

## 1) Multiple Choice

12 / 20

- Each correctly answered multiple choice question gives 1 point.
- Each incorrect answer results in -1 point.
- However, the minimum scores in each of the four groups are 0 points.

### Data processing and feature spaces

3 / 5

Temperature in Celsius is of ☒ interval scale ☐ ratio scale.

The idea of TF-IDF text representation is that globally frequent terms are

☒ more ☐ less informative than less frequent terms.

Every categorical feature with  $m$  possible values is convertible into  $m$  binary features.

☒ True ☐ False

Manhattan (L1) distance between two points is always smaller or equal than Euclidean (L2) distance between the same two points.

☐ True ☒ False

When the mean is larger than the median, the underlying distribution is

☒ positively skewed ☐ negatively skewed.

### Association Rule Mining

4 / 5

What is the size of the itemset generated from the transaction {bread, butter, beer, milk, butter, beer}? ☐ 2 ☒ 4 ☐ 6

If confidence of the rule  $A \rightarrow BCD$  is below the confidence threshold then all other rules created from the itemset  $\{ABCD\}$  are necessarily below the threshold as well.

☐ True ☒ False

Which algorithm requires that the items are sorted by frequency?

☐ Apriori ☒ FP Growth

A closed frequent itemset has ☐ no frequent supersets ☐ has no immediate frequent supersets ☒ has no immediate frequent superset of the same support.

Downward closing property:

When an item set  $X$  is frequent, all its ☐ subsets ☐ supersets are also frequent.

**Classification**

0 / 5

Classification is ☐ an unsupervised task ☒ a supervised task.Decision tree classifiers are of the ☐ lazy ☒ eager learning type.An attribute, in which all training examples have different values is ☐ problematic  
☒ helpful for decision trees based on information gain in terms of generalization error.The accuracy / error rate evaluation metric is robust with respect to class imbalance.  
☐ True ☐ FalseIn a k-fold cross validation, each data point belongs to the test set exactly  
☐ once ☒ k times.**Clustering**

5 / 5

Clustering is ☒ an unsupervised task ☐ a supervised task.The k-Means algorithm is ☒ sensitive to outliers ☐ not sensitive to outliers.DBSCAN is ☐ a partitioning-based clustering ☒ a density-based clustering approach.Two points A and B are density-connected via another point C only if  
☐ A and B are core points ☒ C is a core point ☐ A, B, and C are core points.

Agnes and Diana both refer to hierarchical clustering approaches.

☒ True ☐ False



The following table shows a list of transactions:

T1	Burger, Wrap
T2	Burger, Coke, Fries
T3	Burger, Coke, Fries
T4	Burger, Fries, Wrap
T5	Burger, Coke, Wrap
T6	Coke, Fries

a) Apply the Apriori algorithm with a minimum support of 0.5. Construct for each step the candidate set  $C_k$  and the frequent itemset list  $L_k$  starting with  $k = 1$  until all frequent itemsets are generated. For each step, also list the itemsets that are pruned based on the apriori property and list the itemsets that are pruned due to the application of the minimum support threshold.

$$\text{minimum support} = 0.5 \times 6 = 3$$

$$C_1 = \text{Burger}^{(5)}, \text{Coke}^{(3)}, \text{Fries}^{(4)}, \text{Wrap}^{(3)} \checkmark$$

$$L_1 = \text{Burger}^{(5)}, \text{Coke}^{(4)}, \text{Fries}^{(4)}, \text{Wrap}^{(3)} \checkmark$$

no itemsets are pruned in this step  $\checkmark$

$$C_2 = \{\text{Burger, Coke}\}^{(3)}, \{\text{Burger, Fries}\}^{(3)}, \{\text{Burger, Wrap}\}^{(3)}$$

$$\{\text{Coke, Fries}\}^{(3)}, \{\text{Coke, Wrap}\}^{(1)}, \{\text{Fries, Wrap}\}^{(1)} \checkmark$$

$\{\text{Coke, Wrap}\}$  and  $\{\text{Fries, Wrap}\}$  should be pruned due to minimum support threshold.  $\checkmark$

$$L_2 = \{\text{Burger, Coke}\}^{(3)}, \{\text{Burger, Fries}\}^{(3)}, \{\text{Burger, Wrap}\}^{(3)}, \{\text{Coke, Fries}\}^{(3)} \checkmark$$

$$C_3 = \{\text{Burger, Coke, Fries}\}^{(2)}, \{\text{Burger, Coke, Wrap}\}^{(1)},$$

$\{\text{Burger, Fries, Wrap}\}^{(1)}$ ,  $\{\text{Coke, Fries, Wrap}\}^{(1)}$  will be pruned based on apriori property, and all the 3 itemsets are pruned based on the minimum support threshold.  $\checkmark$

b) Generate all possible rules from the frequent itemsets and calculate their confidence.

