

Technical University of Denmark

Written examination: 16 December 2016, 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
A	A	B	C	D	B	A	C	A	B
11	12	13	14	15	16	17	18	19	20
D	D	A	C	D	D	D	B	A	B
21	22	23	24	25	26	27			
B	D	B	D	D	C	B			

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	Area	A
x_2	Perimeter	P
x_3	Length of kernel	L
x_4	Width of kernel	W
y	Seed type	

Table 1: The attributes of the Seeds data set taken from <http://archive.ics.uci.edu/ml/datasets/seeds>. The output is given by the type of seed, i.e. $y=1$ corresponds to Kama, $y=2$ corresponds to Rosa, and $y=3$ corresponds to Canadian.

Question 1. We will consider the data of wheat kernels based on 70 observations of each class of three seed types, i.e., Kama, Rosa, and Canadian. The original data contains seven attributes, however, we presently only consider four of these attributes given in Table 1. Considering the attributes described in the table and visualized using boxplots in Figure 1 which one of the following statements is *correct*?

- A. All the attributes x_1 , x_2 , x_3 , and x_4 are continuous and ratio.
- B. The output variable y is ordinal.
- C. Rosa and Canadian do not appear to differ in terms of area (A).
- D. The observations pertaining to Kama appear to contain clear outliers that must be removed.
- E. Don't know.

Solution 1. As zero means absence of the attribute for all the attributes, i.e. zero area means no area etc. and it makes sense to talk about an attribute value being twice as large etc. than another attribute value whereas all the values are continuous, the first answer is correct. y is nominal indicating class, but not ordinal we in general cannot argue that one type of seed is better/higher than another, only whether a seed is different or not from another. Indeed it appears from the boxplot that Rosa and Canadian differ in terms of Area such that all Rosa seeds have larger area than Canadian seeds. Finally, although boxplots indicate observations that fall beyond the whiskers these observations should not be removed unless there are strong justifications for doing so which the boxplot does not provide reasons for.

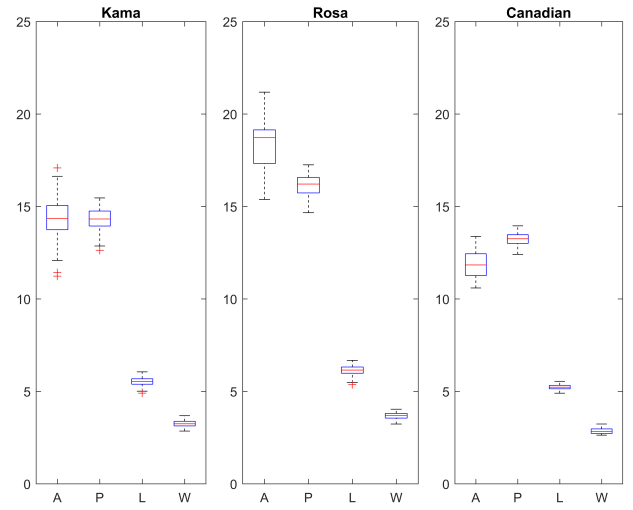


Figure 1: Boxplot of the data visualized separately for each of the three types of seeds; Kama, Rosa, and Canadian.

Question 2. A principal component analysis (PCA) is carried out on the standardized attributes x_1-x_4 , forming the standardized matrix \tilde{X} , resulting in the following S and V matrices obtained from a singular value decomposition:

$$S = \begin{bmatrix} 28.4 & 0 & 0 & 0 \\ 0 & 5.5 & 0 & 0 \\ 0 & 0 & 1.2 & 0 \\ 0 & 0 & 0 & 0.5 \end{bmatrix},$$

$$V = \begin{bmatrix} -0.51 & 0.11 & -0.39 & -0.76 \\ -0.51 & -0.13 & -0.58 & 0.62 \\ -0.49 & -0.69 & 0.53 & -0.05 \\ -0.49 & 0.71 & 0.47 & 0.19 \end{bmatrix}.$$

Which one of the following statements is *correct*?

