

Partial Exam
Introduction to ML
Spring 2024

Duration: 70 minutes
Dr. Abbas Rammal

Problem 1: Association Mining Rules (12 points)

We provide the following database of store and customer transactions. Create all association rules with a minimum support of 60% and a minimum confidence of 60%.

1. Find all frequent itemsets using
 - a) A priori
 - b) FP-Growth.
2. We have the following rule: {Bananas, Bread} => {Milk}. Calculate the Lift, Leverage and conviction coefficients of this rule. Interpret.

Transaction ID Items

T1	Apples, Milk, Bread, Cheese
T2	Bananas, Milk, Bread, Yogurt
T3	Apples, Bread, Eggs, Cheese
T4	Bananas, Milk, Eggs, Yogurt
T5	Apples, Milk, Bread, Yogurt

Problem 2: Decision Tree Classifier (8 points)

We consider a real dataset with three categorical attributes (Weather, Temperature, Day) and a binary target variable (Play: Yes or No). Here are 10 records in the dataset:

Weather	Day	Temperature	Play
Sunny	Weekday	30	Yes
Sunny	Weekend	32	Yes
Rainy	Weekday	20	No
Sunny	Weekend	28	Yes
Rainy	Weekday	18	No
Rainy	Weekend	22	No
Sunny	Weekday	26	Yes
Sunny	Weekend	35	No
Rainy	Weekday	17	No

- a) Take the optimal split for the "Temperature" attribute in your decision tree: 27 and >27. Construct a decision tree based on this training data. For splitting, use information gain as measure for impurity. Build a separate branch for each attribute.
- b) Split the Temperature attribute in two way using Gini Index.

Problem 1:

1) a) Apriori

min Sup = 60% = 3/5

TID	Items	Items	Support
T1	A, M, Br, C	A	3/5 ✓
T2	Ba, M, Br, Y	M	4/5 ✓
T3	A, Br, E, C	Br	4/5 ✓
T4	Ba, M, E, Y	C	2/5
T5	A, M, Br, Y	Ba	2/5
		Y	3/5 ✓
		E	2/5
		C	

$$L_1 = \{A, M, Br, Y\}$$

ItemSet	Support
---------	---------

A, M	2/5
A, Br	3/5 ✓
A, Y	1/5
M, Br	3/5 ✓
M, Y	3/5 ✓
Br, Y	2/5

$$L_2 = \{ \{A, Br\}, \{M, Br\}, \{M, Y\} \}$$

ItemSet	Support
---------	---------

A, M, Br	2/5 x
A, M, Y	1/5 x
A, Br, Y	1/5 x
M, Br, Y	2/5 x

Frequent Itemsets:

