0 \end{array}\right) = \left(\begin{array}{c} 8 \\ 10 \end{array}\right)$ (if $\sqrt{x_1^2 + x_2^2} \geq 2$)

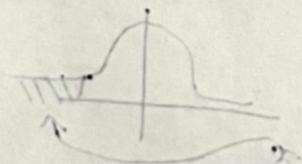
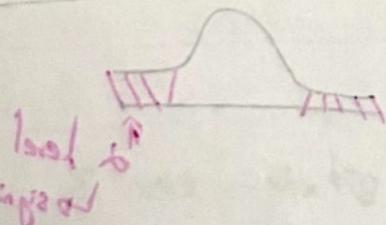
$\left[\begin{array}{c} 0 \\ 2 \end{array}\right] \rightarrow \phi\left(\begin{array}{c} 0 \\ 2 \end{array}\right) = \left(\begin{array}{c} 10 \\ 8 \end{array}\right)$ (if $\sqrt{x_1^2 + x_2^2} < 2$)

STEP-2

$$\text{SVM} = \frac{\sum x_i y_i}{n}$$

no tracing tan xi solve diff if
larger off bind next point exist

(transformed)



next step is calculate test most $\bar{x} - \mu_0 = 5$ PP
classified 1st or higher right not foreign

$$\frac{\bar{x} - \mu_0}{\sigma} \Rightarrow \text{calculate test-T}$$