- A. The first principal component accounts for more than 95 % of the variance.
- B. The two first principal components account for more than 99.9 % of the variance.
- C. The fourth principal component accounts for more than 0.05% of the variance.
- D. The attributes are not correlated as the data has been standardized.
- E. Don't know.

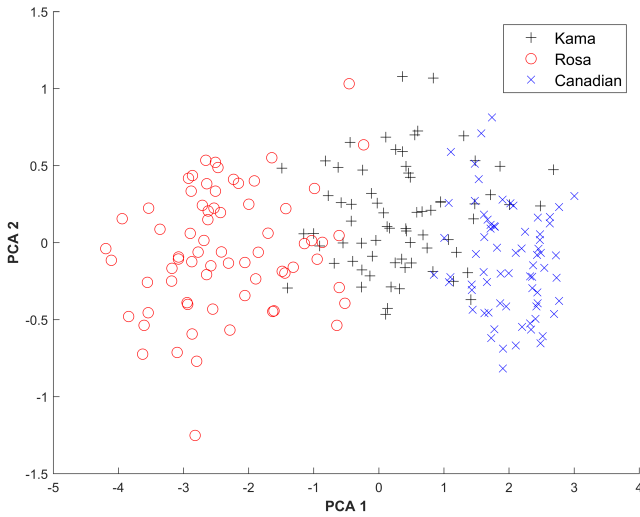


Figure 2: Data projected onto the first and second principal components.

Solution 2. The variation explained by each principal component is given by $\frac{\sigma_i^2}{\sum_{i'} \sigma_{i'}^2}$. As such we find:

$$VarExpPC1 = \frac{28.4^2}{28.4^2 + 5.5^2 + 1.2^2 + 0.5^2} = 0.9619 \quad (1)$$

$$VarExpPC2 = \frac{5.5^2}{28.4^2 + 5.5^2 + 1.2^2 + 0.5^2} = 0.0361 \quad (2)$$

$$VarExpPC3 = \frac{1.2^2}{28.4^2 + 5.5^2 + 1.2^2 + 0.5^2} = 0.0017 \quad (3)$$

$$VarExpPC4 = \frac{0.5^2}{28.4^2 + 5.5^2 + 1.2^2 + 0.5^2} = 0.0003 \quad (4)$$

As such the first PC accounts for more than 95% of the variance, the first two principal components accounts for $0.9619 + 0.0361 = 0.9980$ which is less than 99.9% of the variance. The fourth principal component accounts for 0.03% which is less than 0.05%. As the first principal component accounts for more than 95% of the variance the attributes are indeed very correlated and not the opposite.

Question 3. The data projected onto the two first principal components (as defined in Question 2) is given in Figure 2 where each class is indicated using different markers and colors. which one of the following statements pertaining to the PCA is *correct*?

- A. A relatively long and narrow seed kernel will provide a large positive projection onto the second principal component.
- B. The first principal component pertains to the general size of seeds.**
- C. A seed that has relatively small area and perimeter but large length and width of kernel will have a negative projection onto the third principal component.
- D. As the third and fourth principal components account for a low amount of the variance in the data this is a difficult classification task.
- E. Don't know.

Solution 3. As we for the second principal component have $\mathbf{v}_2^\top = [0.11 \ -0.13 \ -0.69 \ 0.71]$ a relatively long (large positive x_3) and relatively narrow (large negative x_4) will have a large negative projection onto this component. As the coefficients for the first principal component all are negative and generally have same magnitude for the four attributes, this appears to capture the general property of size of the seed, such that relatively large area, perimeter, kernel length and width will provide a negative projection and vice versa. The third principal component is defined by $\mathbf{v}_3^\top = [-0.39 \ -0.58 \ 0.53 \ 0.47]$, thus relatively small area and periphery (negative x_1 and x_2) but large kernel (i.e. positive x_3 , and x_4) will have a positive projection onto this component. The PCA does not take information about the classes into account and therefore does not necessarily reflect features relevant for classification. In particular, the singular values are uninformed by the classes as this information is not available in the PCA analysis.

