

Exam DM868/DM870/DS804 spring 2022

The exam contains 18 questions summing up to 100 points. The point value of each question is stated at the beginning of its text in parentheses. All questions ask for evaluation of multiple statements with (yes/no) or (true/false) answers. You will gain the full point value of a question if you evaluate all of its statements correctly. Penalties apply to wrong evaluations. Hence, skipping to evaluate a statement is very likely to be more advantageous for you than making a random guess. The points earned from an X-point question are calculated by:

$$X*(C-W)/A$$

where

C : Number of correctly evaluated statements,

W: Number of wrongly evaluated statements,

A: Total number of statements in a question.

Note that $C+W = A$ if you evaluate all statements and $C+W < A$ if you skip evaluating at least one statement.

(5 points)

Given the items $I = \{A, B, C, D, E, F, G, H, I\}$
and the set of transactions T :

| TransID | Items |
|---------|-------------|
| 1 | A B C F G H |
| 2 | A B D E G H |
| 3 | A B E F H I |
| 4 | A C F G |
| 5 | A C F G H I |
| 6 | A D E G H I |
| 7 | A I |
| 8 | B E F H I |
| 9 | B E I |
| 10 | D E H |
| 11 | G |

For the minimum support of 3, we already determined the frequent 3-itemsets with the APRIORI algorithm:

$$L_3 = \{ABH, ACF, ACG, AEH, AFG, AFH, AGH, AHI, \\ BEH, BEI, BFH, CFG, DEH, EHI, FHI\}$$

Which of the following 4-itemsets are preliminary candidates in the next step of APRIORI (i.e., after the merging step but before pruning)?

| | True | False |
|------|-------------------------------------|-------------------------------------|
| ACFG | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| AEFH | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| AFHI | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| AFGH | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| BEHI | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| BEFH | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| BFHI | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| CEGI | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

| | | |
|------|-----------------------|-------------------------------------|
| DEHI | <input type="radio"/> | <input checked="" type="checkbox"/> |
| EFHI | <input type="radio"/> | <input checked="" type="checkbox"/> |

(5 points)

For some transaction database we found that the rule $\{A, B, C, D\} \Rightarrow \{E, F, G, H\}$ has a confidence below the confidence threshold.

Which of the following rules will therefore have a confidence below the confidence threshold as well?

| | True | False |
|---|-------------------------------------|-------------------------------------|
| $\{A\} \Rightarrow \{B, C, D, E, F, G, H\}$ | <input checked="" type="checkbox"/> | <input type="radio"/> |
| $\{A, B, C, D\} \Rightarrow \{E, F\}$ | <input type="radio"/> | <input checked="" type="checkbox"/> |
| $\{A, C, D\} \Rightarrow \{B, E, F\}$ | <input type="radio"/> | <input checked="" type="checkbox"/> |
| $\{A, C\} \Rightarrow \{B, D, E, F, G, H\}$ | <input checked="" type="checkbox"/> | <input type="radio"/> |
| $\{A, D\} \Rightarrow \{B, E, F, G\}$ | <input type="radio"/> | <input checked="" type="checkbox"/> |
| $\{A, D\} \Rightarrow \{B, E, F, G, H\}$ | <input type="radio"/> | <input checked="" type="checkbox"/> |
| $\{A, D\} \Rightarrow \{B, C, E, F, G, H\}$ | <input checked="" type="checkbox"/> | <input type="radio"/> |
| $\{B, E, F\} \Rightarrow \{A, C, D\}$ | <input type="radio"/> | <input checked="" type="checkbox"/> |
| $\{C\} \Rightarrow \{A, B, D, E, F, G, H\}$ | <input checked="" type="checkbox"/> | <input type="radio"/> |
| $\{C, D\} \Rightarrow \{A, B, E, F, G, H\}$ | <input checked="" type="checkbox"/> | <input type="radio"/> |

(5 points)

We have the following one-dimensional dataset:

| ID | Value |
|----|-------|
| A | 2 |
| B | 4 |
| C | 6 |
| D | 10 |
| E | 14 |
| F | 16 |
| G | 18 |

