

⇒ Correlation Coefficient (Pearson's Product moment Coefficient)

$$\gamma_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n \sqrt{V_A} \sqrt{V_B}}$$

$$V_A = \sqrt{\frac{\sum (A - \bar{A})^2}{n}}, \quad V_B = \sqrt{\frac{\sum (B - \bar{B})^2}{n}}$$

⇒ Covariance (Numeric Data)

$$\text{Cov}(A, B) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

$$\gamma_{A,B} = \frac{\text{Cov}(A, B)}{\sqrt{V_A} \sqrt{V_B}}$$

⇒ Correlation Analysis (Nominal Data)

(Chi-square ( $\chi^2$ ) test)

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

⇒ Normalization (Min-Max Normalization)

$$v'_i = \frac{v_i - \text{min}_A}{\text{Max}_A - \text{Min}_A} (\text{newmax}_A - \text{newmin}_A) + \text{newmin}_A$$

⇒ Z-score Normalization

$$v' = \frac{v - \mu_A}{\sigma_A}$$

$\mu_A$  = mean of A

$\sigma_A$  = standard deviation

⇒ Normalization By Decimal Scaling

$$v' = \frac{v}{10^j}$$

$$\Rightarrow (\text{Max}|v'|) < 1$$

$$(10^2)$$

$$\left\{ \begin{array}{l} w_1 = 99 \\ w_2 = 97 \end{array} \right.$$

- Binning :- By Equiwidth Partitioning  
 width =  $\frac{\max - \min}{K}$   $K \rightarrow$  intervals of same size  
 $K=3$
- By Equidepth Partitioning → Number of values in each bin  $K=3$   
Equidepth Smoothing by Bin Means :- Replace each value in bin with mean of that bin
- Eqwideth (by Median) :- Replace each value with the middle value of given range.
- Eqwideth by Boundaries :- Replace each value in bin with nearest boundary value.

### → Mean (Measuring Central tendency)

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow \frac{x_1 + x_2 + \dots + x_n}{n}$$

### → Weighted Arithmetic Mean

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

→ Trimmed Mean :- Mean obtained after chopping off values at high and low extremes.

→ Median :- Middle value if odd number of values.

→ Average of middle two values if even number of values.

### → For grouped data :-

$$\text{Median} = L_1 + \left( \frac{\frac{n}{2} - (\sum \text{freq})_L}{\text{freq median}} \right) \times \text{width}$$

$L_1 \rightarrow$  lower boundary of median interval

$n \rightarrow$  Number of values in dataset.

$\sum \text{freq}_L \rightarrow$  Sum of frequencies of all intervals that are lower than the median interval.

$\text{freq median} \rightarrow$  freq of median interval, width  $\rightarrow$  width of median interval.

$\Rightarrow$  Mode: - Value that occurs most frequently in data.

For unimodal

$$\text{mean-mode} = 3 \times (\text{mean}-\text{median})$$

$\Rightarrow$  Range: - Difference between Max and min values in set of observations

$\Rightarrow$  Quartiles :-  $Q_1 = 25^{\text{th}} \text{ percentile}$ ,  $Q_3 (75^{\text{th}} \text{ percentile})$

$Q_2 \rightarrow 50^{\text{th}} \text{ percentile (Median)}$

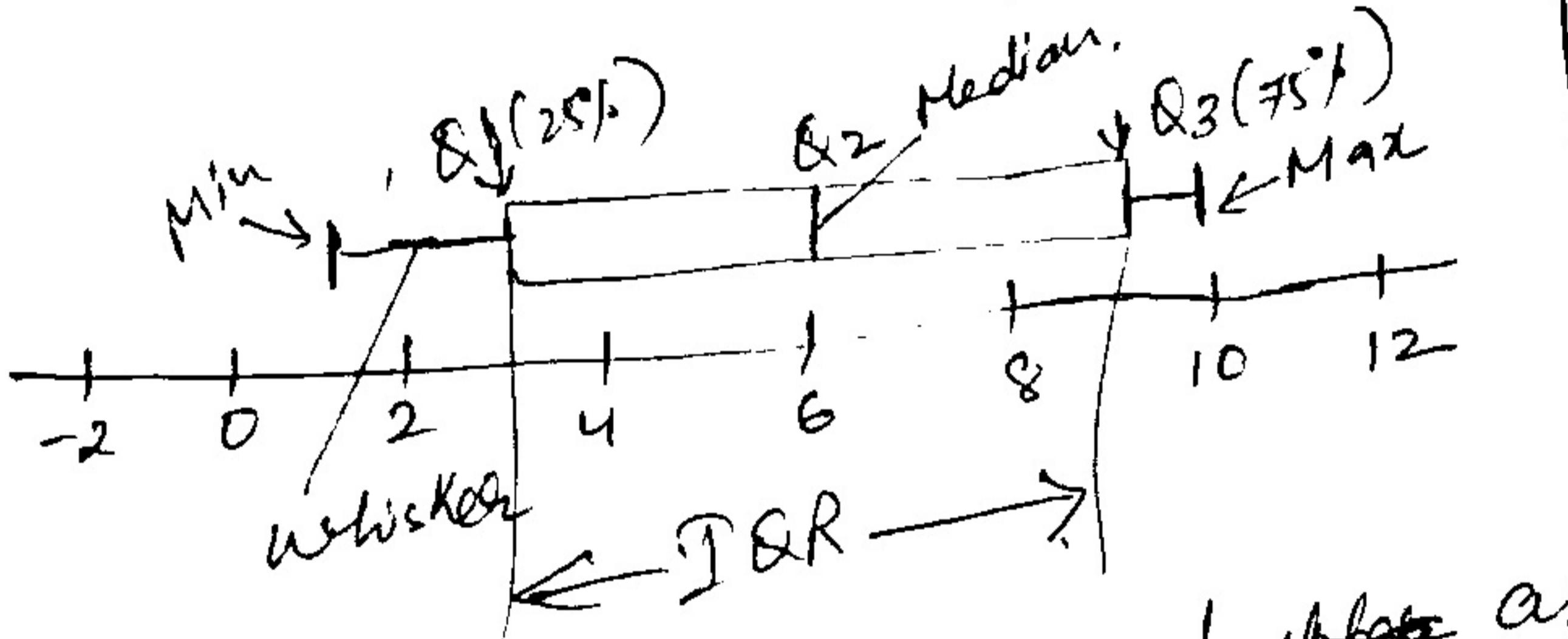
$\Rightarrow$  Interquartile Range :-  $IQR = Q_3 - Q_1$

$\Rightarrow$  Five Number Summary :- min,  $Q_1$ , median,  $Q_3$ , max

$$Q_1 = 25^{\text{th}} \text{ percentile} = \frac{25}{100}(n-1) \text{ th term element}$$

$$Q_3 = 75^{\text{th}} \text{ percentile} = \frac{75}{100}(n-1) \text{ th term element}$$

$\Rightarrow$  Boxplot



height of box = IQR

Outliers  $\rightarrow$  data points  $\rightarrow$  below and ~~above~~ above the lower

and upper limit - lower =  $Q_1 - 1.5 \times IQR$

upper =  $Q_3 + 1.5 \times IQR$

$\Rightarrow$  Variance  $\sigma^2 = V^2 = \frac{1}{n} \sum_{i=1}^N (x_i - \bar{x})^2$   $[\bar{x} \rightarrow \text{mean}]$