$$C: \{\text{Burger}\} \rightarrow \{\text{Coke}\} = \frac{S\{\text{Burger, Coke}\}}{S\{\text{Burger}\}} = \frac{3}{5} \checkmark$$

$$C: \{\text{Coke}\} \rightarrow \{\text{Burger}\} = \frac{S\{\text{Burger, Coke}\}}{S\{\text{Coke}\}} = \frac{3}{4} \checkmark$$

$$C: \{\text{Burger}\} \rightarrow \{\text{Fries}\} = \frac{3}{5} \checkmark, C: \{\text{Fries}\} \rightarrow \{\text{Burger}\} = \frac{3}{4} \checkmark$$

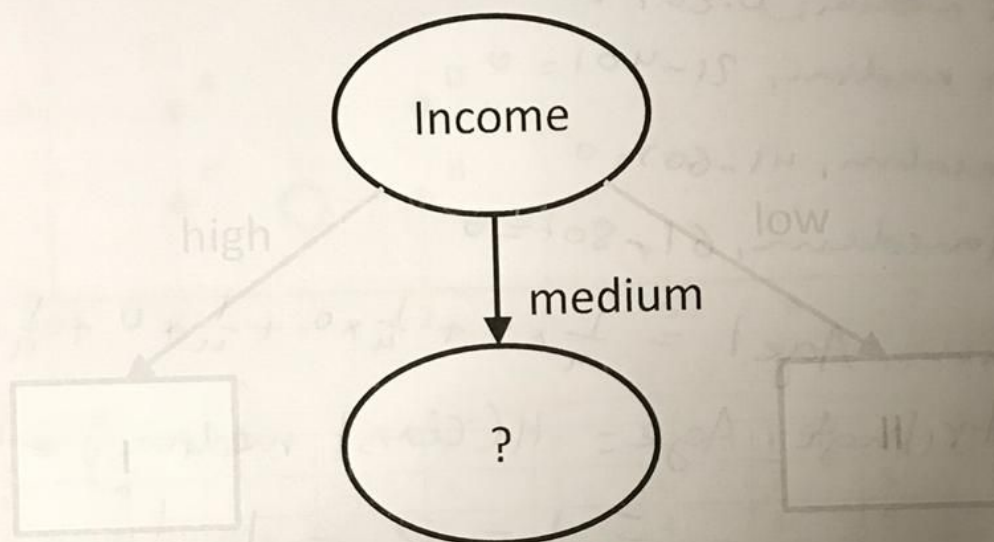
$$C: \{\text{Burger}\} \rightarrow \{\text{Wrap}\} = \frac{3}{5} \checkmark, C: \{\text{Wrap}\} \rightarrow \{\text{Burger}\} = \frac{3}{3} \checkmark$$

$$C: \{\text{Coke}\} \rightarrow \{\text{Fries}\} = \frac{3}{4} \checkmark, C: \{\text{Fries}\} \rightarrow \{\text{Coke}\} = \frac{3}{4} \checkmark$$



Given the following dataset and partial decision tree:

Age	Car	Income	Class
0-20	no	high	I
0-20	no	medium	II
21-40	yes	medium	I
21-40	no	low	II
41-60	yes	low	II
41-60	no	medium	I
61-80	yes	high	I
61-80	yes	medium	II



Calculate the information gain for the remaining attributes (car, age) to complete the decision tree for the *medium* branch. Decide which attribute should be used for the next split.

$$H(\text{class}) = 1, H(\text{class}|\text{car})?$$

$\frac{2}{4}$  no cars,  $\frac{2}{4}$  yes cars  
 $\frac{1}{4}$  class I,  $\frac{1}{4}$  class II  
 $\frac{1}{4}$  class I,  $\frac{1}{4}$  class II

$$H(\text{car}) = -\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} = -\frac{2}{4}(-1) - \frac{2}{4}(-1) = 1$$

$$H(\text{class}|\text{car}) = -\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} = -\frac{2}{4}(-1) - \frac{2}{4}(-1) = 1$$

$$H(\text{class}|\text{car}) = \frac{1}{2} \times 1 + \frac{1}{2} \times 1 = 1, H(\text{class}|\text{car}) = 1$$

$$\text{IG of car attribute} = H(\text{class}) - H(\text{class}|\text{car}) = 1 - 1 = 0$$

Medium, car

For age:

$$\begin{aligned} \frac{1}{4} (0-20) &\rightarrow \frac{1}{2} \text{ class I} \\ \frac{2}{4} (21-40) &\rightarrow \frac{1}{2} \text{ class I} \\ \frac{1}{4} (41-60) &\rightarrow \frac{1}{2} \text{ class I} \\ \frac{2}{4} (61-80) &\rightarrow \frac{1}{2} \text{ class II} \end{aligned}$$