Question 4. A decision tree is fitted to the data projected onto the four principal components. At the root of the tree a split according to the projection of the standardized data onto the first principal component being larger than 0 is considered, i.e. $\tilde{x}_n v_1 \geq 0$. For impurity we will use the classification error given by $I(v) = 1 - \max_c p(c|v)$. Before the split we have 70 Kama, 70 Rosa, and 70 Canadian and after the split:

- 24 Kama, 70 Rosa, 0 Canadian below zero in the projection onto v_1 .
- 46 Kama, 0 Rosa, 70 Canadian above or equal to zero in the projection onto v_1 .

What is the purity gain of this split?

- A. -1.0148
- B. 0.0148
- C. 0.3333**
- D. 0.6666
- E. Don't know.

Solution 4. The purity gain is given by

$$\Delta = I(\text{parent}) - \sum_{j=1}^2 \frac{N(v_j)}{N} I(v_j),$$

where

$$I(v) = 1 - \max_c p(c|v).$$

Inserting for the split defined by the projection onto the first PCA being greater or equal to zero we obtain

$$\begin{aligned} \Delta &= (1 - (\frac{70}{210})) \\ &\quad - [\frac{94}{210} (1 - (\frac{70}{94})) \\ &\quad + \frac{116}{210} (1 - (\frac{70}{116}))] \\ &= \frac{2}{3} - \frac{1}{3} = 1/3. \end{aligned}$$

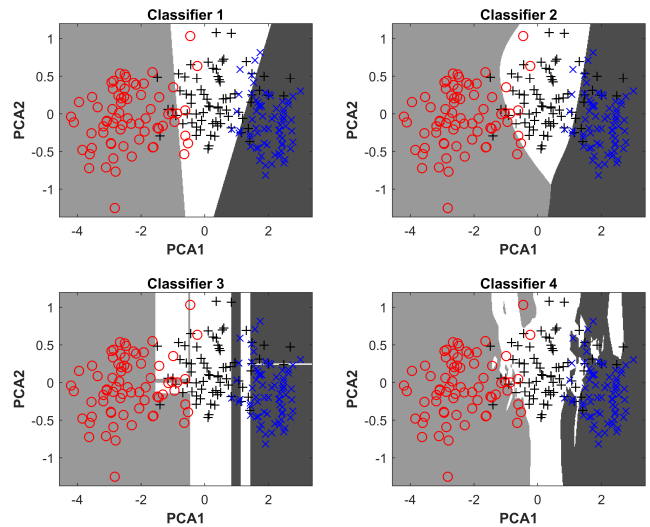


Figure 3: Decision boundaries for four different classifiers trained on the Seeds data projected onto the first two principal components.

Question 5. Four different classifiers are trained on the data projected onto the first two principal components (i.e., using the first and second principal components as features) and the decision boundary for each of the four classifiers is given in Figure 3. Which one of the following statements is *correct*?

- A. Classifier 1 is a decision tree, Classifier 2 is an artificial neural network with three hidden units, Classifier 3 is a multinomial regression model, and Classifier 4 is a 3-nearest neighbor classifier.
- B. Classifier 1 is an artificial neural network with three hidden units, Classifier 2 is a multinomial regression model, Classifier 3 is a 3-nearest neighbor classifier, and Classifier 4 is a decision tree.
- C. Classifier 1 is an artificial neural network with three hidden units, Classifier 2 is a multinomial regression model, Classifier 3 is a decision tree, and Classifier 4 is a 3-nearest neighbor classifier.
- D. Classifier 1 is a multinomial regression model, Classifier 2 is an artificial neural network with three hidden units, Classifier 3 is a decision tree, and Classifier 4 is a 3-nearest neighbor classifier.**
- E. Don't know.

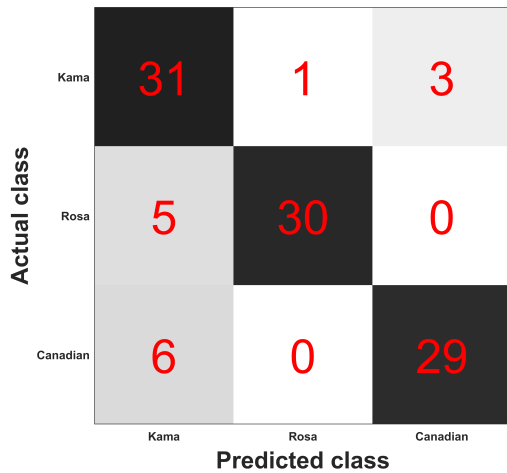


Figure 4: The confusion matrix of a 3-nearest neighbor classifier used to predict the Seeds data.