$\Rightarrow$  Standard Deviation :-  $\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^N (x_i - \bar{x})^2}$

$$\sigma = \sqrt{V}$$

$\Rightarrow$  Dissimilarity Matrix

|          |          |   |
|----------|----------|---|
| 0        |          |   |
| $d(2,1)$ | 0        |   |
| $d(3,1)$ | $d(3,2)$ | 0 |
| !        | !        | ! |
| $d(n,1)$ | $d(n,2)$ | 0 |

$$d(i,j) = d(j,i)$$

$$\text{sim}(i,j) = 1 - d(i,j)$$

## Proximity measure of nominal attributes

$$d(i,i) = \frac{P-m}{P} \Rightarrow P \rightarrow \text{Total Number of attributes}$$

$M \rightarrow \text{No. of matches}$

| Id | Type 1 | Type 2 |
|----|--------|--------|
| 1  | 10     | A      |
| 2  | 30     | B      |
| 3  | 10     | A      |
| 4  | 20     | C      |

$$d(2,1) = \frac{2-0}{2} = 1$$

$$d(3,1) = \frac{2-2}{0} = 0$$

$$d(4,1) = \frac{2-0}{2} = 1$$

$$d(3,2) = \frac{2-0}{2} = 1$$

$$d(4,2) = \frac{2-0}{2} = 1$$

$$d(4,3) = \frac{2-0}{2} = 1$$

dissimilarity matrix =  $\begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix}$

$d(3,1)$  are similar.  
3,1

$$\text{Sim}(i,j) = 1 - d(i,j) = \frac{m}{P}$$

attributes

| Name | T <sub>1</sub> | T <sub>2</sub> | T <sub>3</sub> | T <sub>4</sub> | T <sub>5</sub> | T <sub>6</sub> |
|------|----------------|----------------|----------------|----------------|----------------|----------------|
| X    | 1              | 0              | 1              | 0              | 1              | 0              |
| Y    | 1              | 0              | 1              | 0              | 0              | 1              |
| Z    | 1              | 1              | 0              | 0              | 0              | 0              |

| Object j |     |
|----------|-----|
| Object i |     |
|          | 1 0 |
| 1        | q Y |
| 0        | s t |

$$q \rightarrow \{i=1, n=1\}$$

$$q \rightarrow \{j=1, j=0\}$$

$$s \rightarrow \{i=0, j=1\}$$

$$t \rightarrow \{i=0, j=0\}$$

$$\text{Coherence}(i,j) = \frac{q}{(q+r)+(q+s)-q}$$

$$i \rightarrow x, j \rightarrow y$$

$$q \rightarrow 2, r \rightarrow 0, s \rightarrow 1$$

$$d(X,Y) = \frac{q+s}{q+r+s} = \frac{0+1}{2+0+1} = 0.33$$

$$d(X,Z) = 0.67$$

$$d(Y,Z) = 0.75$$

$d(X,Y)$  are more similar

→ Minkowski Distance :-

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

where  $i = (x_{i1}, x_{i2}, \dots, x_{ip})$  and  $j = (x_{j1}, x_{j2}, \dots, x_{jp})$  are two  $p$ -dimensional data objects. ' $h$ ' is the order (the distance so defined is also called  $h$ -norm)

→  $d(i, j) > 0$  if  $i \neq j$  and  $d(i, i) = 0$

→  $d(i, j) = d(j, i) \rightarrow$  Symmetry

→  $d(i, j) \leq d(i, k) + d(k, j)$

,  $\therefore h=1$  :- Manhattan (City block,  $L_1$  norm) distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

$h=2$  :- Euclidean distance

$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

$h \rightarrow \infty$  :- Supremum ( $L_{\max}$  norm,  $L_{\infty}$  norm) → Max difference

$$d(i, j) = \lim_{h \rightarrow \infty} \left( \sum_{p=1}^P |x_{ip} - x_{jp}|^h \right)^{1/h}$$

$$= \max_p (x_{ip} - x_{jp})$$

Example

| Point | attribute1 | attribute2 |
|-------|------------|------------|
| $x_1$ | 1          | 2          |
| $x_2$ | 3          | 5          |
| $x_3$ | 2          | 0          |
| $x_4$ | 4          | 6          |

Manhattan  
( $L_1$ )

| $L_1$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ |
|-------|-------|-------|-------|-------|
| $x_1$ | 0     |       |       |       |
| $x_2$ | 5     | 6     | 0     |       |
| $x_3$ | 3     | 6     | 1     | 7     |
| $x_4$ | 6     | 1     | 7     | 0     |

Euclidean ( $L_2$ )

| $L_2$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ |
|-------|-------|-------|-------|-------|
| $x_1$ | 0     |       |       |       |
| $x_2$ | 3.61  | 0     |       |       |
| $x_3$ | 2.24  | 5.1   | 0     |       |
| $x_4$ | 4.24  | 1     | 5.39  | 0     |

Supremum

| $L_{\infty}$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ |
|--------------|-------|-------|-------|-------|
| $x_1$        | 0     |       |       |       |
| $x_2$        | 3     | 0     |       |       |
| $x_3$        | 2     | 5     | 0     |       |
| $x_4$        | 3     | 1     | 5     | 0     |

## $\Rightarrow$ Proximity Measure of ordinal attributes

| Obj | T <sub>1</sub> |
|-----|----------------|
| 1   | High           |
| 2   | low            |
| 3   | medium         |
| 4   | high           |

High  $\rightarrow 1$   
 Medium  $\rightarrow 2$   
 low  $\rightarrow 3$   
 Normalize rank  $[0, 1]$

$$\rightarrow M_f \rightarrow 3$$

for high =  $\frac{1-1}{3-1} = 0/3 = 0$

$$Z_{if} = \frac{x_{if}-1}{M_f-1}$$

$$1 \rightarrow 0.0, 2 \rightarrow 0.5, 3 \rightarrow 1.0$$

dissimilarity Matrix  $d(2,1) = (1-0)=1, d(3,1) = (0.5-0)=0.5$   
 $d(4,1) = (0-0)=0, d(3,2) = (0.5-0.5)=0.5$

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1 & 0.5 & 0 \end{bmatrix}$$

$$d(4,2) = (1-0)=1$$

$$d(4,3) = (0.5-0)=0.5$$

$\Rightarrow$  Proximity Measure of Mixed attributes :-  
 A database may contain all attribute types of nominal, symmetric  
 binary asymmetric binary numeric, ordinal.

Combined  $\rightarrow d(i,j) = \frac{\sum_{f=1}^P \delta_{if}(f) d_{if}(f)}{\sum_{f=1}^P \delta_{if}(f)}$   $\rightarrow f \rightarrow$  is binary or nominal  
 $d_{if}(f) = 0$  if  $x_{if} = x_{jf}$   
 $\text{and } d_{if}(f) = 1 \text{ otherwise}$

$\rightarrow f$  is numeric  $\rightarrow$  use the normalized  
 distance

$\rightarrow f$  is ordinal :-

compute ranks  $\delta_{if}$  and  
 Treat  $Z_{if}$  as interval-scaled

$$Z_{if} = \frac{x_{if}-1}{M_f-1}$$

Example :- Based on the information given, find most similar and most dissimilar persons among them. Apply min-max normalization on income to obtain [0, 1] range. Consider profession and mother tongue as nominal and consider native place as ordinal variable with ranking order of [Village, small town, Suburban, Metropolitan]. Give equal weight to each attribute.

| Name   | Income | Profession    | Mother tongue | Native place |
|--------|--------|---------------|---------------|--------------|
| Ram    | 70000  | Doctor        | Bengali       | Village      |
| Balram | 50000  | DataScientist | Hindi         | Smalltown    |
| Bharat | 60000  | Carpenter     | Hindi         | Suburban     |
| Kishan | 80000  | Doctor        | Bhojpuri      | Metropolitan |

