



**Problem 1: Association Mining Rules (7 points)**

Giving the following database with four transactions and a minimum support threshold of 60% and a minimum confidence threshold of 80%.

1. Find all frequent itemsets and association rules using
  - a) A priori
  - b) FP-Growth.
2. Compare the efficiency of both processes.

| Transaction ID | Items           |
|----------------|-----------------|
| T1             | {K, A, D, B}    |
| T2             | {D, A, C, E, B} |
| T3             | {C, A, B, E}    |
| T4             | {B, A, D}       |

**Problem 2: Decision Tree Classifier (7 points)**

Given the training data in the table below:

| Outlook ( $X_1$ ) | Temperature ( $X_2$ ) | Humidity ( $X_3$ ) | Play Tennis? ( $Y$ ) |
|-------------------|-----------------------|--------------------|----------------------|
| sunny             | hot                   | high               | no                   |
| overcast          | hot                   | high               | yes                  |
| rain              | mild                  | high               | yes                  |
| rain              | cool                  | normal             | yes                  |
| sunny             | mild                  | high               | no                   |
| sunny             | mild                  | normal             | yes                  |
| rain              | mild                  | normal             | yes                  |
| overcast          | hot                   | normal             | yes                  |

1. Using the dataset above, calculate the mutual information for each feature ( $X_1$ ,  $X_2$ ,  $X_3$ ) to determine the root node for a Decision Tree trained on the above data.
2. Calculate what the next split should be.
3. Draw the resulting tree.

**Problem 3: Data preprocessing (6 points)**

~~1)~~ Use the methods below to normalize the following group of data:

200, 300, 400, 600, 1000

~~a)~~ min-max normalization by setting min = 0 and max = 1

~~b)~~ z-score normalization

~~c)~~ normalization by decimal scaling

2. Given the following data (in increasing order) for the attribute age: 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 36, 40, 45, 46, 52, 70.

~~a)~~ Use smoothing by bin means to smooth the above data, using a bin depth of 3. Illustrate your steps. Comment on the effect of this technique for the given data.

b) How might you determine outliers in the data?

---

$$\text{minSup} = 60\% = 0.6$$

| TID | Items         | Item | Support      |
|-----|---------------|------|--------------|
| T1  | K, A, D, B    | A    | 4/4 = 1 ✓    |
| T2  | D, A, C, E, B | B    | 4/4 = 1 ✓    |
| T3  | C, A, B, E    | C    | 2/4 = 0.5    |
| T4  | B, A, D       | D    | 3/4 = 0.75 ✓ |
|     |               | E    | 2/4 = 0.5    |
|     |               | K    | 1/4 = 0.25   |

C1

$$L1 = \{A, B, D\}$$

| Itemset | Support |
|---------|---------|
| A, B    | 4/4 ✓   |
| A, D    | 3/4 ✓   |
| B, D    | 3/4 ✓   |

