

Technical University of Denmark

Written examination: 23 May 2017, 9 AM - 1 PM.

Course name: Introduction to Machine Learning and Data Mining.

Course number: 02450.

Aids allowed: All aids permitted.

Exam duration: 4 hours.

Weighting: The individual questions are weighted equally.

You must either use the electronic file or the form on this page to hand in your answers but not both. **We strongly encourage that you hand in your answers digitally using the electronic file.** If you hand in using the form on this page, please write your name and student number clearly.

The exam is multiple choice. All questions have four possible answers marked by the letters A, B, C, and D as well as the answer “Don’t know” marked by the letter E. Correct answer gives 3 points, wrong answer gives -1 point, and “Don’t know” (E) gives 0 points.

The individual questions are answered by filling in the answer fields with one of the letters A, B, C, D, or E.

Answers:

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20
21	22	23	24	25	26	27			

Name: _____

Student number: _____

PLEASE HAND IN YOUR ANSWERS DIGITALLY.

**USE ONLY THIS PAGE FOR HAND IN IF YOU ARE
UNABLE TO HAND IN DIGITALLY.**

No.	Attribute description	Abbrev.
x_1	Number of cylinders	cyl
x_2	Horsepower	hp
x_3	Weight	wt
x_4	Transmission (0=automatic, 1=manual)	am
x_5	Number of forward gears	gear
y	Miles pr. gallon	mpg

Table 1: The attributes of the Motor Trend Car Road Tests dataset taken from <https://vincentarelbundock.github.io/Rdatasets/csv/datasets/mtcars.csv>. The output y is given by the miles pr. gallon the car drives. The dataset has 32 observations and we presently consider the five input features x_1 – x_5 .

Question 1. We will consider the data of Motor Trend Car Road Tests based on 32 automobiles (observations) taken from <https://vincentarelbundock.github.io/Rdatasets/csv/datasets/mtcars.csv> for brevity denoted the Cars dataset in the following. The original data contains eleven attributes, however, we presently consider only five of these attributes given in Table 1 as well as the output attribute y given by how many miles the car drives pr. gallon of fuel.

Considering the attributes described in Table 1 which one of the following statements is *correct*?

- A. The attribute x_5 is continuous.
- B. The output variable y is ratio.
- C. The attribute x_4 is ordinal.
- D. The attribute x_1 is nominal.
- E. Don't know.

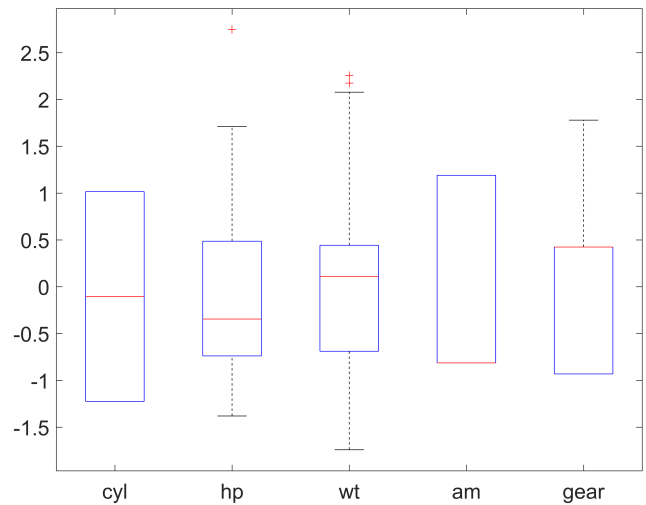


Figure 1: Boxplot of the five attributes x_1 – x_5 after standardizing the data (i.e., subtracting the mean of each attribute and dividing the attribute by its standard deviation).

Question 2. In Figure 1 is given a boxplot of the five attributes x_1 – x_5 after standardizing the data, i.e. subtracting the mean of each attribute and dividing each attribute by its standard deviation. Which one of the following statements is *correct*?

- A. The majority of cars have automatic transmission.
- B. The attribute x_5 (i.e., number of forward gears (gear)) appears to be normal distributed.
- C. The attribute x_2 (i.e., horse power (hp)) has a clear outlier that should be removed.
- D. From the boxplot it is clear that some of the attributes are highly correlated with each other.
- E. Don't know.

Question 3. A principal component analysis (PCA) is carried out on the standardized attributes x_1 – x_5 , forming the standardized matrix $\tilde{\mathbf{X}}$, resulting in the following \mathbf{S} and \mathbf{V} matrices obtained from a singular value decomposition of $\tilde{\mathbf{X}}$:

$$\mathbf{S} = \begin{bmatrix} 10.2 & 0 & 0 & 0 & 0 \\ 0 & 6.1 & 0 & 0 & 0 \\ 0 & 0 & 2.8 & 0 & 0 \\ 0 & 0 & 0 & 2.2 & 0 \\ 0 & 0 & 0 & 0 & 1.6 \end{bmatrix},$$

$$\mathbf{V} = \begin{bmatrix} 0.49 & -0.31 & 0.42 & -0.14 & 0.69 \\ 0.39 & -0.62 & 0.05 & -0.24 & -0.63 \\ 0.51 & -0.06 & -0.55 & 0.66 & 0.08 \\ -0.44 & -0.46 & 0.42 & 0.65 & -0.02 \\ -0.40 & -0.55 & -0.59 & -0.27 & 0.35 \end{bmatrix}.$$

The data projected onto the first two principal components are given in Figure 2. Which one of the following statements is *correct*?