→ After normalizing and quantifying Native place

| Name   | Income | Profession    | Mother tongue | Native |
|--------|--------|---------------|---------------|--------|
| Ram    | 0.67   | Doctor        | Bengali       | 1      |
| Balram | 0      | DataScientist | Hindi         | 2      |
| Bharat | 0.33   | Carpenter     | Hindi         | 3      |
| Kishan | 1      | Doctor        | Bhojpuri      | 4      |

$$d(Ram, Balram) = \frac{0.67 + 1 + 1 + (2-1)}{4-1} = 3, \quad d(Ram, Bharat) = 0.33 + 1 + 1 + (3-1) = 3 \quad \frac{4-1}{4-1}$$

$$d(Ram, Kishan) = 0.33 + 0 + 1 + \frac{(4-1)}{4-1} = 2.33,$$

$$d(Balram, Kishan) = 1 + 1 + 1 + \frac{(4-2)}{4-1} = 3.67$$

Most Similar → Balram and Bharat

Most dissimilar → Balram and Kishan

$$\Rightarrow \text{Cosine Similarity} : - \quad \text{Cos}(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|}$$

$$\text{Ex:- } d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$$

$$d_1 \cdot d_2 = 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 0 + 2 \times 1 + 0 \times 0 + 0 \times 1$$

$$\|d_1\| = \sqrt{5^2 + 0^2 + 3^2 + 0^2 + 2^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2} = 6.481$$

$$\|d_2\| = \sqrt{3^2 + 0^2 + 2^2 + 0^2 + 1^2 + 1^2 + 0^2 + 1^2 + 0^2 + 1^2} = 4.12$$

$$\text{Cos}(d_1, d_2) = 0.94$$

## Classification

### GINI INDEX

$$\text{At any node } t, GINI(t) = 1 - \sum_j [P(j|t)]^2$$

$P(j|t)$  → relative frequency of class  $j$  at node  $t$ .

→ Maximum ( $1 - 1/n_c$ ) when records are equally distributed among all classes implying least interesting information.

→ Minimum ( $0, 0$ ) when all records belong to one class implying most interesting information.

|                |    |
|----------------|----|
| C <sub>1</sub> | 10 |
| C <sub>2</sub> | 6  |

$$GINI = 1 - \sum_j [P(j|t)]^2$$

$$\Rightarrow P(C_1) = 0\% = 0, P(C_2) = 6\% = 1$$

$$GINI = 1 - (0)^2 + (1)^2 \\ \Rightarrow 1 - 0 = 1$$

|                |   |
|----------------|---|
| C <sub>1</sub> | 1 |
| C <sub>2</sub> | 5 |

$$P(1) = \frac{1}{6}, P(2) = \frac{5}{6}$$

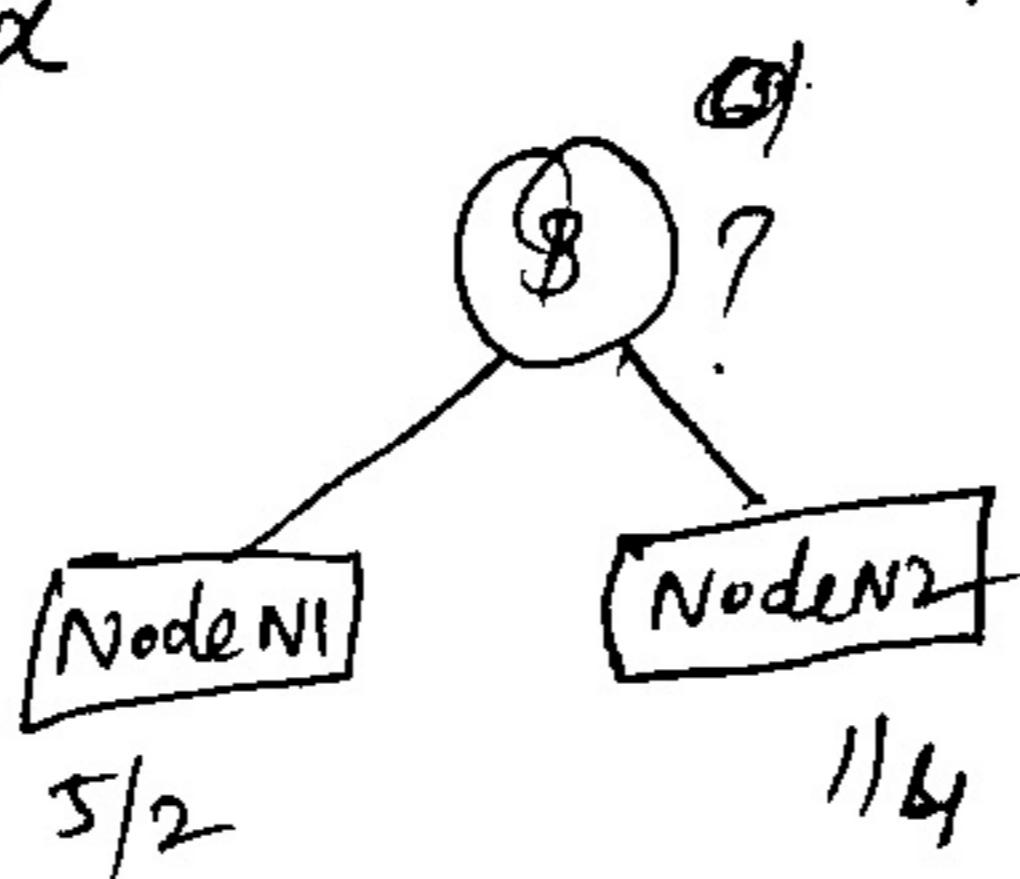
$$GINI = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{5}{6}\right)^2 = 0.278$$

### SPLITTING Based on GINI

$$GINI_{SPLIT} = \sum_{i=1}^K \frac{n_i}{n} (GINI(i)) \rightarrow \text{when a node } P \text{ is split into } K \text{ partitions (children)}$$

$\rightarrow n_i \rightarrow \text{number of records at child } i$   
 $\rightarrow n \rightarrow \text{number of records at node } P$

### Computing Gini Index



|                | Parent |
|----------------|--------|
| C <sub>1</sub> | 6      |
| C <sub>2</sub> | 6      |
| $GINI = 0.500$ |        |

$$GINI(N_1) = 1 - \left(\frac{5}{7}\right)^2 - \left(\frac{2}{7}\right)^2 = 0.408$$

$$GINI(N_2) = 1 - \left(\frac{11}{16}\right)^2 - \left(\frac{5}{16}\right)^2 = 0.32$$

|                | N <sub>1</sub> | N <sub>2</sub> |
|----------------|----------------|----------------|
| C <sub>1</sub> | 5              | 1              |
| C <sub>2</sub> | 9              | 4              |
| $GINI = 0.37$  |                |                |

$$GINI(\text{children}) = \frac{7}{12} \times 0.408 + \frac{5}{12} \times 0.32 = 0.37$$

⇒ Example 1:- Which attribute should be chosen as first splitting attribute (b/w  $q_1$  &  $q_2$ ) according to GINI INDEX?