$$L2 = \{ \{A, B\}, \{A, D\}, \{B, D\} \}$$

| Itemset | Support |
|---------|---------|
| A, B, D | 3/4 ✓   |

$$L3 = \{ \{A, B, D\} \}$$

$$\text{min Conf} = 0.8$$

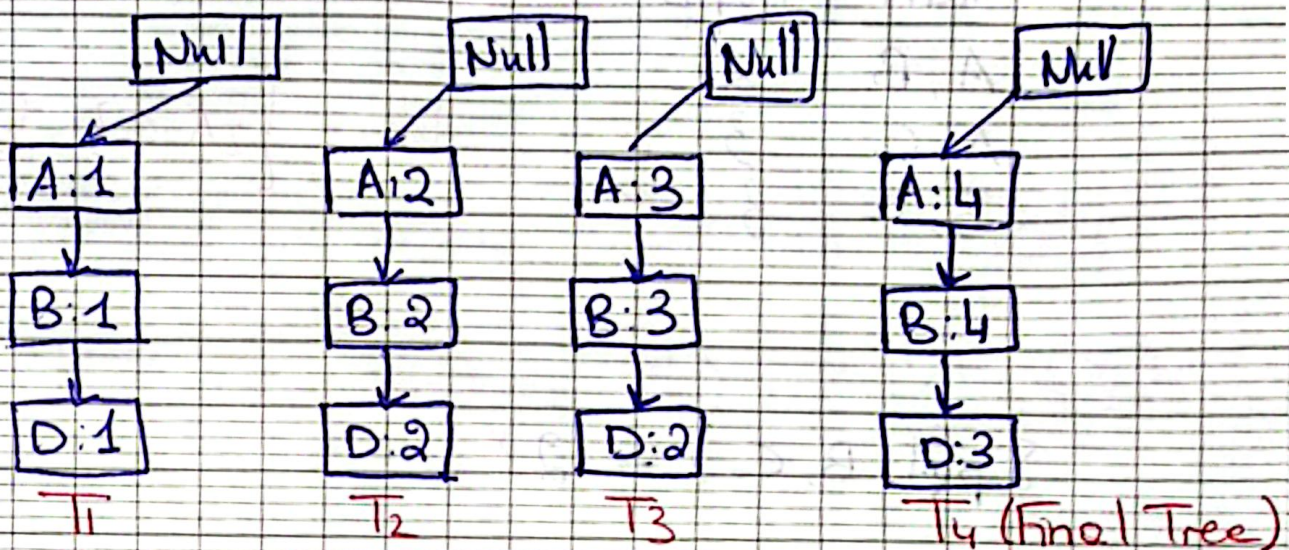
Rules

|       |       |        |       |
|-------|-------|--------|-------|
| A → B | 4/4 ✓ | A → BD |       |
| B → A | 4/4 ✓ | BD → A | 3/3 ✓ |
| A → D | 3/4   | B → AD | 3/4   |
| D → A | 3/3 ✓ | AD → B | 3/3 ✓ |
| B → D | 3/4   | D → AB | 3/3 ✓ |
| D → B | 3/3 ✓ | AB → D | 3/4   |

## Step 1: Count and Sort

| Item | Supp | TID            | Items         | Sorted items |
|------|------|----------------|---------------|--------------|
| A    | 4/4  | T <sub>1</sub> | K, A, D, B    | A, B, D      |
| B    | 4/4  | T <sub>2</sub> | D, A, C, E, B | A, B, D      |
| D    | 3/4  | T <sub>3</sub> | C, A, B, E    | A, B         |
|      |      | T <sub>4</sub> | B, A, D       | A, B, D      |

## Step 2: Build FP-Tree



## Step 3: Confidence Pattern Base and Confidence FP-Tree

| Item | Conf. Pattern Base | Conf. FP Tree | F.P Generated             |
|------|--------------------|---------------|---------------------------|
| D(3) | {A, B: 3}          | {A: 3, B: 3}  | {D: 3}, {A, B: 3}         |
| B(4) | {A: 4}             | {A: 4}        | {B, D: 3}, {A, B, D: 3}   |
| A(4) | { }                | { }           | {B: 4}, {A, B: 4}, {A: 4} |

### 4) Rules

A → D 3/4  
 D → A 3/3 ✓  
 B → D 3/4  
 D → B 3/3 ✓

AB → D 3/4  
 A → B 4/4 ✓  
 B → A 4/4 ✓  
 AB → BD 3/4

BD → A 3/3 ✓  
 B → AD 3/4  
 AD → B 3/3 ✓  
 D → AB 2/3 ✓

Problem 2:

Before Splitting  $\{Y^6, N^2\}$

$$G.I(B) = 1 - \left[ \left(\frac{6}{8}\right)^2 + \left(\frac{2}{8}\right)^2 \right] = 0.625$$

G.I for each attrb.

1) Outlook:

Sunny  $\{Y^1, N^2\}$     0:  $\{Y^2, N^0\}$     R:  $\{Y^3, N^0\}$

$$G.I(S) = 1 - \left[ \left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] = 0.44$$

$$G.I(0) = 0$$

$$G.I(R) = 0$$

$$G.I(\text{Outlook}) = \frac{3}{8} (0.44) = 0.165$$

$$\begin{aligned} \text{Gain}(S, \text{Outlook}) &= G.I(B) - G.I(\text{Outlook}) \\ &= 0.625 - 0.165 = 0.46 \end{aligned}$$

2) Temperature: H:  $\{Y^2, N^1\}$     M:  $\{Y^3, N^1\}$     C:  $\{Y^1, N^0\}$

$$G.I(H) = 1 - \left[ \left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right] = 0.44$$

$$G.I(M) = 1 - \left[ \left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right] = 0.375$$

$$G.I(C) = 0$$

$$\begin{aligned} G.I(\text{Temp}) &= \frac{3}{8} (0.44) + \frac{4}{8} (0.375) \\ &= 0.3525 \end{aligned}$$

$$\text{Gain}(S, \text{Temp}) = 0.625 - 0.3525 = 0.27$$

3) Humidity:  $H: \{Y^2, N^2\}$   $N: \{Y^4, N^0\}$

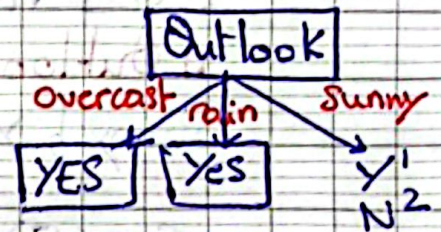
$$G.I(H) = 1 - \left[ \left(\frac{2}{4}\right)^2 + \left(\frac{2}{4}\right)^2 \right] = 0.5$$

$$G.I(N) = 0$$

$$G.I(\text{Humidity}) = \frac{4}{8} (0.5) = 0.25$$

$$\text{Gain}(S, \text{Hum.}) = 0.625 - 0.25 = 0.37$$

Highest Gain: Outlook



Filter the dataset on 'overcast' and 'rain'

| $\theta$ | T    | H      | PY? |
|----------|------|--------|-----|
| Sunny    | hot  | high   | N   |
| Sunny    | mild | high   | N   |
| Sunny    | mild | normal | Y   |