Solution 5. The decision boundary of classifier 1 is based on lines thus using multinomial regression. Classifier two has smooth boundaries that are non-linear thus based on ANN. Classifier three has axis aligned boundaries corresponding to the decision tree, leaving classifier four as the 3-nearest neighbour due to its very complex and non-smooth boundaries.

Question 6. The data is split in half and a KNN classifier used to predict the test-set based on the training set for $K=3$. The confusion matrix of the KNN classifier is given in Figure 4. What is the accuracy of the classifier?

- A. 0.1429
- B. 0.8571**
- C. 0.8911
- D. 0.9574
- E. Don't know.

Solution 6. The accuracy is given by the number of correctly classified observations out of the total classified observations which is $accuracy = \frac{31+30+29}{31+30+29+1+3+5+6} = 90/105 = 0.8571$.

	O1	O2	O3	O4	O5	O6	O7	O8	O9
O1	0	0.534	1.257	1.671	1.090	1.315	1.484	1.253	1.418
O2	0.534	0	0.727	2.119	1.526	1.689	1.214	0.997	1.056
O3	1.257	0.727	0	2.809	2.220	2.342	1.088	0.965	0.807
O4	1.671	2.119	2.809	0	0.601	0.540	3.135	2.908	3.087
O5	1.090	1.526	2.220	0.601	0	0.331	2.563	2.338	2.500
O6	1.315	1.689	2.342	0.540	0.331	0	2.797	2.567	2.708
O7	1.484	1.214	1.088	3.135	2.563	2.797	0	0.275	0.298
O8	1.253	0.997	0.965	2.908	2.338	2.567	0.275	0	0.343
O9	1.418	1.056	0.807	3.087	2.500	2.708	0.298	0.343	0

Table 2: Pairwise Euclidean distance between nine observations in the Seeds data. Black observations (i.e., O1, O2, O3) are observations corresponding to Kama seeds, red observations (i.e., O4, O5, O6) are observations corresponding to Rosa seeds, and blue observations (i.e., O7, O8, O9) are observations corresponding to Canadian seeds.

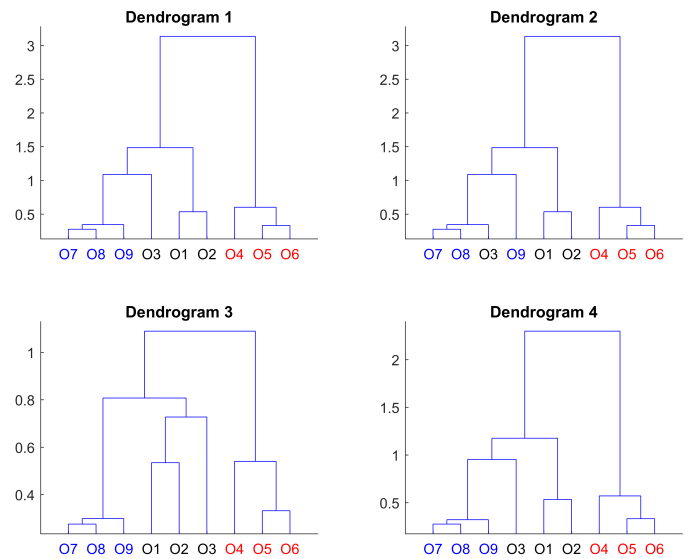


Figure 5: Four different dendrograms derived from the distances between the nine observation in Table 2.

Question 7. In Table 2 is given the pairwise Euclidean distances between nine observations of the Seeds data. A hierarchical clustering is used to cluster these nine observations using complete (i.e., maximum) linkage. Which one of the dendrograms given in Figure 5 corresponds to the clustering?