- A. The first principal component accounts for less than 70 % of the variance.
- B. The two first principal components account for less than 90 % of the variance.
- C. The fifth principal component accounts for less than 1% of the variance.
- D. As can be observed in Figure 2 there is a positive correlation between the projection of the data to the first and second principal component.
- E. Don't know.

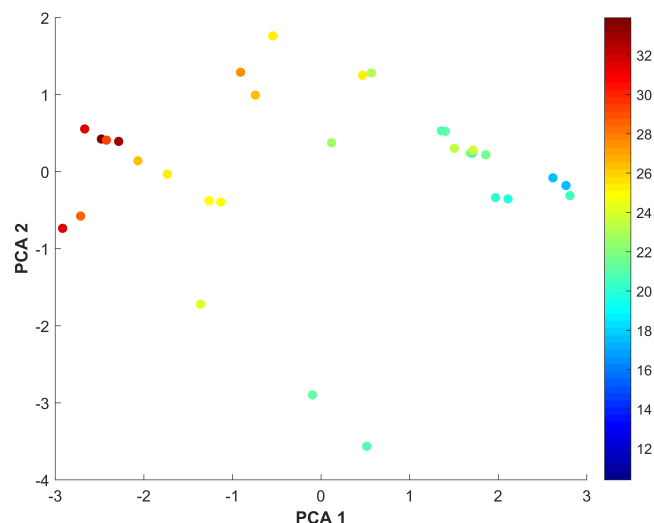


Figure 2: Data projected onto the first and second principal component. Each observation is color coded according to y , i.e. how many miles pr. gallon the car drives.

Question 4. The data projected onto the two first principal components (as defined in Question 3) is given in Figure 2 where the output variable y is indicated by the color of each observation. which one of the following statements pertaining to the PCA is *correct*?

- A. Cars with a relatively small number of cylinders, low horsepower, low weight, that have manual transmission, and many forward gears tend to drive longer pr. gallon of fuel.
- B. The second principal component appears to provide a better description of how far a car drives pr. gallon of fuel than principal component direction one.
- C. From the PCA plot it appears to be very difficult to predict fuel consumption based on the attributes x_1 – x_5 .
- D. Cars with a relatively small number of cylinders, low horsepower, that are automatic, and with few forward gears will have a large negative projection onto the second principal component.
- E. Don't know.

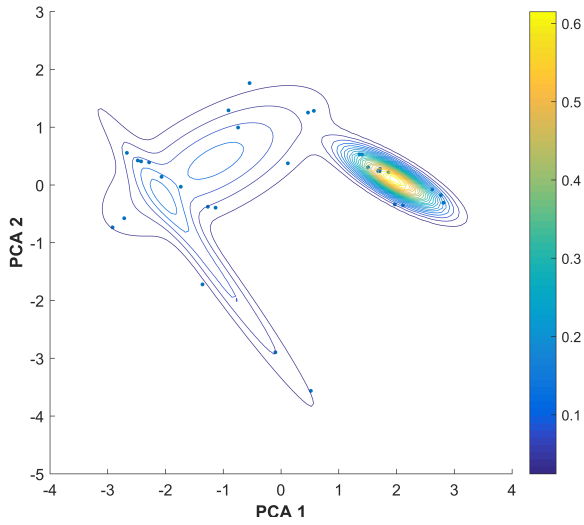


Figure 3: A Gaussian Mixture Model (GMM) with three clusters fitted to the Cars dataset projected onto the first two principal components.

Question 5. Consider the Gaussian Mixture Model (GMM) fitted to the Cars dataset projected onto the first two principal components for which the contours of the fitted GMM is given in Figure 3. We will in the following use:

$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{M/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$ to denote the multivariate normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Which one of the following GMM densities corresponds to the fitted GMM?

A.

$$\begin{aligned} p(\mathbf{x}) = & 0.449 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -1.129 \\ 0.417 \end{bmatrix}, \begin{bmatrix} 1.307 & 0.557 \\ 0.557 & 0.575 \end{bmatrix}\right) \\ & + 0.176 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -1.298 \\ -1.261 \end{bmatrix}, \begin{bmatrix} 0.248 & -0.128 \\ -0.128 & 0.101 \end{bmatrix}\right) \\ & + 0.375 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 1.957 \\ 0.091 \end{bmatrix}, \begin{bmatrix} 1.458 & -1.982 \\ -1.982 & 2.771 \end{bmatrix}\right) \end{aligned}$$

B.

$$\begin{aligned} p(\mathbf{x}) = & 0.449 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -1.129 \\ 0.417 \end{bmatrix}, \begin{bmatrix} 1.307 & 0.557 \\ 0.557 & 0.575 \end{bmatrix}\right) \\ & + 0.176 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -1.298 \\ -1.261 \end{bmatrix}, \begin{bmatrix} 1.458 & -1.982 \\ -1.982 & 2.771 \end{bmatrix}\right) \\ & + 0.375 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 1.957 \\ 0.091 \end{bmatrix}, \begin{bmatrix} 0.248 & -0.128 \\ -0.128 & 0.101 \end{bmatrix}\right) \end{aligned}$$

C.

$$\begin{aligned} p(\mathbf{x}) = & 0.449 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -1.129 \\ 0.417 \end{bmatrix}, \begin{bmatrix} 1.458 & -1.982 \\ -1.982 & 2.771 \end{bmatrix}\right) \\ & + 0.176 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -1.298 \\ -1.261 \end{bmatrix}, \begin{bmatrix} 1.307 & 0.557 \\ 0.557 & 0.575 \end{bmatrix}\right) \\ & + 0.375 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 1.957 \\ 0.091 \end{bmatrix}, \begin{bmatrix} 0.248 & -0.128 \\ -0.128 & 0.101 \end{bmatrix}\right) \end{aligned}$$