In three attempts, k-means delivered the following three clustering solutions:

$$S_1 = \{A, B, C\}, \{D, E, F, G\}$$

$$S_2 = \{A, B\}, \{C, D\}, \{E, F, G\}$$

$$S_3 = \{A, B, C, D\}, \{E, F, G\}$$

We want to compare the solutions using TD^2 . Which of the following statements are correct?

| | True | False |
|---|-------------------------------------|-------------------------------------|
| S_1 is better than S_2 in terms of TD^2 . | <input type="radio"/> | <input checked="" type="checkbox"/> |
| S_2 is better than S_3 in terms of TD^2 . | <input checked="" type="checkbox"/> | <input type="radio"/> |
| S_1 and S_3 are equally good in terms of TD^2 . | <input type="radio"/> | <input checked="" type="checkbox"/> |
| S_3 is better than S_1 in terms of TD^2 . | <input checked="" type="checkbox"/> | <input type="radio"/> |

(8 points)

We have the following one-dimensional dataset:

| ID | Value |
|----|-------|
| A | 2 |
| B | 4 |
| C | 6 |
| D | 10 |
| E | 14 |
| F | 16 |
| G | 18 |

In three attempts, k-means delivered the following three clustering solutions:

$$S_1 = \{A, B, C\}, \{D, E, F, G\}$$

$$S_2 = \{A, B\}, \{C, D\}, \{E, F, G\}$$

$$S_3 = \{A, B, C, D\}, \{E, F, G\}$$

We want to compare the solutions using simplified Silhouette. Which of the following statements are correct?

| | True | False |
|---|-------------------------------------|-------------------------------------|
| S_1 is better than S_2 in terms of simplified Silhouette. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| S_2 is better than S_3 in terms of simplified Silhouette. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| S_1 and S_3 are equally good in terms of simplified Silhouette. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| S_3 is better than S_1 in terms of simplified Silhouette. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

(4 points)

EM-clustering - which of the following statements are true:

| | True | False |
|---|-------------------------------------|-------------------------------------|
| EM clustering makes use of Bayes' rule. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| EM clustering assumes independence between attributes. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| EM clustering is a generalization of k -means. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| The principle of EM clustering is to find the MAP hypothesis. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| Compared to k -means, EM has to fit more parameters with the same value of k . | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| Compared to k -means, EM works with stricter assumptions regarding the cluster distributions. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

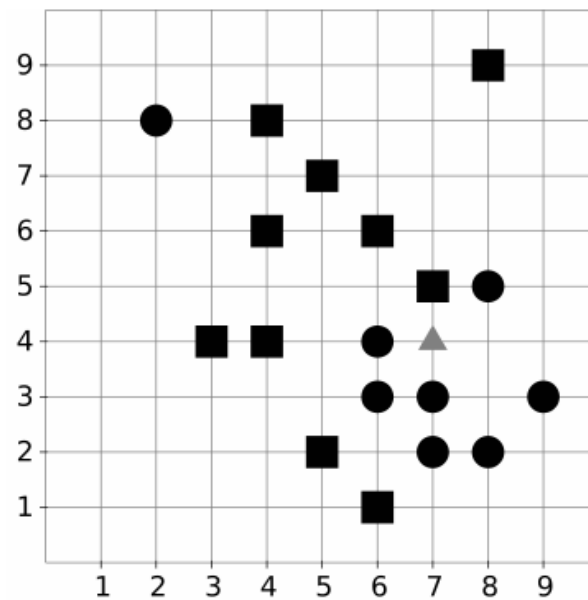
(4 points)

Which of the following statements are correct?

| | True | False |
|--|-------------------------------------|-------------------------------------|
| The number of parameters to describe a d -dimensional normal distribution grows quadratically with the number of dimensions. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| In 10-fold cross-validation, each object is used exactly ten times for training. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| A non-parametric clustering algorithm is an algorithm that does not require any user-specified parameters. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| The larger a decision tree grows, the more accurate will it be on the training data. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| k -means, k -nearest neighbor classifiers, and EM clustering - all these three methods are parametric learning methods. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |

(7 points)

We have the following two-dimensional dataset:



Using Manhattan distance, for which choices of k would the kNN classifier classify the query point (triangle) as square?