Sol:- Contingency table for  $q_1$  &  $q_2$

| Class | $q_1 = T$ | $q_1 = F$ |
|-------|-----------|-----------|
| +     | 3         | 1         |
| -     | 1         | 4         |

| Class | $q_2 = T$ | $q_2 = F$ |
|-------|-----------|-----------|
| +     | 2         | 2         |
| -     | 3         | 2         |

| Instance | $q_1$ | $q_2$ | $q_3$ | target class |
|----------|-------|-------|-------|--------------|
| 1        | T     | T     | 1.0   | +            |
| 2        | T     | T     | 6.0   | +            |
| 3        | T     | F     | 5.0   | -            |
| 4        | F     | P     | 4.0   | +            |
| 5        | F     | T     | 7.0   | -            |
| 6        | F     | T     | 3.0   | -            |
| 7        | F     | F     | 8.0   | -            |
| 8        | T     | F     | 7.0   | +            |
| 9        | F     | T     | 5.0   | -            |

$$\text{Overall Gini} = 1 - \sum_{i=1}^2 [P(i/t)]^2$$

$$\text{Before splitting} = 1 - \left(\frac{4}{9}\right)^2 - \left(\frac{5}{9}\right)^2 = 1 - 0.197 = 0.495$$

$$Gini(\text{Split}) = \sum_{i=1}^K \frac{n_i}{n} Gini(i)$$

$$\text{For } q_1 : Gini_{\text{split}} \circ Gini(q_1) = \frac{4}{9} \left[ 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \right] + \frac{5}{9} \left[ 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{4}{5}\right)^2 \right]$$

$$Gini(q_1) = 0.3441$$

$$\Delta Gini(q_1) = 0.495 - 0.3441 = 0.151$$

$$\text{For } q_2 : Gini(q_2) = \frac{5}{9} Gini(i_1) + \frac{4}{9} Gini(i_2)$$

$$= \frac{5}{9} \left( 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 \right) + \frac{4}{9} \left[ 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \right]$$

$$= 0.4884$$

The attribute that maximises the reduction in impurity (or has minimum Gini index) is selected as splitting attribute.

$$\Delta Gini(q_2) = 0.495 - 0.4884 = 0.007$$

Hence,  $q_1$ , will be selected as the splitting attribute.

## $\Rightarrow$ Classification Error [Splitting Criteria]

Classification Error at node  $t$  :-  $i = \text{class}$

$$\text{Error}(t) = 1 - \max\{P(i|t)\}$$

Error before splitting :-

$$\text{Error Orig} : - 1 - \max\left(\frac{4}{9}, \frac{5}{9}\right) = 1 - \max(0.444, 0.555) \\ = 1 - 0.555 = 0.444.$$

Error after splitting on  $q_1$

$$\begin{aligned} \text{Error}_{q_1} &= 1 - \frac{4}{9} \left[ 1 - \max\left(\frac{3}{4}, \frac{1}{4}\right) \right] - \frac{5}{9} \left[ 1 - \max\left(\frac{1}{5}, \frac{4}{5}\right) \right] \\ &= 1 - \frac{4}{9} \left[ 1 - \max(0.75, 0.25) \right] - \frac{5}{9} \left[ 1 - \max(0.2, 0.8) \right] \\ &\Rightarrow 1 - 0.444(1 - 0.75) - 0.555(1 - 0.8) \\ &\Rightarrow 1 - 0.111 - 0.111 \\ &\Rightarrow 0.778. \end{aligned}$$

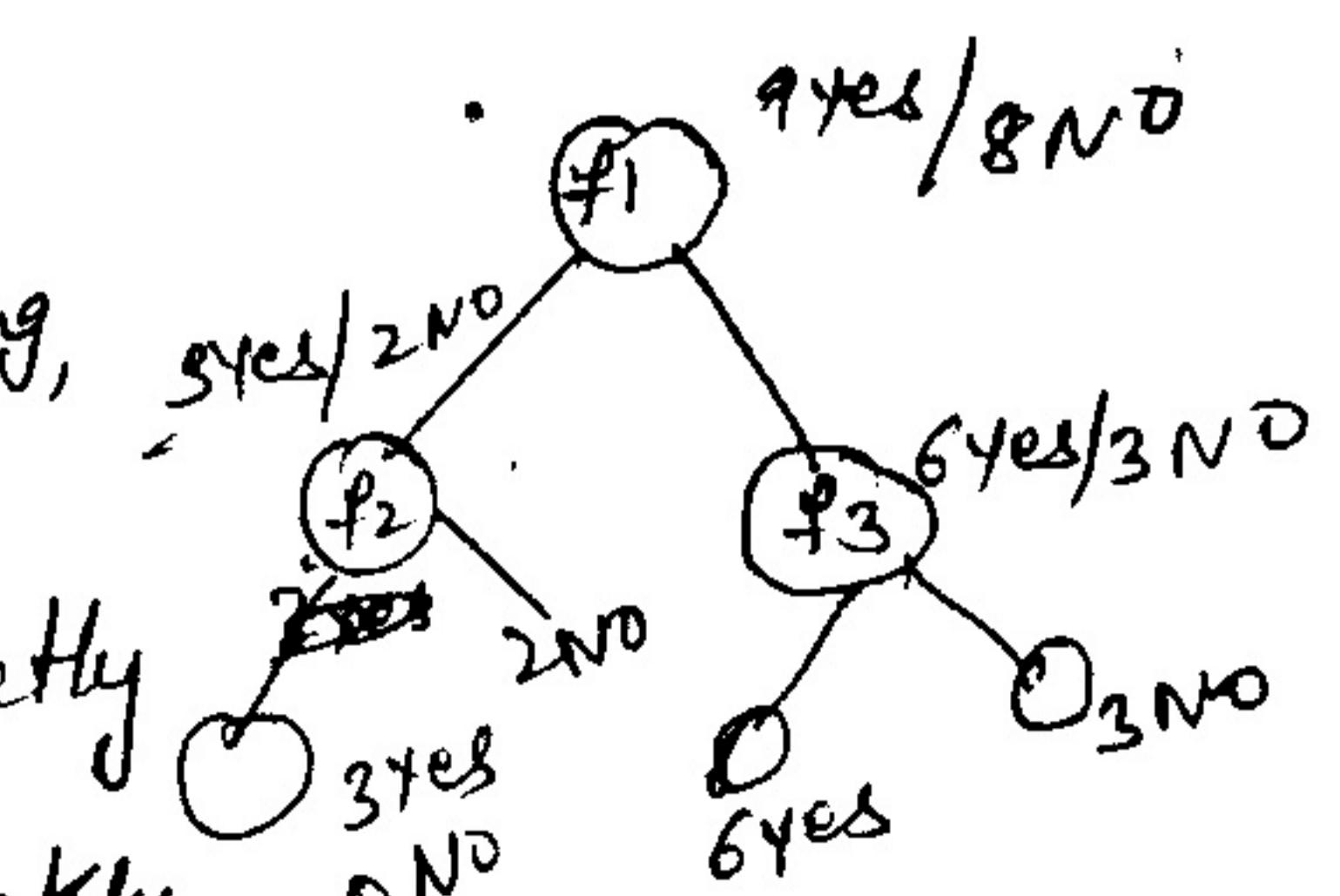
Error after splitting on  $q_2$

$$\begin{aligned} \text{Error}_{q_2} &= 1 - \frac{5}{9} \left[ 1 - \max\left(\frac{2}{5}, \frac{3}{5}\right) \right] - \frac{4}{9} \left[ 1 - \max\left(\frac{2}{4}, \frac{2}{4}\right) \right] \\ &= 1 - 0.555 \left[ 1 - \max(0.4, 0.6) \right] - 0.444 \left[ 1 - \max(0.5, 0.5) \right] \\ &= 1 - 0.222 - 0.222 = 0.556 \end{aligned}$$

The attribute that minimises classification error will be selected as splitting attribute, hence  $q_2$  will be selected as splitting attribute.

## $\Rightarrow$ ENTROPY :-

To select right attribute/feature for splitting, we need to calculate Entropy. How can we select the right node, if we select correctly we can reach leaf node quickly



Suppose features  $f_1, f_2, f_3$

$\Rightarrow$  Entropy helps to measure the purity of split.

$\Rightarrow$  Pure split whenever we get, that becomes leaf node.