D.

$$\begin{aligned} p(\mathbf{x}) = & 0.449 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -1.129 \\ 0.417 \end{bmatrix}, \begin{bmatrix} 1.458 & -1.982 \\ -1.982 & 2.771 \end{bmatrix}\right) \\ & + 0.176 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -1.298 \\ -1.261 \end{bmatrix}, \begin{bmatrix} 0.248 & -0.128 \\ -0.128 & 0.101 \end{bmatrix}\right) \\ & + 0.375 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 1.957 \\ 0.091 \end{bmatrix}, \begin{bmatrix} 1.307 & 0.557 \\ 0.557 & 0.575 \end{bmatrix}\right) \end{aligned}$$

E. Don't know.

Question 6. A least squares linear regression model is trained using different combinations of the five attributes x_1 , x_2 , x_3 , x_4 , and x_5 . Table 2 gives the training and test root-mean-square error (RMSE= $\sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$) performance of the least squares linear regression model when trained using different combinations of the five attributes. Which one of the following statements is *correct*?

- A. Forward selection will select the same optimal feature set as backward selection.
- B. For this problem backward selection will identify the optimal feature combination.
- C. Forward selection will end up selecting all the features, i.e., terminate at the feature set x_1, x_2, x_3, x_4, x_5 .
- D. Forward selection will as first feature select x_3 .
- E. Don't know.

Feature(s)	Training RMSE	Test RMSE
x_1	3.2522	3.0343
x_2	3.1721	5.5072
x_3	2.907	3.4486
x_4	4.4608	5.8157
x_5	4.4637	9.0155
x_1 and x_2	3.043	3.8668
x_1 and x_3	2.4192	2.7692
x_1 and x_4	2.8918	3.3297
x_1 and x_5	3.2465	3.3177
x_2 and x_3	2.5325	2.5853
x_2 and x_4	2.8066	3.2509
x_2 and x_5	3.1646	4.6723
x_3 and x_4	2.8601	3.7105
x_3 and x_5	2.8230	4.3690
x_4 and x_5	3.9742	8.4346
x_1 and x_2 and x_3	2.4098	2.4869
x_1 and x_2 and x_4	2.7253	2.6258
x_1 and x_2 and x_5	3.0341	4.6252
x_1 and x_3 and x_4	2.3834	2.9486
x_1 and x_3 and x_5	2.3709	2.9676
x_1 and x_4 and x_5	2.7956	3.4894
x_2 and x_3 and x_4	2.4703	2.3423
x_2 and x_3 and x_5	2.5319	2.7017
x_2 and x_4 and x_5	2.7847	4.2531
x_3 and x_4 and x_5	2.8028	4.4864
x_1 and x_2 and x_3 and x_4	2.3679	2.5095
x_1 and x_2 and x_3 and x_5	2.3623	3.0234
x_1 and x_2 and x_4 and x_5	2.6249	4.3059
x_1 and x_3 and x_4 and x_5	2.2937	3.0251
x_2 and x_3 and x_4 and x_5	2.4561	2.8221
x_1 and x_2 and x_3 and x_4 and x_5	2.2759	3.1368

Table 2: Root-mean-square error (RMSE) for the training and test set when using least squares regression to predict fuel consumption of a car, i.e., mpg, using different combinations of the five attributes (x_1 – x_5).

Question 7. Using the 32 observations of the Cars dataset we would like to predict the fuel consumption of cars (y) based on the five features ($x_1 - x_5$). For this purpose we consider regularized least squares regression which minimizes with respect to \mathbf{w} the following cost function:

$$E(\mathbf{w}) = \sum_n (y_n - [1 \ x_{n1} \ x_{n2} \ x_{n3} \ x_{n4} \ x_{n5}] \mathbf{w})^2 + \lambda \mathbf{w}^\top \mathbf{w},$$

where x_{nm} denotes the m 'th feature of the n 'th observation, and 1 is concatenated the data to account for the bias term. We consider nine different values of λ and use leave-one-out cross-validation to quantify the performance of each of these different values of λ . The results of the leave-one-out cross-validation performance is given in Figure 4 where the optimal value of lambda is found to be $\lambda = 10^{-1.75}$ indicated with a black cross in the figure. Which one of the following statements is correct?

- A. The identified test performance having RMSE=2.8 is an unbiased estimator of how the optimal model identified will generalize to new data.
- B. To create the test error curve given in Figure 4 requires training 288 regularized least squares regression models.
- C. Leave-one-out cross-validation is computationally more efficient than 10-fold cross-validation.
- D. As λ increases the fitted models will have smaller and smaller bias.
- E. Don't know.