$$\{A\}, \{M\}, \{Br\}, \{Y\}$$

$$\{A, Br\}, \{M, Br\}, \{M, Y\}$$

Association Rules

minConf: 0.6

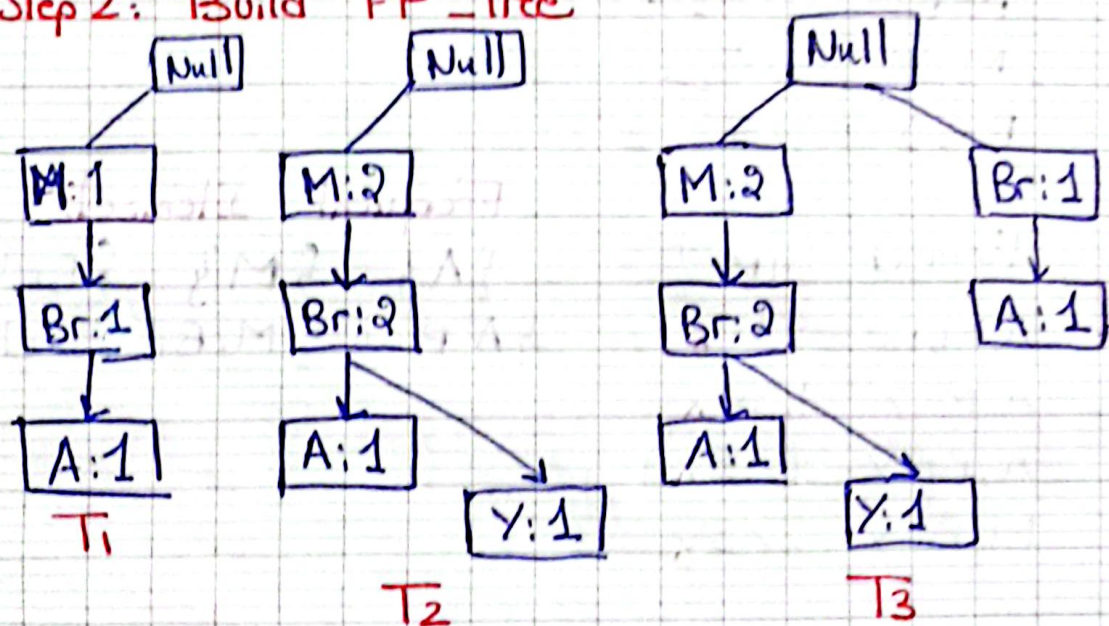
- $(A \rightarrow Br) = 3/3$ ✓
- $(Br \rightarrow A) = 3/4$ ✓
- $(M \rightarrow Br) = 3/4$ ✓
- $(Br \rightarrow M) = 3/4$ ✓
- $(M \rightarrow Y) = 3/4$ ✓
- $(Y \rightarrow M) = 3/3$ ✓

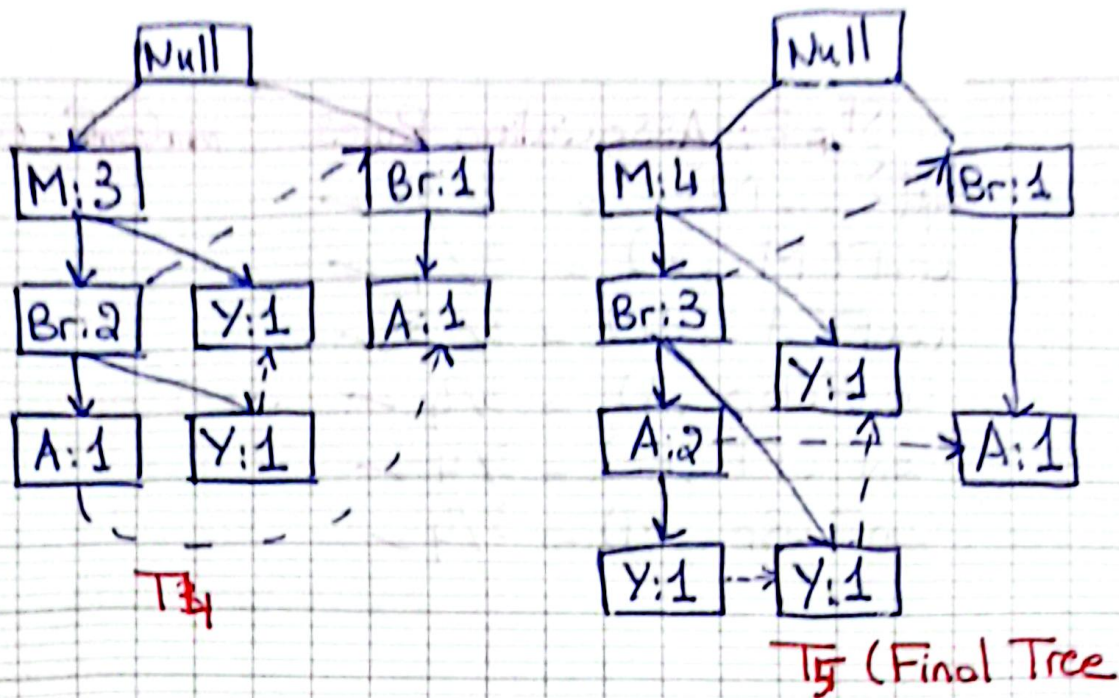
b) FP_Growth

Step 1: Count and Sort

Item	Supp	TID	Items	Ordered Items
M	4	T1	A, M, Br, C	M, Br, A
Br	4	T2	Ba, M, Br, Y	M, Br, Y
A	3	T3	A, Br, E, C	Br, A
Y	3	T4	Ba, M, E, Y	M, Y
		T5	A, M, Br, Y	M, Br, A, Y