O/P - Yes  
No

$$\boxed{\text{Entropy } H(S) = -P(+)\log_2 P(+) - P(-)\log_2 P(-)}$$

$P_+/P_-$  = % of +ve class / % of -ve class

$S$  = subset of training Example

$$\boxed{\text{Entropy}(+) = - \sum_j P(j/t) \log_2 P(j/t)}$$

→ when we have complete impure split like 3 Yes and 3 No, the entropy will be 1 (worst). Entropy will range from 0 to 1, hence the entropy value, better the split.

$$\Rightarrow \text{Ex!} - \begin{array}{|c|c|} \hline C_1 & 0 \\ \hline C_2 & 6 \\ \hline \end{array} \quad P(C_1) = -\frac{0}{6} \log_2 0 = 0$$

$$P(C_2) = \frac{6}{6} = 1$$

$$\text{Entropy } H(S) = 0 \log_2 0 + 1 \log_1 = 0$$

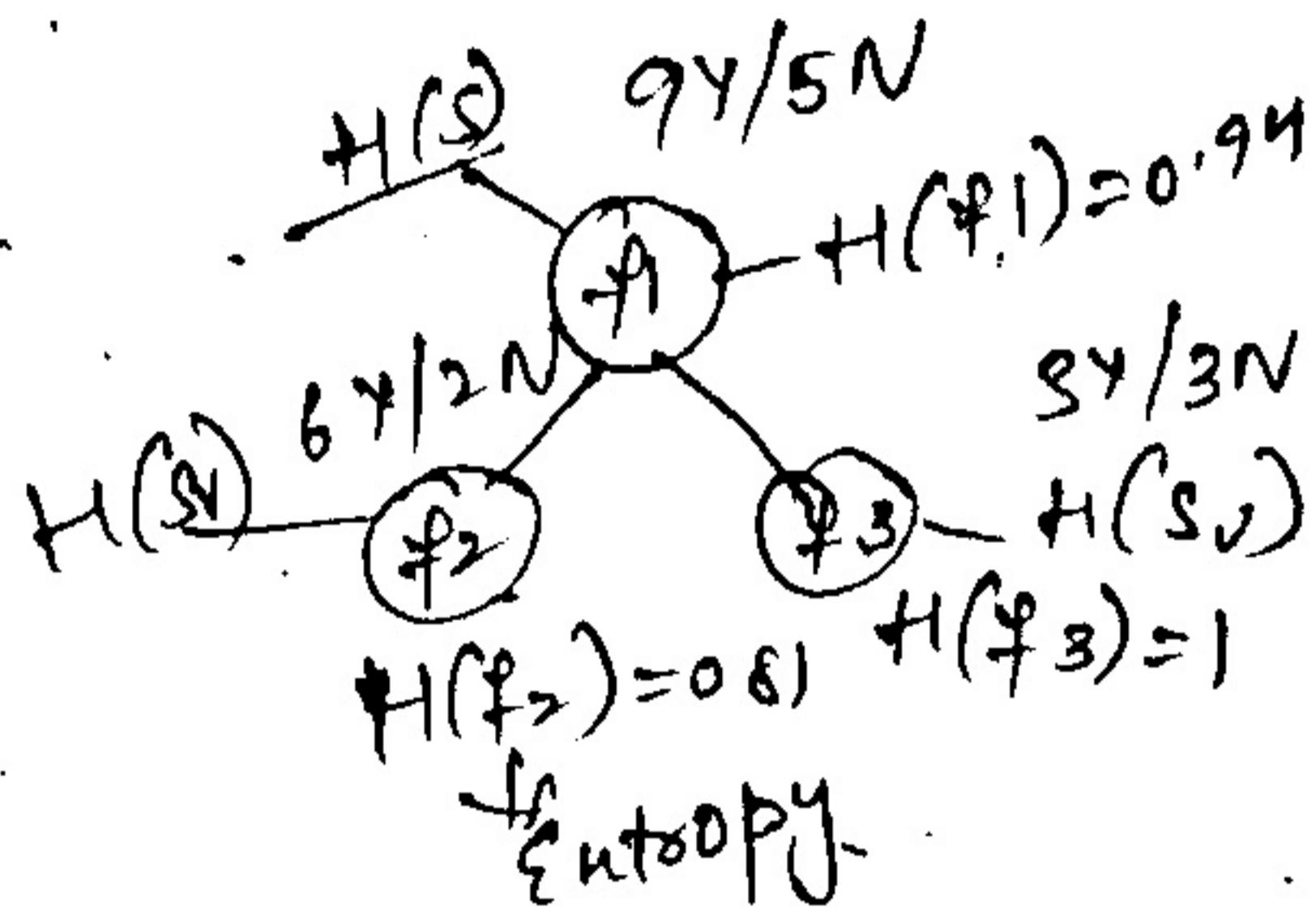
### Information Gain

$$\text{Gain}(S, A) = H(S) - \frac{1}{|S|} \sum |S_v| H(S_v)$$

$$\text{Gain}(S, f_1) = H(S) - \frac{8}{14} H(f_2) - \frac{6}{14} H(f_3)$$

$$\Rightarrow 0.91 - \frac{8}{14} \times 0.81 - \frac{6}{14} \times 1 = 0.049$$

Higher the Information gain, better split



→ Select the attribute with highest Information Gain  $f_1$

→ let  $P_i$  be the probability that an arbitrary tuple in  $D$

belongs to class  $C_i$ , estimated by  $|C_i \cap D| / |D|$

$$\rightarrow \text{Expected information (Entropy)} \text{ needed to classify a tuple in } D: \quad \text{Info}(D) = - \sum_{i=1}^m P_i \log_2 (P_i)$$

→ Information needed to split  $D$  into  $v$  partitions to classify  $D$ :  $\text{Info}_v(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j)$

$$\text{Information Gained} \boxed{Gain(A) = \text{Info}(D) - \text{Info}_v(D)}$$

So Example! - Using Information Gain as impurity measure

Contingency Table  $g_1 \rightarrow$

| Class | $g_1=T$ | $g_1=F$ |
|-------|---------|---------|
| +     | 3       | 1       |
| -     | 1       | 4       |

| Class | $g_2=T$ | $g_2=F$ |
|-------|---------|---------|
| +     | 2       | 2       |
| -     | 3       | 2       |

$\Rightarrow$  Core All information needed

to classify a tuple D is

$$\text{Info}_A(D) = \sum_{j=1}^k \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2 p_i$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$\Rightarrow$  Overall information needed to classify a tuple

$$\text{Info}(D) = -\frac{4}{9} \log_2 \frac{4}{9} - \frac{5}{9} \log_2 \frac{5}{9} = 0.989$$

$\Rightarrow$  Expected Information needed to classify a tuple D acc to  $g_1$

$$\text{Info}_{g_1}(D) = \frac{4}{9} \left[ -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right] + \frac{5}{9} \left[ -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} \right]$$

$$= 0.760$$

$\Rightarrow$  Information Gain ( $g_1$ )  $= \text{Info}(D) - \text{Info}_{g_1}(D) = 0.989 - 0.760 = 0.229$

$\Rightarrow$  Information Gain ( $g_1$ )  $> \text{Info}(D) - \text{Info}_{g_2}(D)$

$\Rightarrow$  Expected Information needed to classify tuple D acc to  $g_2$

$$\text{Info}_{g_2}(D) = \frac{4}{9} \left[ \frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right] + \frac{5}{9} \left[ -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right]$$

$$= 0.985$$

$\Rightarrow$  Information Gain ( $g_2$ )  $= \text{Info}(D) - \text{Info}_{g_2}(D) = 0.989 - 0.985 = 0.004$

Information Gain ( $g_2$ )  $<$  Information Gain ( $g_1$ )  $\Rightarrow g_2$  will be selected as

$\Rightarrow$  Gain( $g_2$ ) is more often  $g_2$ , it will be selected as first splitting attribute.