Question 8. We will build a model to determine if a car has a relatively high or low fuel consumption using the original data (i.e., the data is no longer standardized). For this purpose we will split the output y in two classes forming the new output variable z defined by thresholding at the median value of y , i.e. if $y_i > \text{median}(y)$ then $z_i = 1$, otherwise $z_i = 0$. We fit a logistic regression model using the features $x_1 - x_5$ and use as output for the logistic regression model the binary variable z indicating if a car has relatively high ($z = 0$) or low ($z = 1$) fuel consumption. The predicted output probability is given by:

$$\hat{z} = \sigma(1257.6 - 46.8x_1 + 0.6x_2 - 271.1x_3 - 31.9x_4 - 44.7x_5),$$

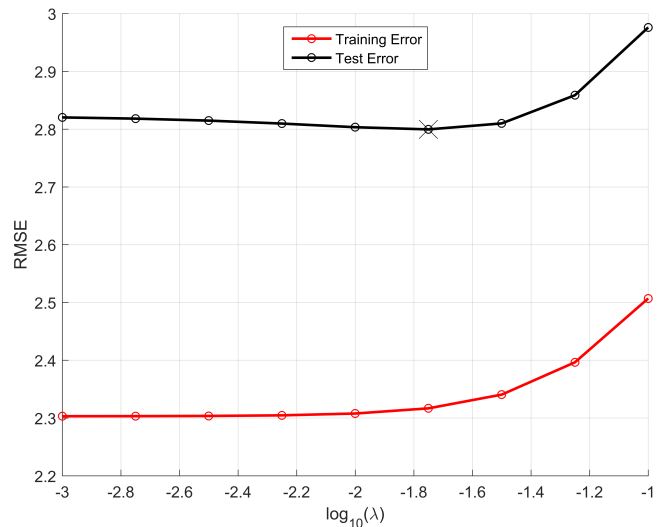


Figure 4: Regularized least squares regression applied to the Cars data in order to predict fuel consumption. Given is training and test performance as quantified by the root mean square error (RMSE) in red and black respectively as a function of the value of λ based on leave-one-out cross-validation.

where $\sigma(\cdot)$ is the logistic sigmoid function. Which one of the following statements regarding the logistic regression model is correct?

- A. x_2 is the least important attribute for defining whether a car has high (i.e., $z = 0$) or low (i.e., $z = 1$) fuel consumption.
- B. According to the model increasing a car's number of forward gears will make it more likely to have high fuel consumption (i.e., $z = 0$).
- C. A new car with the following observation vector: $\mathbf{x}^* = [6 \ 120 \ 3.2 \ 0 \ 4]$ will be more likely to have high (i.e., $z = 0$) than low (i.e., $z = 1$) fuel consumption.
- D. Logistic regression is not very suited for classification as it is a regression method.
- E. Don't know.

	3 gears ($x_5 = 3$)	4 gears ($x_5 = 4$)	5 gears ($x_5 = 5$)
Low mpg ($z = 0$)	13	2	2
High mpg ($z = 1$)	2	10	3

Table 3: Number of low mpg and high mpg cars (i.e. $z = 0$ and $z = 1$) according to the number of gears, i.e. $x_5 = 3$, $x_5 = 4$, or $x_5 = 5$.

Question 9. In order to improve the performance of the logistic regression model we will use ensembling based on the Adaboost algorithm using the 32 observations of the Cars dataset (note that for the Adaboost algorithm $\log(\cdot)$ is based on the natural logarithm). The first trained logistic regression classifier (i.e., for boosting round $t = 1$) has an error rate of $1/16$. What will be the updated weight for each of the correctly classified observations?

- A. 0.0081
- B. 0.0167
- C. 0.0332
- D. 0.2500
- E. Don't know.

Question 10. A decision tree is fitted to the data considering as output whether the car has a relatively high (i.e., $z = 0$) or low (i.e., $z = 1$) fuel consumption. At the root of the tree three different splits according to the number of forward gears are considered based on the data given in Table 3. For impurity we will use the classification error given by $I(v) = 1 - \max_c p(c|v)$. For the three considered splits we have:

- Split A: 3 gear vs. 4 or 5 gears.
- Split B: 3 or 4 gears vs. 5 gears.
- Split C: 3 gears vs. 4 gears vs. 5 gears.

Thus, split A and B have two branches whereas split C has three branches. Which statement about the splits is correct?

- A. Split B provides a higher purity gain than split A.
- B. Split C provides a higher purity gain than split A.
- C. The best obtainable purity gain is $9/32$.
- D. Split B provides a higher purity gain than split C.
- E. Don't know.

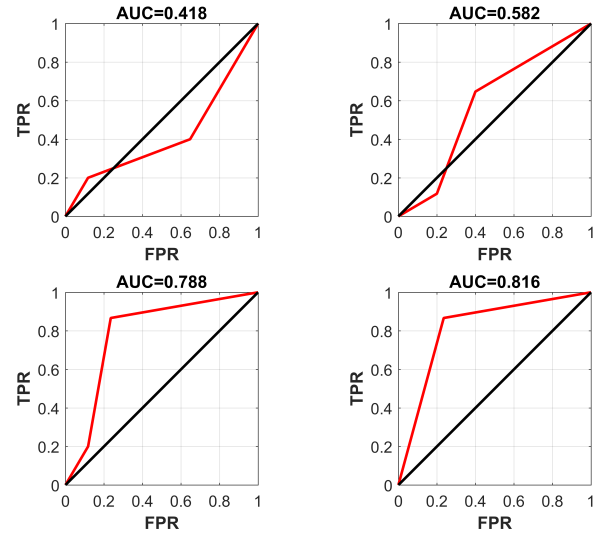


Figure 5: Four different receiver operator characteristic (ROC) curves and their area under curve (AUC) value.

Question 11. We will evaluate the feature x_5 (gear) in its ability to discriminate low mpg, i.e., $z = 0$, (considered the negative class) from high mpg, i.e. $z = 1$, (considered the positive class) based on the data given in Table 3. For this purpose, we will evaluate the area under curve (AUC) of the receiver operator characteristic (ROC) using the feature x_5 to discriminate cars with high (i.e., $z = 0$) from cars with low (i.e., $z = 1$) fuel consumption. Which one of the ROC curves given in Figure 5 corresponds to using x_5 to discriminate between high fuel consumption ($z = 0$) and low fuel consumption ($z = 1$)?