| | True | False |
|----|-------------------------------------|-------------------------------------|
| 3 | <input type="radio"/> | <input checked="" type="checkbox"/> |
| 6 | <input type="radio"/> | <input checked="" type="checkbox"/> |
| 10 | <input type="radio"/> | <input checked="" type="checkbox"/> |
| 13 | <input checked="" type="checkbox"/> | <input type="radio"/> |
| 15 | <input checked="" type="checkbox"/> | <input type="radio"/> |
| 16 | <input checked="" type="checkbox"/> | <input type="radio"/> |
| 17 | <input checked="" type="checkbox"/> | <input type="radio"/> |
| 18 | <input checked="" type="checkbox"/> | <input type="radio"/> |

(5 points)

Given the true class of 10 test objects and the predictions of some classifier h , which statements are correct w.r.t. to the class specific evaluation measures recall and precision?

| o | true class ($f(o)$) | prediction ($h(o)$) |
|----------|-----------------------|-----------------------|
| o_1 | A | A |
| o_2 | A | A |
| o_3 | A | C |
| o_4 | A | B |
| o_5 | B | B |
| o_6 | A | A |
| o_7 | B | B |
| o_8 | B | B |
| o_9 | B | C |
| o_{10} | C | B |
| o_{11} | C | A |
| o_{12} | C | C |

| | True | False |
|---|----------------------------------|----------------------------------|
| recall for class A $>$ recall for class B | <input type="radio"/> | <input checked="" type="radio"/> |
| precision for class A $>$ precision for class C | <input checked="" type="radio"/> | <input type="radio"/> |
| precision for class A $>$ recall for class A | <input checked="" type="radio"/> | <input type="radio"/> |
| precision for class B $>$ recall for class B | <input type="radio"/> | <input checked="" type="radio"/> |
| precision for class C $>$ recall for class B | <input type="radio"/> | <input checked="" type="radio"/> |
| recall for class A $>$ precision for class C | <input checked="" type="radio"/> | <input type="radio"/> |

(8 points)

We have a classification problem with two classes “+” and “−”, three trained classifiers h_1 , h_2 , and h_3 , with the following probabilities of the classifiers, given the training data D :

$$\Pr(h_1|D) = 0.2$$

$$\Pr(h_2|D) = 0.1$$

$$\Pr(h_3|D) = 0.7$$

For the three test instances o_1 , o_2 , o_3 , the classifiers give the following class probabilities:

| | |
|--------------------------|--------------------|
| $o_1 : \Pr(+ h_1) = 0.4$ | $\Pr(- h_1) = 0.6$ |
| $\Pr(+ h_2) = 0.4$ | $\Pr(- h_2) = 0.6$ |
| $\Pr(+ h_3) = 0.6$ | $\Pr(- h_3) = 0.4$ |
| $o_2 : \Pr(+ h_1) = 0.8$ | $\Pr(- h_1) = 0.2$ |
| $\Pr(+ h_2) = 0.4$ | $\Pr(- h_2) = 0.6$ |
| $\Pr(+ h_3) = 0.5$ | $\Pr(- h_3) = 0.5$ |
| $o_3 : \Pr(+ h_1) = 0.6$ | $\Pr(- h_1) = 0.4$ |
| $\Pr(+ h_2) = 0.8$ | $\Pr(- h_2) = 0.2$ |
| $\Pr(+ h_3) = 0.5$ | $\Pr(- h_3) = 0.5$ |

We combine the three classifiers to get a Bayes optimal classifier. Which of the following class probabilities will we get from this Bayes optimal classifier?

| | | True | False |
|-------------------------------------|----------|-------------------------------------|-------------------------------------|
| $o_1 : \Pr(+ \text{Bayes optimal})$ | $= 0.54$ | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| $o_1 : \Pr(- \text{Bayes optimal})$ | $= 0.45$ | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| $o_2 : \Pr(+ \text{Bayes optimal})$ | $= 0.45$ | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| $o_2 : \Pr(- \text{Bayes optimal})$ | $= 0.45$ | <input checked="" type="checkbox"/> | <input type="checkbox"/> |

$o_3 : \Pr(+|\text{Bayes optimal})$ $= 0.55$  $o_3 : \Pr(-|\text{Bayes optimal})$ $= 0.55$ 

(8 points)