$$G.I(B) = 1 - \left[ \left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right]$$

$$= 0.44$$

1) Temperature:  $H: \{Y^0, N^1\}$   $M: \{Y^1, N^1\}$

$$G.I(H) = 0 / G.I(M) = 0.5$$

$$G.I(\text{Temp}) = \frac{2}{3} (0.5) = 0.33$$

$$\text{Gain}(\text{Temp}) = 0.44 - 0.33 = 0.11$$

2) Humidity:  $H: \{Y^0, N^2\}$   $N: \{Y^1, N^0\}$

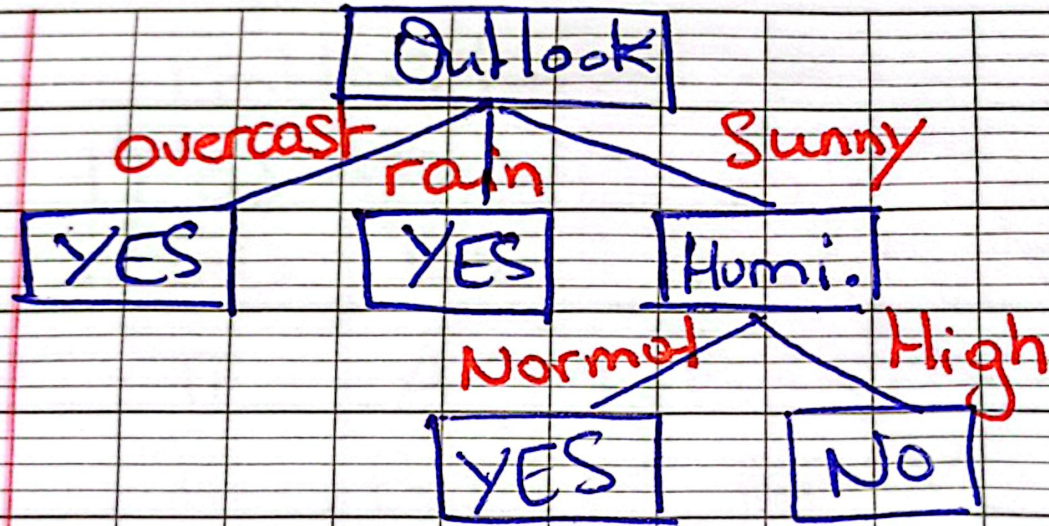
$$G.I(H) = 0$$

$$G.I(N) = 0$$

$$\text{Gain}(\text{Humidity}) = 0.44$$

$$G.I(\text{Humidity}) = 0$$

(Highest)  
next split



### Problem 3

1) a) Min - Max Normalization

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

$$\min = 200, \max = 1000$$

$$200' = \frac{200 - 200}{1000 - 200} = 0$$

$$300' = \frac{300 - 200}{1000 - 200} = 0.125$$

$$400' = \frac{400 - 200}{1000 - 200} = 0.25$$

$$500' = \frac{500 - 200}{1000 - 200} = 0.5$$

$$1000' = \frac{1000 - 200}{1000 - 200} = 1$$

0, 0.125, 0.25, 0.5, 1

### b) z-Score normalization

$$x = \frac{x - \mu}{\sigma}$$

$$\text{Mean } (\mu) = 500$$

$$\text{Std } (\sigma) = 282.84$$

$$200' = 200 - 500 / 282.84 = -1.06$$

$$300' = 300 - 500 / 282.84 = -0.71$$

$$400' = 400 - 500 / 282.84 = -0.35$$

$$600' = 600 - 500 / 282.84 = 0.35$$

$$1000' = 1000 - 500 / 282.84 = 1.77$$

### c) Decimal Scaling

$$x' = x / 10^j$$

$$j = \log_{10}(\max|x|)$$

$$\max = 1000$$

$$j = 3$$

$$200' = 200 / 10^3 = 0.2$$

$$300' = 300 / 10^3 = 0.3$$

$$400' = 400 / 10^3 = 0.4$$

$$600' = 600 / 10^3 = 0.6$$

$$1000' = 1000 / 10^3 = 1$$

step:

2) Group into Bins of Depth 3

13 15 16/16 19 20/20 21 22/22 25 25/

25 25 30/33 33 35/35 35 35/36 40 45

46 52 70/

Step 2: Compute the mean of each Bin

- 1) 14.66
- 2) 18.33
- 3) 21
- 4) 24
- 5) 26.6
- 6) 33.66
- 7) 35
- 8) 40.33
- 9) 56

Step 3: Replace each value with the mean of the corresponding bin

14.66 14.66 14.66 18.33 18.33 18.33  
21 21 21 24 24 24 26.6 26.6 26.6 ...

Effect: Reduces noise

b) **Z-score Method** values with  $z \text{ score} > 3$  or  $< -3$  are outliers  
**IQR Method**

Values less than  $Q1 - 1.5 \times IQR$  or greater than  $Q3 + 1.5 \times IQR$  are considered outliers