~~For medium:~~

$$H(\text{class} | \text{medium}, 0-20) = 0$$

$$H(\text{class} | \text{medium}, 21-40) = 0$$

$$H(\text{class} | \text{medium}, 41-60) = 0$$

$$H(\text{class} | \text{medium}, 61-80) = 0$$

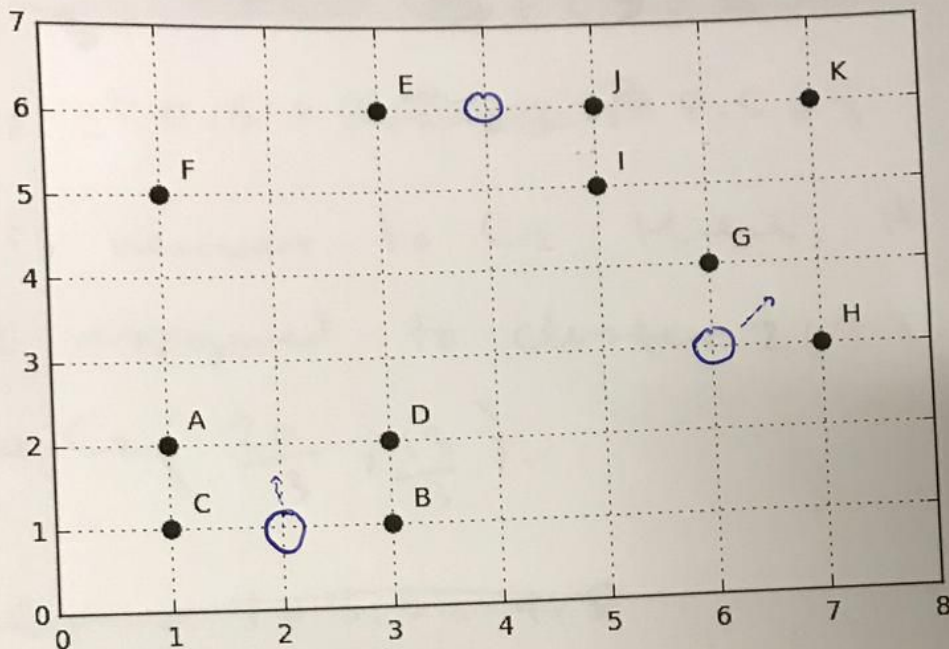
$$H(\text{class} | \text{medium}, \text{Age}) = \frac{1}{4} \times 0 + \frac{1}{4} \times 0 + \frac{1}{4} \times 0 + \frac{1}{4} \times 0 = 0$$

$$\begin{aligned} \text{IG of attribute Age} &= H(\text{class} | \text{medium}) - H(\text{class} | \text{medium}, \text{Age}) \\ &= 1 - 0 = 1 \quad \checkmark \end{aligned}$$

IG of Age larger than IG of Car, then the attribute that should be used on next split is Age.  $\checkmark$



The following datapoints are given:



	A	B	C	D	E	F	G	H	I	J	K
x	1	3	1	3	3	1	6	7	5	5	7
y	2	1	1	2	6	5	4	3	5	6	6

Cluster the datapoints with k-Means using Manhattan distance. The initial centroids for k-Means are (2,1), (6,3), and (4,6), i.e. the parameter  $k = 3$ . Specify for each iteration, to which centroids each point got assigned to and the calculation of the new centroids.

Iteration 1:

$C_1 = (2,1)$  : assigned objects = {

~~C, B, A, C, B~~ f.  $C_1 \left( \frac{1+3+1}{3}, \frac{2+1+1}{3} \right) = \left( \frac{5}{3}, \frac{4}{3} \right)$  f.