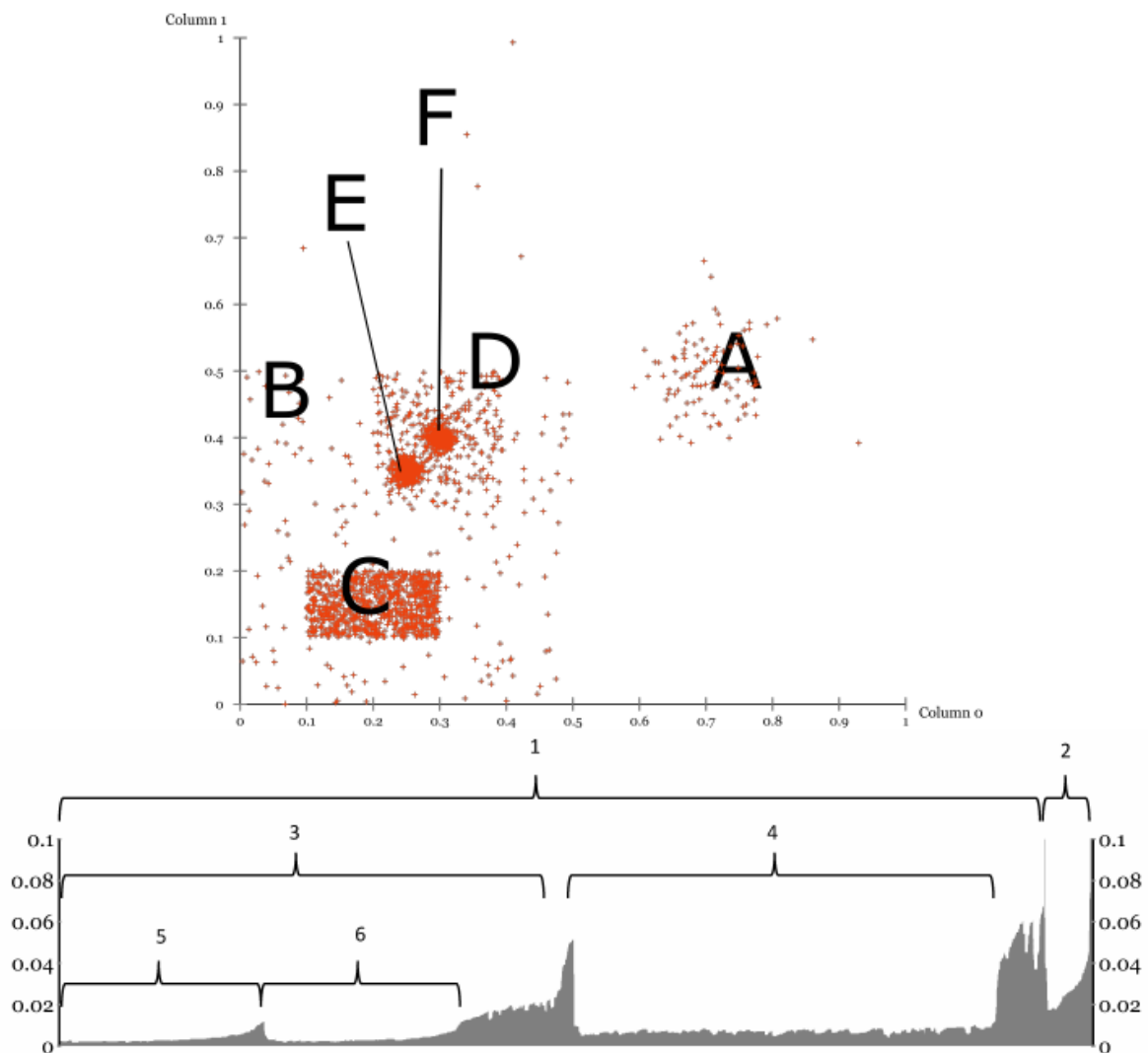
Navigating the society out of the lock-down, health inspector Mortensen shall give his opinion on allowing or forbidding several social activities. He takes orientation on the recommendations of medical, economical, and social science experts Olsen, Frandsen, and Jensen. As these three are not always in agreement, Mortensen wants to use a naïve Bayes classifier to combine their opinions. As training data, he uses the following observations from the past:

| activity | Olsen | Frandsen | Jensen | allow? |
|-------------------------------|--------|----------|--------|--------|
| <i>activity</i> ₁ | forbid | forbid | allow | yes |
| <i>activity</i> ₂ | forbid | forbid | forbid | yes |
| <i>activity</i> ₃ | forbid | forbid | allow | yes |
| <i>activity</i> ₄ | allow | forbid | forbid | yes |
| <i>activity</i> ₅ | allow | forbid | allow | yes |
| <i>activity</i> ₆ | allow | allow | forbid | yes |
| <i>activity</i> ₇ | allow | allow | allow | no |
| <i>activity</i> ₈ | forbid | allow | forbid | no |
| <i>activity</i> ₉ | forbid | allow | allow | no |
| <i>activity</i> ₁₀ | allow | forbid | forbid | no |