## $\Rightarrow$ Gain Ratio

→ Information Gain measure is biased towards attributes with large number of values. C4.5 (ID3 successor) uses gain ratio to overcome the problem (normalization of information gain)

$$\text{Gain Ratio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}(A)}$$

$$\text{SplitInfo}(D) = -\sum_{j=1}^V \frac{|D_j|}{|D|} \times \log_2 \frac{|D_j|}{|D|}$$

Sof!-  $\text{SplitInfo}_{q_1}(D) = -\frac{4}{9} \log_2 \frac{4}{9} - \frac{5}{9} \log_2 \frac{5}{9} = 0.989$

$$\text{splitInfo}_{q_2}(D) = -\frac{5}{9} \log_2 \frac{5}{9} - \frac{4}{9} \log_2 \frac{4}{9}$$

$$= -0.555 * (-0.847) - 0.444 (-1.169) = 0.989$$

$$\text{Gain Ratio}(q_1) = \frac{\text{Gain}(q_1)}{\text{SplitInfo}_{q_1}(D)} = \frac{0.229}{0.989} = 0.23$$

$$\text{Gain Ratio}(q_2) = \frac{\text{Gain}(q_2)}{\text{SplitInfo}_{q_2}(D)} = \frac{0.004}{0.989} = 0.004$$

The attribute with maximum gain ratio is selected as splitting attribute, hence  $q_1$  will be selected as first splitting attribute

Example!-

| age       | income             | student | Credit-rating | buys-computer |
|-----------|--------------------|---------|---------------|---------------|
| $\leq 30$ | high               | no      | fair          | no            |
| $\leq 30$ | high               | no      | excellent     | yes           |
| $31-40$   | high               | no      | fair          | yes           |
| $> 40$    | medium             | no      | fair          | yes           |
| $> 40$    | <del>low</del> low | yes     | fair          | yes           |
| $> 40$    | low                | yes     | excellent     | no            |
| $31-40$   | low                | yes     | excellent     | yes           |
| $\leq 20$ | medium             | no      | fair          | no            |
| $\leq 20$ | low                | yes     | fair          | yes           |
| $> 40$    | medium             | yes     | fair          | yes           |
| $\leq 30$ | medium             | yes     | excellent     | yes           |
| $31-40$   | medium             | no      | excellent     | yes           |
| $31-40$   | high               | yes     | fair          | yes           |
| $> 40$    | medium             | no      | excellent     | no            |

- Class P<sub>0</sub> - buys - Computer "Yes"
- class N<sub>1</sub> - buys - Computer = "NO"

$$\Rightarrow \text{Info}(D) = I(9,5) = \frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14}$$

$$I(9,5) = 0.940$$

| age   | P <sub>i</sub> | n <sub>i</sub> | I(P <sub>i</sub>  n <sub>i</sub> ) |
|-------|----------------|----------------|------------------------------------|
| ≤ 30  | 2              | 3              | 0.971                              |
| 31-40 | 4              | 0              | 0                                  |
| > 40  | 3              | 2              | 0.971                              |

$$\Rightarrow \text{Info}_{\text{age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{3}{14} I(3,2)$$

$$= 0.694$$

$\frac{5}{14} I(2,3)$  → means "age  $\leq 30$ " → has 5 out of 14 samples with 2 yes's and 3 NO's hence

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D)$$

$$= 0.940 - 0.694$$

$$= 0.246$$

Similarly, Gain(Income) = 0.029

$$\text{Gain(Student)} = 0.015$$

$$\text{Gain(Credit-rating)} = 0.048$$

Gain Ratio:

$$\text{Split Info}_{\text{Income}}(D) = \frac{-4}{14} \log_2 \frac{4}{14} - \frac{6}{14} \log_2 \frac{6}{14} - \frac{4}{14} \log_2 \frac{4}{14}$$

$$= 1.557$$

$$\text{gain-ratio}(\text{Income}) = \frac{0.029}{1.557} = 0.019$$

The attribute with maximum gain ratio is selected as the splitting attribute.

## $\Rightarrow$ Rule based Classification

$\Rightarrow$  Given a tuple  $X$ , from a class labeled data set  $D$ ,  
 let  $n_{\text{covers}}$  be the number of tuples covered by  $R$   
 $n_{\text{correct}}$  be the number of tuples correctly classified by  $R$   
 $|D| \rightarrow$  be the number of tuples in  $D$ , then

$$\text{Coverage}(R) = \frac{n_{\text{covers}}}{|D|}$$

$$\text{accuracy}(R) = \frac{n_{\text{correct}}}{n_{\text{covers}}}$$

$\Rightarrow$  let  $(\text{Pos}) \rightarrow$  be the Positive and  $(\text{neg})$  be the negative tuples  
 covered by  $R$

let  $(\text{Pos}')$  → be the Positive and  $(\text{neg}')$  be the negative tuples  
 covered by  $R'$ , FOIL assesses the information gained by

extending condition as

$$\text{FOIL-Gain} = \text{Pos} \times \left( \log_2 \frac{\text{Pos}'}{\text{Pos} + \text{neg}'} - \log_2 \frac{\text{Pos}}{\text{Pos} + \text{neg}} \right)$$

$$\text{FOIL-PRUNE}(R) = \frac{\text{Pos} - \text{neg}}{\text{Pos} + \text{neg}} \quad \begin{cases} \text{if FOIL-Gain value is} \\ \text{higher, then we} \\ \text{Prune } R \end{cases}$$

$$\Rightarrow \text{Likelihood-Ratio} = 2 * \sum_{i=1}^m f_i \log_2 \left( \frac{f_i}{e_i} \right) \rightarrow \text{Higher} \rightarrow \text{better}$$

- $\rightarrow m \rightarrow$  number of classes
- $\rightarrow f_i \rightarrow$  observed freq of each class among tuples
- $\rightarrow e_i \rightarrow$  expected freq if the rule made random predictions
- $\rightarrow$  used by CN2

Example 1 :- Derive the possible rules for classification of bank customer into defaulter or not defaulter classes based on the following data set:

| Tid | Home-owner | MaritalStatus | Defaulter |
|-----|------------|---------------|-----------|
| 1   | Yes        | Single        | NO        |
| 2   | NO         | Married       | NO        |
| 3   | NO         | Single        | NO        |
| 4   | Yes        | Married       | NO        |
| 5   | NO         | Divorced      | Yes       |
| 6   | Yes        | Married       | NO        |
| 7   | NO         | Divorced      | Yes       |
| 8   | NO         | Single        | NO        |
| 9   | NO         | Married       | NO        |
| 10  | NO         | Single        | Yes       |

few direct possible Rules by looking at data set

$$r_1 = (\text{Home-owner} = \text{Yes}) \Rightarrow (\text{Defaulter} = \text{NO})$$

$$r_2 = (\text{Marital-status} = \text{married}) \Rightarrow (\text{Defaulter} = \text{NO})$$

$$r_3 = (\text{Home-owner} = \text{NO}) \& (\text{Marital-status} = \text{Married}) \Rightarrow (\text{Defaulter} = \text{NO})$$

$r_i : A \rightarrow y$        $\rightarrow A$  is rule antecedent  
 $y$  is rule consequent

$$\text{Coverage}(r_1) = \frac{|A|}{|D|} = \frac{n_{\text{covers}}}{|D|} = \frac{3}{10} = 0.3 = 30\%$$

$$\text{Accuracy}(r_1) = \frac{n_{\text{correct}}}{n_{\text{covers}}} = \frac{3}{3} = 1 = 100\%$$

$$\text{Coverage}(r_2) = \frac{4}{10} = 0.4 = 40\% ; \text{Accuracy}(r_2) = \frac{4}{4} = 1 = 100\%$$

$$\text{Coverage}(r_3) = \frac{2}{10} = 0.2 = 20\% ; \text{Accuracy}(r_3) = \frac{2}{2} = 1 = 100\%$$

$\Rightarrow r_2 \rightarrow$  has high accuracy with high coverage  $\rightarrow$  stronger

⇒ Example :- Consider a training dataset that contains 60 +ve samples & 100 -ve samples. Suppose following two candidate rules are made:-  
 r1: → Covers 50 +ve samples & 5 -resamples  
 r2: → Covers 2 +ve samples & 2 -resamples. , which rule is better. determine Information gain.