- A. The curve having AUC=0.418
- B. The curve having AUC=0.582
- C. The curve having AUC=0.788
- D. The curve having AUC=0.816
- E. Don't know.

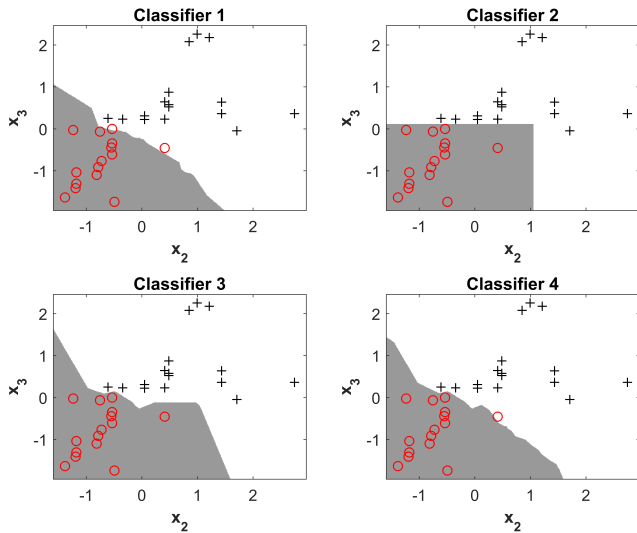


Figure 6: Decision boundaries for four different classifiers trained on the Cars dataset using only the two features x_2 and x_3 .

Question 12. Four different classifiers are trained on the Cars dataset using only x_2 and x_3 as features (the features have been standardized) and the decision boundary for each of the four classifiers is given in Figure 6. Which one of the following statements is *correct*?

- A. Classifier 1 is an artificial neural network with one hidden unit in the hidden layer.
- B. Classifier 2 is a 3-nearest neighbor classifier.
- C. Classifier 3 is a 1-nearest neighbor classifier.
- D. Classifier 4 is a logistic regression classifier.
- E. Don't know.

Question 13. In an attempt to make a stronger classifier the first three classifiers (i.e., Classifier 1, Classifier 2, and Classifier 3) are combined using majority voting. Which one of the decision boundaries given in Figure 7 corresponds to the combined classifier?

- A. Decision boundary A.
- B. Decision boundary B.
- C. Decision boundary C.
- D. Decision boundary D.
- E. Don't know.

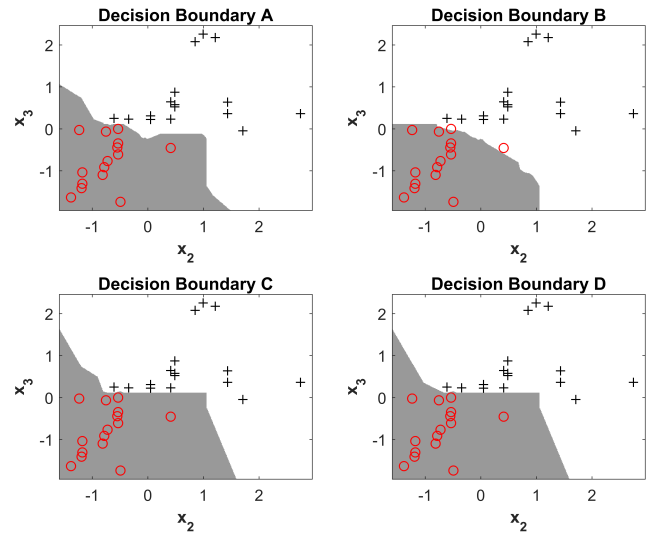


Figure 7: Decision boundaries for which one of the four decision boundaries corresponds to combining Classifier 1, Classifier 2, and Classifier 3 in Figure 6 using majority voting.

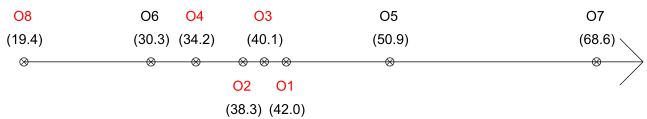


Figure 8: The eight first observations of the Cars dataset considered in regards to the feature $q = x_2/x_3$ (the value of q is given in parenthesis).

Question 14. We will consider the first eight observations of the Cars dataset and the new feature defined by the ratio of horse power to weight, i.e. defined by $q = x_2/x_3$. In Figure 8 is shown the value of q for the first eight observations that are colored according to low (red) and high (black) fuel consumption. We will cluster this data using k-means with Euclidean distance into two clusters (i.e., $k=2$) and initialize the k-means algorithm with centroids located at observation O8, and O6. Which one of the following statements is *correct*?

- A. The converged solution will be {O8}, {O1, O2, O3, O4, O5, O6, O7}.
- B. The converged solution will be {O1, O2, O3, O4, O6, O8}, {O5, O7}.
- C. The converged solution will be {O1, O2, O3, O4, O5, O6, O8}, {O7}.
- D. The converged solution will be {O4, O6, O8}, {O1, O2, O3, O5, O7}.
- E. Don't know.

Question 15. We will consider a clustering given by the two clusters $\{O2, O3, O4, O6, O8\}, \{O1, O5, O7\}$. We will evaluate this clustering in terms of its correspondence with the class label information in which O1,O2,O3,O4, and O8 correspond to cars with low fuel consumption and O5, O6, and O7 correspond to cars with high fuel consumption. We recall that the Jaccard coefficient between the true labels and the extracted clusters is given by:

$$J = \frac{f_{11}}{K - f_{00}},$$

where f_{11} is the number of object pairs in same class assigned to same cluster, f_{00} is the number of object pairs in different class assigned to different clusters, and $K = N(N - 1)/2$ is the total number of object pairs, where N is the number of observations considered. What is the above value of J between the true labeling of the observations in terms of high and low fuel consumption and the two clusters?