$C_2 = (6,3)$  : assigned objects = {

G, H

✓  $C_2 \left( \frac{6+7}{2}, \frac{4+3}{2} \right)$  ✓

$C_3 = (4,6)$  : assigned objects = {

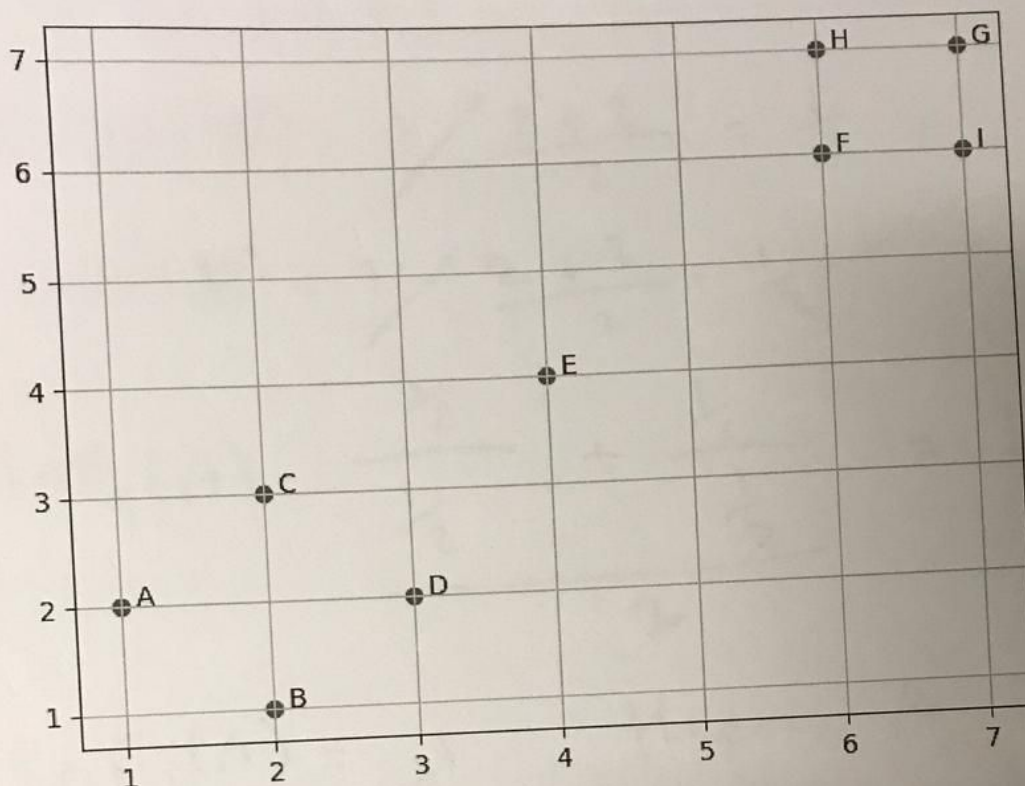
E, J, I, F, K

f.  $C_3 \left( \frac{3+5+5+1+7}{5}, \frac{6+6+5+5+6}{5} \right)$  f.

Iteration 2:

~~Assigned to objects~~ ~~A, B, D assigned to cluster~~  
 ~~$C_1 \left( \frac{5}{3}, \frac{4}{3} \right)$~~   ~~$C_2 \left( \frac{13}{2}, \frac{7}{2} \right)$~~   ~~$C_3 \left( \frac{13}{3}, \frac{17}{2} \right)$~~   
 no assignment

The following dataset is given. Use Manhattan distance for your calculations.



Use the local outlier factor method ( $LOF_2$ ) to calculate the scores and decide whether the points A and E are outliers given a threshold of 1.

For point E,  $LOF_2(E)$

$$N_2(E) = \{C, D\}$$

$$N_2(C) = \{A, D\}$$

$$N_2(D) = \{A, B, C\}$$

$$LOF_2(E) = \frac{1}{\frac{3+3}{2}} = \frac{2}{6} = \frac{1}{3}$$

$$LOF_2(C) = \frac{1}{\frac{2+2}{2}} = \frac{1}{2}$$

$$LOF_2(D) = \frac{1}{\frac{2+2}{2}} = \frac{1}{2}$$

$$LOF_2(E) = \frac{\left(\frac{1}{2}\right) + \left(\frac{1}{2}\right)}{\frac{1}{\frac{1}{3}}} = \frac{\frac{3}{2} + \frac{3}{2}}{2} = \frac{3}{2} = \frac{3}{2}$$

~~E is not outlier~~

$LOF_2(E) = \frac{3}{2} > 1 \rightarrow E$  is outlier



# 5) Outlier Detection: Local Outlier Factor

$$N_2(A) = (C, B) \quad \checkmark$$

$$N_2(B) = (A, D) \quad \checkmark$$

$$lrd_2(A) = \frac{1}{2} \quad \checkmark$$

$$lrd_2(A) = 1 / \frac{2+2}{2} = \frac{1}{2} \quad \checkmark$$

$$lrd_2(B) = 1 / \frac{2+2}{2} = \frac{1}{2} \quad \checkmark$$

$$lof_2(A) = \frac{\frac{1}{2}}{\frac{1}{2}} + \frac{\frac{1}{2}}{\frac{1}{2}} = \frac{1+1}{2} = 1 \quad \checkmark$$

$$lof_2(A) = 1 \quad \checkmark$$

then A is not outlier

$$lrd_2(C)? \quad \checkmark$$

\$ 10