Sol ! - Coverage (r<sub>1</sub>) =  $\frac{55}{160} = 0.343 = 34.3\%$ , Accuracy(r<sub>1</sub>) =  $\frac{50}{55} = 0.909 = 90.9\%$

Coverage (r<sub>2</sub>) =  $\frac{2}{160} = 0.0125 = 1.25\%$ , Accuracy (r<sub>2</sub>) =  $\frac{2}{2} = 100\%$

r<sub>1</sub> is better despite its lower accuracy, the high accuracy for r<sub>2</sub> is potentially dangerous because of low coverage of rule.

$$\text{FOL-Infogain}(r_1) = P_1 \times \left[ \log_2 \frac{P_1}{P_1 + N_1} - \log_2 \frac{P_0}{P_0 + N_0} \right]$$

$$\text{FOL-gain}(r_1) = 50 \times \left[ \log_2 \frac{50}{50+55} - \log_2 \frac{60}{60+100} \right]$$

$$= 50 [-0.137 + 1.415] = 63.9$$

$$\text{FOL-gain}(r_2) = 2 \times \left[ \log_2 \frac{2}{2} - \log_2 \frac{60}{60+160} \right] = 2.83$$

As FOL gain for r<sub>1</sub> is better than r<sub>2</sub>, r<sub>1</sub> is better rule.

### Model Evaluation and Selection :-

#### Confusion Matrix

| Predicted Class | C <sub>1</sub>       | -C <sub>1</sub>      |
|-----------------|----------------------|----------------------|
| Actual class    |                      |                      |
| C <sub>1</sub>  | True Positives (TP)  | False Negatives (FN) |
| -C <sub>1</sub> | False Positives (FP) | True Negatives (TN)  |

$$\Rightarrow \text{Accuracy} \Rightarrow \frac{TP+TN}{P+N}$$

$TP \Rightarrow$  no. of True Positive  
 $TN \Rightarrow$  No. of True Negative  
 $FP \Rightarrow$  false Positive  
 $P \rightarrow$  Total Positive  
 $N \rightarrow$  Total Negative

$$\Rightarrow \text{Error rate} \Rightarrow 1 - \text{accuracy}$$

Or :-  $\frac{(FP+FN)}{P+N}$

$$\Rightarrow \text{Sensitivity} = \frac{TP}{P} \quad \begin{matrix} \text{True Positive} \\ \text{recognition} \\ \text{rate} \end{matrix}$$

$$\Rightarrow \text{Specificity} = \frac{TN}{N} \quad \begin{matrix} \text{True Negative} \\ \text{recognition} \\ \text{rate} \end{matrix}$$

$$\Rightarrow \text{Precision} : - \frac{TP}{TP+FP}$$

(exactness)

$$\Rightarrow \text{Recall} = \frac{TP}{TP+FN} = \frac{TP}{P}$$

(Completeness)

Perfect score = 1.0

$$\Rightarrow F \text{ measure } \{F_1 \text{ or f score}\} =$$

(harmonic mean of Precision & Recall)

$$F = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\Rightarrow F_B = \frac{(1+\beta^2) \times \text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

$\rightarrow$  weight measure

$$\beta^2 = (1-\alpha)/\alpha$$

Example :-

|              |   | Predicted Class $\rightarrow$ |     | $\downarrow$ | Actual class |      |  |
|--------------|---|-------------------------------|-----|--------------|--------------|------|--|
|              |   | Y                             | N   |              |              |      |  |
| Actual class | Y | 900                           | 100 | 1000         | 1000         | 2000 |  |
|              | N | 200                           | 800 | 1000         | 1000         | 2000 |  |
|              |   | 1100                          | 900 | 2000         |              |      |  |

$$TP = 900, FP = 200, FN = 100, TN = 800$$

$$\text{Precision} P = \frac{TP}{TP+FP} = \frac{900}{900+200} = 0.818$$

$$\text{Recall } R = \frac{TP}{TP+FN} = \frac{900}{900+100} = 0.9$$

$$F\text{Score} = \frac{2 \times TP}{2 \times TP + FP + FN} = \frac{2 \times R \times P}{R + P} = \frac{2 \times 900}{2 \times 900 + 200 + 100} = 0.85$$

$$\begin{aligned} \text{Sensitivity} \\ (\text{True +ve Rate}) &= \frac{TP}{TP+FN} = \frac{900}{900+100} = 0.9 \end{aligned}$$

$$\begin{aligned} \text{Specificity} \\ (\text{True -ve Rate}) &= \frac{TN}{TN+FP} = \frac{800}{800+200} = 0.8 \end{aligned}$$

$$\text{False +ve Rate} = \frac{FP}{TN+FP} = \frac{200}{800+200} = 0.2$$

$$\text{False -ve Rate} = \frac{FN}{TP+FN} = \frac{100}{900+100} = 0.1$$

### Simple Linear Regression

$$Y = w_0 + w_1 X \rightarrow \begin{aligned} w_0 &\rightarrow \text{x-intercept} \\ w_1 &\rightarrow \text{slope} \end{aligned}$$

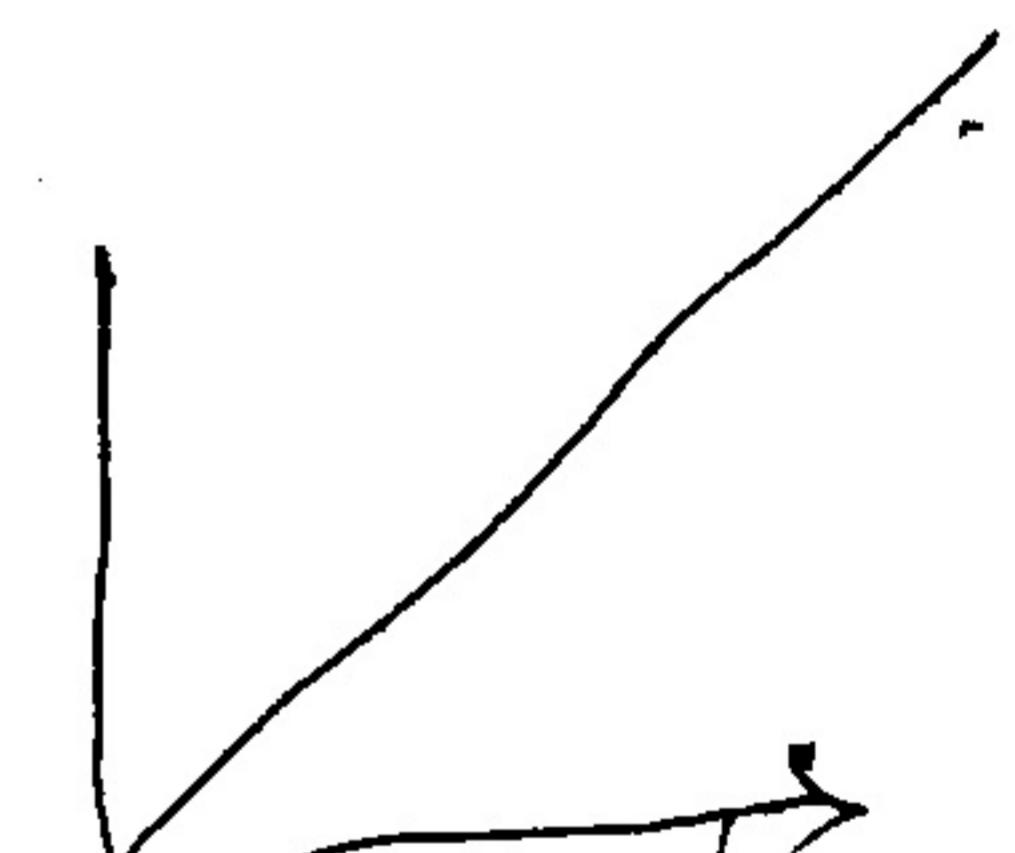
$$w_0 = \frac{(\sum Y)(\sum X^2) - (\sum X)(\sum XY)}{n(\sum X^2) - (\sum X)^2}$$