Which of the following activities would the classifier recommend to allow (=yes)?

| | | | | Yes | No |
|------------------------------|--------------|-----------------|---------------|--------------------------|-------------------------------------|
| activity | Olsen | Frandsen | Jensen | | |
| <i>activity</i> _A | allow | allow | forbid | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| activity | Olsen | Frandsen | Jensen | | |
| <i>activity</i> _B | allow | forbid | forbid | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| activity | Olsen | Frandsen | Jensen | | |
| <i>activity</i> _C | forbid | forbid | allow | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| activity | Olsen | Frandsen | Jensen | | |
| <i>activity</i> _D | forbid | allow | allow | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

(5 points)



Given the plotted dataset and the OPTICS plot – which statements are correct?

| | True | False |
|--|-------------------------------------|-------------------------------------|
| The area 1 in the OPTICS plot relates to cluster B. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| The area 2 in the OPTICS plot relates to cluster A. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| The area 3 in the OPTICS plot relates to cluster D. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| The area 4 in the OPTICS plot relates to cluster E or F. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| The area 5 in the OPTICS plot relates to cluster C. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

(8 points)

| ID | forecast | humidity | wind | play tennis? |
|----|----------|----------|--------|--------------|
| 1 | sunny | high | weak | no |
| 2 | sunny | high | strong | no |
| 3 | sunny | high | weak | yes |
| 4 | sunny | normal | weak | yes |
| 5 | sunny | normal | strong | no |
| 6 | rainy | high | weak | no |
| 7 | rainy | normal | weak | yes |
| 8 | rainy | normal | weak | yes |
| 9 | rainy | normal | strong | yes |
| 10 | rainy | high | strong | no |

A decision tree is being trained on the above data set. As root of the tree, the attribute “forecast” was already selected.

Which attributes are selected as test nodes at the next level based on the Gini index?

| | True | False |
|---|-------------------------------------|-------------------------------------|
| For the branch of forecast=sunny, we test wind. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| For the branch of forecast=sunny, we test humidity. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| For the branch of forecast=rainy, we test wind. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| For the branch of forecast=rainy, we test humidity. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |

(4 points)

Given the distance measure dist :

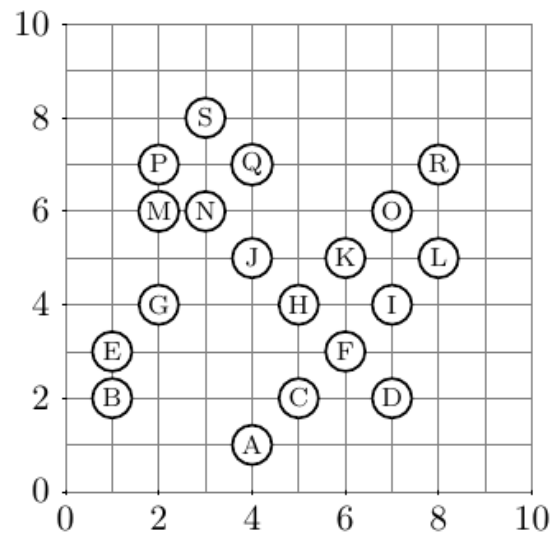
$$\text{dist}(x, y) = \sqrt{(x_1 - y_1, x_2 - y_2) \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix} (x_1 - y_1, x_2 - y_2)^\top}$$

for two-dimensional points and the point $p = (0, 2)$, which of the following points have the same distance as p from the origin $(0, 0)$?