Step 2: Build FP_Tree





Step 3: Conditional Pattern Base and Conditional FP tree

Items	Cond. Pattern Base	Cond. FP Tree
Y(3)	{M, Br, A:1} {M, Br:1} {M:4}	{ M:3, Br:2, A:1 }
A(3)	{M, Br:2} {Br:1}	{ M:2, Br:3 }
Br(4)	{M:3}	{M:3}
M(4)	{}	{}

Step 4: Frequent Pattern Generated

Items	FP Generated
Y(3)	{Y:3} {Y, M:3}
A(3)	{A:3} {A, Br:3}
Br(4)	{Br:4} {M, Br:3}
M(4)	{M:4}

Step 5: Association Rules

minConf = 0.6

$$\text{Sup}(Y \rightarrow M) = 3/3 \checkmark$$

$$\text{Conf}(M \rightarrow Y) = 3/4 \checkmark$$

$$\text{Conf}(A \rightarrow Br) = 3/3 \checkmark$$

$$\text{Conf}(Br \rightarrow A) = 3/4 \checkmark$$

$$\text{Conf}(M \rightarrow Br) = 3/4 \checkmark$$

$$\text{Conf}(Br \rightarrow M) = 3/4 \checkmark$$

$$2) \{Ba, Br\} \rightarrow \{M\}$$

$$\begin{aligned} \rightarrow \text{LIFT}(\{Ba, Br\} \rightarrow \{M\}) &= \frac{\text{Conf}(\{Ba, Br\} \rightarrow \{M\})}{\text{Support } \{M\}} \\ &= \frac{1/5 / 1/5}{4/5} = \frac{5}{4} = 1.25 \end{aligned}$$

Interpretation:

$\text{LIFT} = 1.25 > 1 \Rightarrow \{Bananas, Bread\}$ positively correlates with $\{Milk\}$, implying the presence of Ba, Br increase the likelihood of buying Milk

$$\rightarrow \text{Leverage}(\{Ba, Br\} \rightarrow \{M\}) = \text{Supp}$$

$$\begin{aligned} &\text{Supp}(Ba, Br, M) - \text{Supp}(Ba, Br) \times \text{Supp}(M) \\ &= \frac{1}{5} - \frac{1}{5} \times \frac{4}{5} = 0.04 \end{aligned}$$

Interpretation:

The rule ^{occurs} more frequent than expected if $\{Bananas, Bread\}$ and $\{Milk\}$ were independent

$$\begin{aligned} \rightarrow \text{Conviction}(\{Ba, Br\} \rightarrow \{M\}) &= \frac{1 - \text{Supp}(M)}{1 - \text{Conf}(Ba, Br \rightarrow M)} \\ &= \frac{1 - \frac{4}{5}}{1 - 1} = \infty \end{aligned}$$

Interpretation:

The rule is perfectly confident, as Milk always appears when Bananas and Bread are purchased

problem 2:

	Weather	Day	Temp	Play	Before Splitting
1	Sunny	weekDay	>27	Y	$E(B) = -\frac{4}{9} \log_2 \left(\frac{4}{9}\right) - \frac{5}{9} \log_2 \left(\frac{5}{9}\right)$ $= 0.99$
2	Sunny	weekEnd	>27	Y	
3	Rainy	weekDay	27	N	
4	Sunny	weekEnd	>27	Y	
5	Rainy	weekDay	27	N	
6	Rainy	weekEnd	27	N	
7	Sunny	weekDay	27	Y	
8	Sunny	weekEnd	>27	N	
9	Rainy	weekDay	27	N	