$$w_1 = \frac{n(\sum XY) - (\sum X)(\sum Y)}{n(\sum X^2) - (\sum X)^2}$$

$n \rightarrow$  number of tuples

$$w_0 = \frac{(1)(30) - (10)(54)}{(4)(30) - 100} = 1.5$$

$$w_1 = \frac{(4)(54) - (10)(30)}{(4)(30) - 100} = 1.3$$



| X | Y | XY | X <sup>2</sup> | P    | Error     |
|---|---|----|----------------|------|-----------|
| 1 | 3 | 3  | 1              | 2.8  | 0.2       |
| 2 | 4 | 8  | 4              | 4.1  | 0.1       |
| 3 | 5 | 15 | 9              | 5.4  | 0.4       |
| 4 | 7 | 28 | 16             | 6.7  | 0.3       |
|   |   | 54 | 30             | 19   |           |
|   |   |    |                | Σ 30 |           |
|   |   |    |                |      | Predicted |

$$Y = 1.3X + 1.5$$

for any given  $x \rightarrow$  find  $y$

## Association Rule Analysis

Itemset :-  $\{ \text{Milk, Bread, Diaper} \}$

K-itemset

→ Itemset that contains K items

| ID | Items                       |
|----|-----------------------------|
| 1  | Bread, Milk                 |
| 2  | Bread, Diaper, Butter, Beer |
| 3  | Milk, Diaper, Butter, Coke  |
| 4  | Bread, Milk, Diaper, Butter |
| 5  | Bread, Milk, Diaper, Coke   |

Support Count (C) :- Frequency of an occurrence of Itemset

$$\text{Eg: } D(\{ \text{Milk, Bread, Diaper} \}) = 2$$

Support :- Fraction of transactions that contain an itemset

$$\text{Eg: } S(\text{Milk, Bread, Diaper}) = 2/5$$

Frequent Itemset :

- An Itemset whose support is greater than or equal to minsup threshold.

Association Rule  $X \rightarrow Y$

where X & Y are Itemsets

$$\text{Ex: } \{ \text{Milk, Diaper} \} \rightarrow \text{Butter}$$

- Support (S) :- Fraction of transactions that contain both X & Y

$$S = \frac{D(\text{Milk, Diaper, Butter})}{T} = \frac{2}{5} = 0.4$$

- Confidence (C) :- Measure how often items in Y appear in transactions that contain X

$$C = \frac{D(\text{Milk, Diaper, Butter})}{D(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

→ Support  $\geq$  minsup threshold

confidence  $\geq$  minconf threshold

## ⇒ Apriori Algorithm

Apriori employs an iterative approach known as level-wise search where  $k$ -itemsets are used to explore  $(k+1)$ -itemsets. First set of frequent 1-itemset is found by scanning DB to accumulate the count of each item and collecting those items that satisfy minimum support. The resulting set is denoted by  $L_1$ . Next  $L_1$  is used to find  $L_2$ , the set of frequent 2-itemsets, which is used to find  $L_3$  and so on, ~~until~~ until no more frequent  $k$ -item sets can be found! -

Example:- Min support = 50%

Threshold confidence %.

| TID | Items   |
|-----|---------|
| 100 | 1 3 4   |
| 200 | 2 3 5   |
| 300 | 1 2 3 5 |
| 400 | 2 5     |

(II)  $L_2$

| Itemset | Support     |
|---------|-------------|
| {1, 2}  | 1/4 → 25% X |
| {1, 3}  | 2/4 → 50%   |
| {1, 5}  | 1/4 → 25% X |
| {2, 3}  | 2/4 → 50%   |
| {2, 5}  | 3/4 → 75%   |
| {3, 5}  | 2/4 → 50%   |

(I)  $L_1$

| Itemset | Support   |
|---------|-----------|
| 1       | 2/4 = 50% |
| 2       | 3/4 = 75% |
| 3       | 3/4 → 75% |
| 4       | 1/4 → 25% |
| 5       | 3/4 → 75% |

(III)  $L_3$

| Itemset   | Support     |
|-----------|-------------|
| {1, 3, 5} | 1/4 = 25% X |
| {2, 3, 5} | 2/4 = 50%   |
| {1, 2, 3} | 1/4 = 25% X |

$L_1 \rightarrow \{1, 2, 3, 5\}$

Itemset  $\Rightarrow \{1, 2, 3, 5\}$

this  
can  
be  
dropped

$\Rightarrow \{2, 3, 5\}$

Example: - Min Sup = ~~10%~~ 2

| TID | Items List           |
|-----|----------------------|
| 1   | $I_1, I_2, I_5$      |
| 2   | $I_2, I_4$           |
| 3   | $I_1, I_2, I_4$      |
| 4   | $I_1, I_3$           |
| 5   | $I_2, I_3$           |
| 6   |                      |
| 7   | $I_1, I_3$           |
| 8   | $I_1, I_2, I_3, I_5$ |
| 9   | $I_1, I_2, I_3$      |

~~Step 1~~ ~~Step 2~~  
Step 1 : Get frequency support count  
 $I_1 - 6$   
 $I_2 - 7$   
 $I_3 - 6$   $\Leftarrow C_1$   
 $I_4 - 2$   
 $I_5 - 2$   
 Candidate set of  $C_1$   
 $\min \text{support count} = 2$   
 Nothing will be removed

Step 2 - ItemSet | SuppCount

|                |     |
|----------------|-----|
| $\{I_1, I_2\}$ | 4   |
| $\{I_1, I_3\}$ | 4   |
| $\{I_1, I_4\}$ | 1 X |
| $\{I_1, I_5\}$ | 2   |
| $\{I_2, I_3\}$ | 4   |
| $\{I_2, I_4\}$ | 2   |
| $\{I_2, I_5\}$ | 2   |
| $\{I_3, I_4\}$ | 0 X |
| $\{I_3, I_5\}$ | 1 X |
| $\{I_4, I_5\}$ | 0 X |

→ Compare Candidate Support count with min support count

| Itemset        | SuppCount |
|----------------|-----------|
| $\{I_1, I_2\}$ | 4         |
| $\{I_1, I_3\}$ | 4         |
| $\{I_1, I_5\}$ | 2         |
| $\{I_2, I_3\}$ | 4         |
| $\{I_2, I_4\}$ | 2         |
| $\{I_2, I_5\}$ | 2         |

$L_2$

Step 3! - Itemset | SuppCount

|                     |     |
|---------------------|-----|
| $\{I_1, I_2, I_3\}$ | 2   |
| $\{I_1, I_2, I_5\}$ | 2   |
| $\{I_1, I_2, I_4\}$ | 1 X |
| $\{I_1, I_3, I_5\}$ | 1 X |
| $\{I_2, I_3, I_4\}$ | 0 X |
| $\{I_2, I_3, I_5\}$ | 1 X |
| $\{I_3, I_4, I_5\}$ | 0 X |

| Itemset             | Supp. |
|---------------------|-------|
| $\{I_1, I_2, I_3\}$ | 2     |
| $\{I_1, I_2, I_5\}$ | 2     |