| | True | False |
|-------------------------|-------------------------------------|-------------------------------------|
| $(4, 0)$ | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| $(\sqrt{8}, 0)$ | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| $(8, 0)$ | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| $(\sqrt{2}, 0)$ | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| $(\sqrt{2}, \sqrt{3})$ | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| $(1, 0)$ | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| $(0, -2)$ | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| $(-\sqrt{8}, \sqrt{2})$ | <input checked="" type="checkbox"/> | <input type="checkbox"/> |

(6 points)

Given the following dataset and Manhattan distance as distance function:

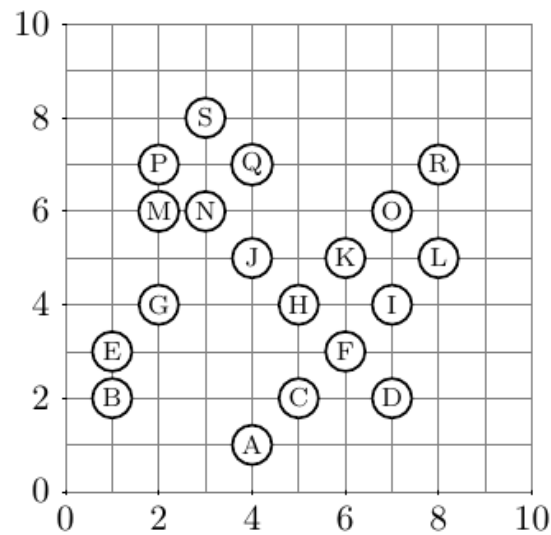


Given the definitions of DBSCAN and counting the query object as a member of its neighborhood: which of the following statements are correct?

| | True | False |
|--|-------------------------------------|-------------------------------------|
| A is core point with $\varepsilon = 2$ and $\text{MinPts} = 4$. | <input type="radio"/> | <input checked="" type="checkbox"/> |
| C is core point with $\varepsilon = 2$ and $\text{MinPts} = 4$. | <input checked="" type="checkbox"/> | <input type="radio"/> |
| G is core point with $\varepsilon = 2$ and $\text{MinPts} = 4$. | <input type="radio"/> | <input checked="" type="checkbox"/> |
| J is core point with $\varepsilon = 2$ and $\text{MinPts} = 4$. | <input checked="" type="checkbox"/> | <input type="radio"/> |
| M is core point with $\varepsilon = 2$ and $\text{MinPts} = 4$. | <input checked="" type="checkbox"/> | <input type="radio"/> |
| P is core point with $\varepsilon = 1$ and $\text{MinPts} = 2$. | <input checked="" type="checkbox"/> | <input type="radio"/> |
| Q is core point with $\varepsilon = 1$ and $\text{MinPts} = 4$. | <input type="radio"/> | <input checked="" type="checkbox"/> |
| R is core point with $\varepsilon = 1$ and $\text{MinPts} = 4$. | <input type="radio"/> | <input checked="" type="checkbox"/> |

(6 points)

Given the following dataset and Manhattan distance as distance function:

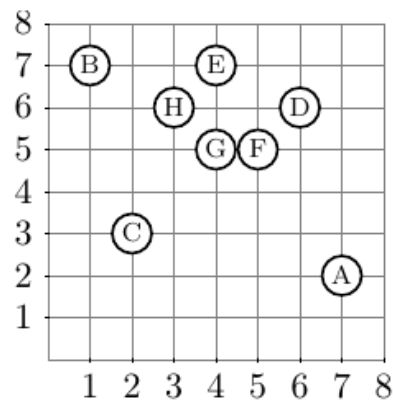


Given the definitions of DBSCAN and counting the query object as a member of its neighborhood: which of the following statements are correct?

| | True | False |
|---|-------------------------------------|-------------------------------------|
| P is directly density-reachable from N with $\varepsilon = 2$ and $\text{MinPts} = 4$. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| S is directly density-reachable from M with $\varepsilon = 2$ and $\text{MinPts} = 4$. | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| A and L are density-connected with $\varepsilon = 2$ and $\text{MinPts} = 2$. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| B and R are density-connected with $\varepsilon = 2$ and $\text{MinPts} = 2$. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |

(4 points)

Given the following dataset and Manhattan distance as distance function:



We have subsets of these points ordered by their outlier score (using the kNN-outlier score, the query point does not count to its own neighborhood here). For a correctly ordered subset of points, the top outlier among the selected points is listed first, later points follow with decreasing outlier score. Which of the following subsets are correctly ordered w.r.t. the given outlier method and parameter?

| | True | False |
|------------------------------|-------------------------------------|-------------------------------------|
| A,C,B w.r.t. kNN ($k = 2$) | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| C,B,E w.r.t. kNN ($k = 2$) | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| C,E,B w.r.t. kNN ($k = 2$) | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| F,E,B w.r.t. kNN ($k = 2$) | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

(4 points)

In a dataset with ten points $\{A, B, C, D, E, F, G, H, I, J\}$, A and B are labeled outliers.