Entropy for each attribute

1) Weather: Sunny: $\{Y^4, N^1\}$ Rainy: $\{Y^0, N^4\}$

$$E(S) = -\frac{4}{5} \log_2 \left(\frac{4}{5}\right) - \frac{1}{5} \log_2 \left(\frac{1}{5}\right)$$

$$= 0.72$$

$$E(R) = 0$$

$$E(\text{Weather}) = \frac{5}{9} (0.72) = 0.4$$

$$\text{Gain (Weather)} = 0.99 - 0.4 = 0.59$$

2) Day: WD: $\{Y^2, N^3\}$ WE: $\{Y^2, N^2\}$

$$E(\text{WD}) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)$$

$$= 0.97$$

$$E(\text{WE}) = -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right) = 1$$

$$E(\text{Day}) = \frac{5}{9} (0.97) + \frac{4}{9} = 0.98$$

$$\text{Gain (Day)} = 0.99 - 0.98 = 0.01$$

3) Temp: $27: \{Y^1, N^4\}$ $>27: \{Y^3, N^1\}$

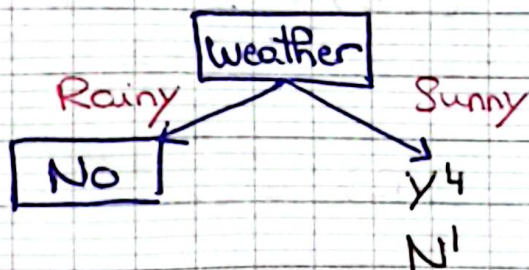
$$E(27) = -\frac{1}{5} \log_2\left(\frac{1}{5}\right) - \frac{4}{5} \log_2\left(\frac{4}{5}\right) = 0.72$$

$$E(>27) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \left(\frac{1}{4}\right) \log_2\left(\frac{1}{4}\right) = 0.81$$

$$E(\text{Temp}) = \frac{5}{9} (0.72) + \frac{4}{9} (0.81) = 0.76$$

$$\text{Gain}(\text{Temp}) = 0.99 - 0.76 = 0.23$$

Weather have the highest Gain we split on it first



New dataset filtered on "Sunny"

Weather	Day	Temp	Play
Sunny	week day	>27	Y
Sunny	week end	>27	Y
Sunny	week end	>27	Y
Sunny	week day	27	Y
Sunny	week end	>27	N

$$E(B) = -\frac{4}{5} \log_2\left(\frac{4}{5}\right) - \frac{1}{5} \log_2\left(\frac{1}{5}\right) = 0.72$$

1) Day: WD: $\{Y^2, N^0\}$ WE: $\{Y^2, N^1\}$

$$E(WD) = 0$$

$$E(WE) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \left(\frac{1}{3}\right) \log_2\left(\frac{1}{3}\right) = 0.918$$

$$E(\text{Day}) = \frac{3}{5} (0.918) = 0.5508$$

$$\text{Gain}(\text{S, Day}) = 0.72 - 0.5508 = 0.169$$

2) Temp: $\leq 27: \{Y^1, N^0\}$ $> 27: \{Y^3, N^1\}$

$$E(\leq 27) = 0$$

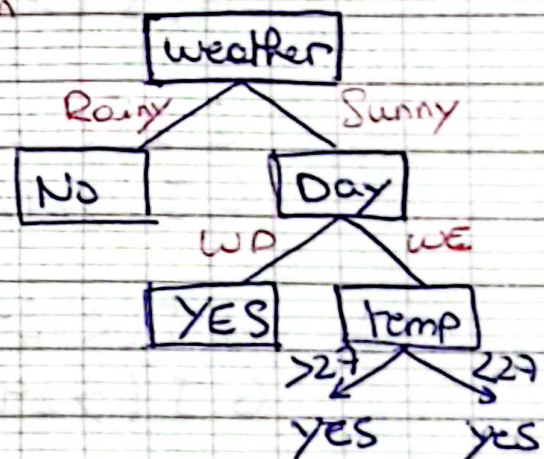
$$E(> 27) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \left(\frac{1}{4}\right) \log_2\left(\frac{1}{4}\right) = 0.81$$