- A. 0.1665
- B. 0.3684
- C. 0.5714
- D. 0.7500
- E. Don't know.

Question 16. According to the Cars dataset we have that 59.38% of the cars have automatic transmission. Furthermore, 63.16% of the cars that have automatic transmission have eight cylinders, whereas 15.38% of the cars that have manual transmission have eight cylinders. According to the Cars dataset what is the probability that a car that has eight cylinders, i.e. $x_1 = 8$ will have automatic transmission, i.e. $x_4 = 0$?

- A. 16.7 %
- B. 37.5 %
- C. 63.2%
- D. 85.7 %
- E. Don't know.

Question 17. We will consider an artificial neural network (ANN) trained to predict mpg (i.e., y) of a car. The ANN is based on the model:

$$f(\mathbf{x}, \mathbf{w}) = \sum_{j=1}^2 w_j^{(2)} h^{(1)}([1 \ \mathbf{x}] \mathbf{w}_j^{(1)}) + w_0^{(2)}.$$

where $h^{(1)}(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ is the hyperbolic tangent function used as activation function in the hidden layer. We will consider an ANN with two hidden units in the hidden layer defined by:

$$\mathbf{w}_1^{(1)} = \begin{bmatrix} -4 \\ 1 \\ 0.01 \\ 1 \\ -1 \\ -1 \end{bmatrix}, \quad \mathbf{w}_2^{(1)} = \begin{bmatrix} -10 \\ 1 \\ -0.02 \\ 1 \\ 1 \\ 1 \end{bmatrix},$$

and $w_0^{(2)} = 7$, $w_1^{(2)} = 8$, and $w_2^{(2)} = 9$.

What is the predicted fuel consumption of a car with observation vector $\mathbf{x}^* = [6 \ 120 \ 3.2 \ 0 \ 4]$?

- A. 17.00
- B. 20.85
- C. 24.00
- D. 33.40
- E. Don't know.

	hp _L	hp _H	wt _L	wt _H	am=0	am=1
O1	1	0	1	0	0	1
O2	1	0	1	0	0	1
O3	1	0	1	0	0	1
O4	1	0	1	0	1	0
O5	0	1	0	1	1	0
O6	1	0	0	1	1	0
O7	0	1	0	1	1	0
O8	1	0	1	0	1	0

Table 4: The first eight observations of the Cars dataset binarized considering the attribute x_2 , x_3 , and x_4 such that x_2 is split according to the median value in terms of low and high horse power (i.e. hp_L and hp_H), low and high weight (i.e. wt_L and wt_H), as well as whether the car has automatic (i.e., am=0) or manual (i.e., am=1) transmission. The eight observations are color coded in terms of low fuel consumption {O1, O2, O3, O4, O8} and high fuel consumption {O5, O6, O7}.

Question 18. Considering the dataset in Table 4 as a market basket problem with observation O1–O8 corresponding to customers and hp_L, hp_H, wt_L, wt_H, am=0, am=1 corresponding to items, what is the confidence of the association rule $\{wt_H, am=0\} \rightarrow \{hp_H\}$?

- A. 0.0%
- B. 25.0
- C. 66.7%
- D. 100.0 %
- E. Don't know.

Question 19. We will again consider the data in Table 4. What are all frequent itemsets with support greater than 30%?

- A. {hp_L}, {wt_L}, {wt_H}, {am=0}, {am=1}.
- B. {hp_L}, {wt_L}, {wt_H}, {am=0}, {am=1}, {hp_L, wt_L}, {hp_L, am=0}, {hp_L, am=1}, {wt_L, am=1}, {wt_H, am=0}.
- C. {hp_L}, {wt_L}, {wt_H}, {am=0}, {am=1}, {hp_L, wt_L}, {hp_L, am=0}, {hp_L, am=1}, {wt_L, am=1}, {wt_H, am=0}, {hp_L, wt_L, am=1},
- D. {hp_L}, {wt_L}, {wt_H}, {am=0}, {am=1}, {hp_L, wt_L}, {hp_L, am=0}, {hp_L, am=1}, {wt_L, am=1}, {wt_H, am=0}, {hp_L, wt_L, am=1}, {hp_L, wt_L, am=0},

- E. Don't know.

Question 20. Considering the data in Table 4, we will calculate the similarity as well as distance between O1 given by the vector $\mathbf{a} = [1 \ 0 \ 1 \ 0 \ 0 \ 1]$ and O4 given by the vector $\mathbf{b} = [1 \ 0 \ 1 \ 0 \ 1 \ 0]$ using respectively the Jaccard (J), simple matching coefficient (SMC), cosine similarity (cos) and p-norm ($\|\cdot\|_p$) given by

$$J(\mathbf{a}, \mathbf{b}) = \frac{f_{11}}{M - f_{00}},$$

$$SMC(\mathbf{a}, \mathbf{b}) = \frac{f_{11} + f_{00}}{M},$$

$$\cos(\mathbf{a}, \mathbf{b}) = \frac{f_{11}}{\|\mathbf{a}\|_2 \|\mathbf{b}\|_2},$$

$$\|\mathbf{a} - \mathbf{b}\|_p = \left(\sum_{m=1}^M |a_m - b_m|^p \right)^{1/p}.$$

Which one of the following statements is correct?