Five outlier detection methods, m_1, \dots, m_5 , deliver the following rankings (from left-to-right: top-rank to bottom-rank):

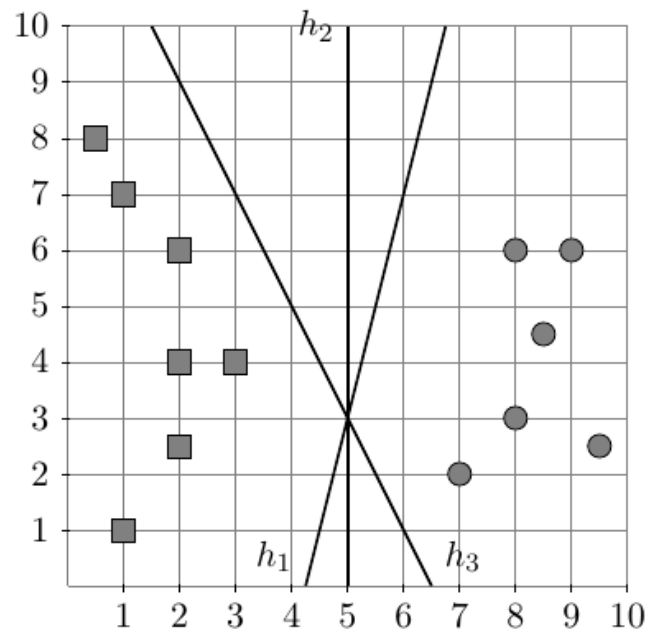
| method | ranking |
|--------|---------------------|
| m_1 | C,D,A,E,F,B,G,H,I,J |
| m_2 | J,A,D,E,F,G,B,H,I,C |
| m_3 | I,D,A,E,F,G,B,H,C,J |
| m_4 | I,J,E,A,B,F,G,H,C,D |
| m_5 | I,A,E,J,H,C,D,B,F,G |

Based on ROC AUC as evaluation measure, which of the following statements are correct?

| | True | False |
|---------------------------------------|-------------------------------------|-------------------------------------|
| m_1 and m_2 perform equally well. | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| m_1 is better than m_5 . | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| m_2 is better than m_3 . | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| m_2 is better than m_4 . | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| m_2 is better than m_5 . | <input checked="" type="checkbox"/> | <input type="checkbox"/> |
| m_3 is better than m_4 . | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| m_5 is better than m_3 . | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| m_5 is better than m_4 . | <input type="checkbox"/> | <input checked="" type="checkbox"/> |

(4 points)

The lines in this plot indicate the decision boundary derived by some classifier for the given training data:



In the following we conjecture, which classifier might have generated a given decision boundary. Which conjecture is possibly true?

| | True | False |
|--|-------------------------------------|-------------------------------------|
| h_1 : decision tree | <input type="radio"/> | <input checked="" type="checkbox"/> |
| h_1 : perceptron | <input checked="" type="checkbox"/> | <input type="radio"/> |
| h_1 : support vector machine (linear kernel) | <input checked="" type="checkbox"/> | <input type="radio"/> |
| h_2 : decision tree | <input checked="" type="checkbox"/> | <input type="radio"/> |
| h_2 : perceptron | <input type="radio"/> | <input checked="" type="checkbox"/> |
| h_2 : support vector machine (linear kernel) | <input checked="" type="checkbox"/> | <input type="radio"/> |
| h_3 : decision tree | <input type="radio"/> | <input checked="" type="checkbox"/> |
| h_3 : perceptron | <input checked="" type="checkbox"/> | <input type="radio"/> |
| h_3 : support vector machine (linear kernel) | <input checked="" type="checkbox"/> | <input type="radio"/> |