$$E(\text{Temp}) = \frac{4}{5} (0.81) = 0.648$$

$$\text{Gain}(\text{Temp}) = 0.72 - 0.648 = 0.072$$

Day have the highest gain

Weather	Day	Temp	Play
Sunny	WE	≤ 27	Y
Sunny	WE	≤ 27	Y
Sunny	WE	≤ 27	N



b) 1) Sort Temp in increasing order

	17	18	20	22	26	28	30	32	35					
	17	19	21	24	27	29	31	33						
	<=	>	<=	>	<=	>	<=	>	<=	>				
Y	0	4	0	4	0	4	1	3	2	2	3	1	4	0
N	1	4	2	3	3	2	4	1	4	1	4	1	4	1
G.I	0.44	0.38	0.29	0.17	0.34	0.44	0.49	0.44	0.44					

least G.I.

$$G.I(<=17) = 0$$

$$G.I(>17) = 1 - \left[\left(\frac{4}{8}\right)^2 + \left(\frac{4}{8}\right)^2 \right] = \frac{1}{2} \quad \left. \begin{array}{l} \text{avg} = \frac{2+2}{2} = 2 \\ \text{avg} = \frac{2+2}{2} = 2 \end{array} \right\}$$

$$G.I(<=19) = 0$$

$$G.I(>19) = 1 - \left[\left(\frac{4}{7}\right)^2 + \left(\frac{3}{7}\right)^2 \right] = 0.49 \quad \left. \begin{array}{l} \text{avg} = \frac{7+9}{2} = 0.49 \\ \text{avg} = 0.38 \end{array} \right\}$$

$$G.I(<=21) = 0$$

$$G.I(>21) = 1 - \left[\left(\frac{4}{6}\right)^2 + \left(\frac{2}{6}\right)^2 \right] = 0.44 \quad \left. \begin{array}{l} \text{avg} = \frac{6+9}{2} = 0.44 \\ \text{avg} = 0.29 \end{array} \right\}$$

$$G.I(<=24) = 0$$

$$G.I(>24) = 1 - \left[\left(\frac{4}{5}\right)^2 + \left(\frac{1}{5}\right)^2 \right] = 0.32 \quad \left. \begin{array}{l} \text{avg} = \frac{5+9}{2} = 0.32 \\ \text{avg} = 0.17 \end{array} \right\}$$

$$G.I(<=27) = 1 - \left[\left(\frac{1}{5}\right)^2 + \left(\frac{4}{5}\right)^2 \right] = 0.32$$

$$G.I(>27) = 1 - \left[\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right] = 0.375 \quad \left. \begin{array}{l} \text{avg} = \frac{5+9}{2} = 0.32 \\ \text{avg} = \frac{4+9}{2} = 0.375 \\ \text{avg} = 0.34 \end{array} \right\}$$

$$G.I(<=29) = 1 - \left[\left(\frac{2}{6}\right)^2 + \left(\frac{4}{6}\right)^2 \right] = 0.44$$

$$G.I(>29) = 1 - \left[\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right] = 0.44 \quad \left. \begin{array}{l} \text{avg} = \frac{6+9}{2} = 0.44 \\ \text{avg} = \frac{3+9}{2} = 0.44 \\ \text{avg} = 0.44 \end{array} \right\}$$

$$G.I(<=31) = 1 - \left[\left(\frac{3}{7}\right)^2 + \left(\frac{4}{7}\right)^2 \right] = 0.49$$

$$G.I(>31) = 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right] = 0.5 \quad \left. \begin{array}{l} \text{avg} = \frac{7+9}{2} = 0.49 \\ \text{avg} = \frac{2+9}{2} = 0.5 \\ \text{avg} = 0.49 \end{array} \right\}$$

$$\left. \begin{aligned} G_I(\leq 33) &= 1 - \left[\left(\frac{4}{3}\right)^2 + \left(\frac{4}{3}\right)^2 \right] = 0.5 \\ G_I(> 33) &= 0 \end{aligned} \right\} \begin{array}{l} 8/9 (0.5) \\ 0.44 \end{array}$$