- A. $\|\mathbf{a} - \mathbf{b}\|_2 = 2$.
- B. $\|\mathbf{a} - \mathbf{b}\|_1 < \|\mathbf{a} - \mathbf{b}\|_2$.
- C. $J(\mathbf{a}, \mathbf{b}) = 2/3$
- D. $\cos(\mathbf{a}, \mathbf{b}) = SMC(\mathbf{a}, \mathbf{b})$.
- E. Don't know.

Question 21. We would like to predict whether a car has high fuel or low fuel consumption using the data in Table 4. We will apply a Naïve Bayes classifier that assumes independence between the attributes given the class label (i.e. high or low fuel consumption corresponding to observation indicated in black and red respectively in the table). Given that a car has hp_L = 1, and am=0 what is the probability that the car will have high fuel consumption according to the Naïve Bayes classifier derived from the data in Table 4?

- A. 3/80
- B. 1/8
- C. 1/3
- D. 3/5
- E. Don't know.

	O1	O2	O3	O4	O5	O6	O7	O8
O1	0	0.2606	1.1873	2.4946	2.9510	2.5682	3.4535	2.4698
O2	0.2606	0	1.2796	2.4442	2.8878	2.4932	3.3895	2.4216
O3	1.1873	1.2796	0	2.8294	3.6892	2.9147	4.1733	2.2386
O4	2.4946	2.4442	2.8294	0	1.4852	0.2608	2.2941	1.8926
O5	2.9510	2.8878	3.6892	1.4852	0	1.5155	1.0296	3.1040
O6	2.5682	2.4932	2.9147	0.2608	1.5155	0	2.3316	1.8870
O7	3.4535	3.3895	4.1733	2.2941	1.0296	2.3316	0	3.7588
O8	2.4698	2.4216	2.2386	1.8926	3.1040	1.8870	3.7588	0

Table 5: Pairwise Euclidean distance between the first eight observations in the Cars dataset. Red observations (i.e., O1, O2, O3, O4, O8) are observations corresponding to low fuel consumption, whereas black observations (i.e., O5, O6, O7) are observations with high fuel consumption.

Question 22. To determine whether the fuel consumption of a car is high or low we will use a k-nearest neighbor (KNN) classifier to predict each of the eight observations based on the Euclidean distance between the observations given in Table 5. We will use leave-one-out cross-validation for the KNN in order to classify the eight considered observations using a one-nearest neighbor classifier, i.e. $K = 1$. The analysis will be based only on the data given in Table 5. Which one of the following statements is *correct*?

- A. None of the observations will be misclassified.
- B. One of the observations will be misclassified.
- C. Two the observations will be misclassified.
- D. Three of the observations will be misclassified.
- E. Don't know.

Question 23. In Table 5 is given the pairwise Euclidean distances between the first eight observations of the Cars data. A hierarchical clustering is used to cluster these observations using average linkage. Which one of the dendrograms given in Figure 9 corresponds to the clustering?

- A. Dendrogram 1.
- B. Dendrogram 2.
- C. Dendrogram 3.
- D. Dendrogram 4.
- E. Don't know.

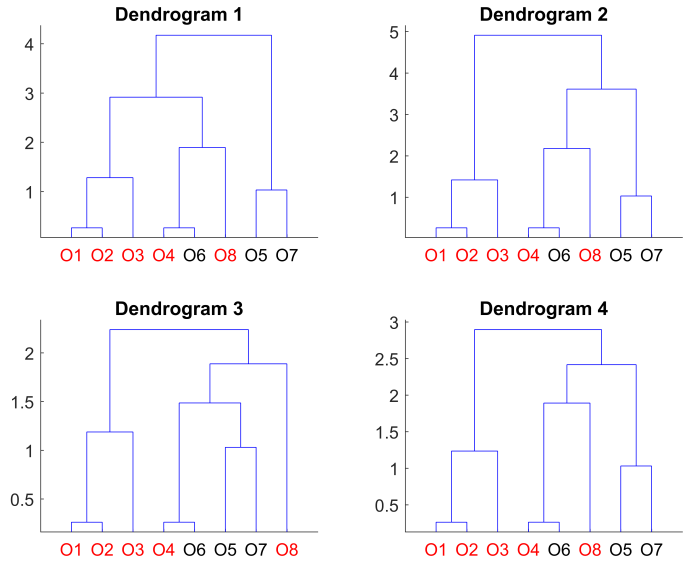


Figure 9: Four different dendrograms derived from the distances between the eight first cars given in Table 5.

Question 24. We suspect that observation O8 may be an outlier. In order to assess if this is the case we would like to calculate the average relative KNN density based on the observations given in Table 5 only. We recall that the KNN density and average relative density (ard) for the observation \mathbf{x}_i are given by:

$$\text{density}_{\mathbf{X}_{\setminus i}}(\mathbf{x}_i, K) = \frac{1}{\frac{1}{K} \sum_{\mathbf{x}' \in N_{\mathbf{X}_{\setminus i}}(\mathbf{x}_i, K)} d(\mathbf{x}_i, \mathbf{x}')},$$

$$\text{ard}_{\mathbf{X}}(\mathbf{x}_i, K) = \frac{\text{density}_{\mathbf{X}_{\setminus i}}(\mathbf{x}_i, K)}{\frac{1}{K} \sum_{\mathbf{x}_j \in N_{\mathbf{X}_{\setminus i}}(\mathbf{x}_i, K)} \text{density}_{\mathbf{X}_{\setminus j}}(\mathbf{x}_j, K)},$$

where $N_{\mathbf{X}_{\setminus i}}(\mathbf{x}_i, K)$ is the set of K nearest neighbors of observation \mathbf{x}_i excluding the i 'th observation, and $\text{ard}_{\mathbf{X}}(\mathbf{x}_i, K)$ is the average relative density of \mathbf{x}_i using K nearest neighbors. Based on the data in Table 5, what is the average relative density for observation O8 for $K = 2$ nearest neighbors?

- A. 0.4660
- B. 0.4800
- C. 0.5292
- D. 1.8898
- E. Don't know.

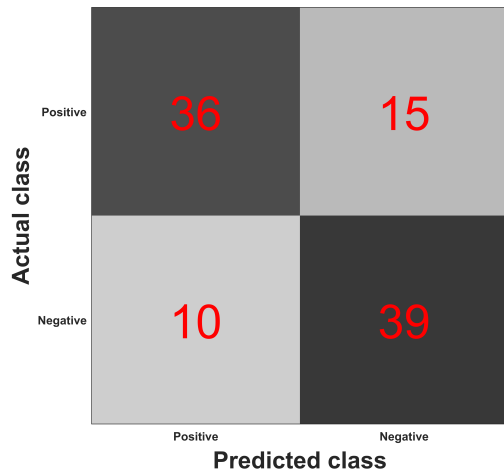


Figure 10: Confusion matrix of a classifier discriminating between 100 positive and negative test observations.

Question 25. We will consider a classifier classifying a dataset with 100 test observations into two classes (positive and negative) with confusion matrix given in Figure 10. Which statement regarding the classifier is correct?

- A. The error rate of the classifier is 33.3 %.
- B. The precision of the classifier is 75.0 %.
- C. The recall of the classifier is 70.6 %.
- D. There are more negative than positive examples in the test set.
- E. Don't know.

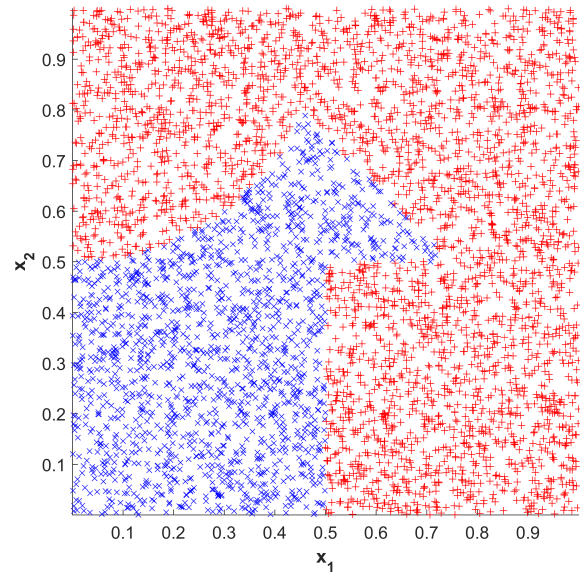


Figure 11: A two class classification problem with red plusses (i.e., +) and blue crosses (i.e., x) constituting the two classes.

Question 26. We will consider the two class classification problem given in Figure 11 in which the goal is to separate red plusses (i.e., +) from blue crosses (i.e., x). Which one of the following procedures will perfectly separate the two classes?

- A. $\|\mathbf{x} - \begin{bmatrix} 0 \\ 1 \end{bmatrix}\|_2 > 0.5$ and $\|\mathbf{x} - \begin{bmatrix} 1 \\ 0 \end{bmatrix}\|_\infty > 0.5$ and $\|\mathbf{x} - \begin{bmatrix} 1 \\ 1 \end{bmatrix}\|_1 > 0.75$ then blue cross, otherwise red plus.
- B. $\|\mathbf{x} - \begin{bmatrix} 0 \\ 1 \end{bmatrix}\|_\infty > 0.5$ and $\|\mathbf{x} - \begin{bmatrix} 1 \\ 0 \end{bmatrix}\|_2 > 0.5$ and $\|\mathbf{x} - \begin{bmatrix} 1 \\ 1 \end{bmatrix}\|_1 > 0.75$ then blue cross, otherwise red plus.
- C. $\|\mathbf{x} - \begin{bmatrix} 0 \\ 1 \end{bmatrix}\|_2 > 0.5$ and $\|\mathbf{x} - \begin{bmatrix} 1 \\ 0 \end{bmatrix}\|_1 > 0.5$ and $\|\mathbf{x} - \begin{bmatrix} 1 \\ 1 \end{bmatrix}\|_\infty > 0.75$ then blue cross, otherwise red plus.
- D. $\|\mathbf{x} - \begin{bmatrix} 0 \\ 1 \end{bmatrix}\|_\infty > 0.5$ and $\|\mathbf{x} - \begin{bmatrix} 1 \\ 0 \end{bmatrix}\|_1 > 0.5$ and $\|\mathbf{x} - \begin{bmatrix} 1 \\ 1 \end{bmatrix}\|_\infty > 0.75$ then blue cross, otherwise red plus.
- E. Don't know.

Question 27. Which one of the following statements is correct?

- A. Unsupervised learning differs from supervised learning in that unsupervised learning both uses the input data and the outputs for training whereas supervised learning only uses the input data.
- B. When using Gaussian Mixture Models (GMM) for outlier detection it is important that the observations evaluated for being outliers are included in the training of the GMM.
- C. When training an artificial neural network for a dataset with very few observations it is important to include many hidden units in order to avoid overfitting.
- D. Cross-validation can both be used for supervised and unsupervised learning.
- E. Don't know.