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

Solution 7. In complete distance clusters are merged according to maximal distance between observations within each cluster. The dendrogram grows by first merging O7 and O8 at 0.275, then O5, O6 at level 0.331, then {O7,O8} with O9 at 0.343, then O1 and O2 at level 0.534, then O4 with {O5,O6} at 0.601, then O3, {O7,O8,O9} at level 1.0882, then {O1, O2} with {O3, O7,O8,O9} at 1.484, and finally {O1, O2, O3 O7,O8,O9} with {O4, O5,O6} at 3.135. Only Dendrogram 1 has these properties.

Question 8. We will consider thresholding Dendrogram 4 at the level of three clusters. We recall that the Rand index also denoted the simple matching coefficient (SMC) between the true labels and the extracted clusters is given by:

$$SMC = \frac{f_{00} + f_{11}}{K},$$

where f_{00} is the number of object pairs in different class assigned to different clusters and f_{11} is the number of object pairs in same class assigned to same cluster, whereas $K = N(N - 1)/2$ is the total number of object pairs where N is the number of observations considered. What is the above SMC between the true labeling of the observations into the three classes Kama, Rosa, and Canadian, and the clustering defined by thresholding Dendrogram 4 at the level of three clusters?

A. 0.7500

B. 0.7778

C. 0.8611

D. 1.0000

E. Don't know.

Solution 8. When thresholding the clustering we obtain: the cluster indices: $[1 \ 1 \ 3 \ 2 \ 2 \ 2 \ 3 \ 3 \ 3]^\top$, whereas the true class labels are $[1 \ 1 \ 1 \ 2 \ 2 \ 2 \ 3 \ 3 \ 3]^\top$. From this, we obtain: Total number of object pairs is: $K = 9(9 - 1)/2 = 36$
 $f_{00} = 2 \cdot (3 + 3) + 3 \cdot (3 + 1) = 24$
 $f_{11} = 2 \cdot (2-1)/2 + 3 \cdot (3-1)/1 + 3 \cdot (3-1)/2 + 1 \cdot (1-1)/2 = 7$
 $SMC = \frac{f_{00} + f_{11}}{K} = \frac{24 + 7}{36} = 0.8611$

Question 9. To determine the type of seed of an observation we will use a k-nearest neighbor (KNN) classifier to predict each of the nine observations based on the Euclidean distance between the observations given in Table 2. We will use leave-one-out cross-validation for the KNN in order to classify the nine considered observations using a two-nearest neighbor classifier, i.e. $K = 2$. For tied classes we will classify the observation according to its closest observation. The analysis will be based only on the data given in Table 2. Which one of the following statements is *correct*?

- A. All the observations will be correctly classified.**
- B. One of the observations will be misclassified.
- C. Two of the observations will be misclassified.
- D. Three of the observations will be misclassified.
- E. Don't know.

Solution 9. $N(O1, 2) = \{O2, O5\}$ as O2 is closest it will be correctly classified as Kama.

$N(O2, 2) = \{O1, O3\}$ and will be correctly classified as Kama.

$N(O3, 2) = \{O2, O9\}$ as O2 is closest it will be correctly classified as Kama.

$N(O4, 2) = \{O6, O5\}$ and will be correctly classified as Rosa.

$N(O5, 2) = \{O6, O4\}$ and will be correctly classified as Rosa.

$N(O6, 2) = \{O5, O4\}$ and will be correctly classified as Rosa.

$N(O7, 2) = \{O8, O9\}$ and will be correctly classified as Canadian.

$N(O8, 2) = \{O7, O9\}$ and will be correctly classified as Canadian.

$N(O9, 2) = \{O7, O8\}$ and will be correctly classified as Canadian.

Question 10. We suspect that observation O4 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 2 only. We recall that the KNN density and average relative density 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 2, what is the average relative density for observation O4 for $K = 1$ nearest neighbors?

- A. 0.54
- B. 0.61**
- C. 1.63
- D. 1.85
- E. Don't know.

Solution 10.

$$\text{density}(\mathbf{x}_{O4}, 1) = \left(\frac{1}{1} \cdot 0.540\right)^{-1} = 1.8519$$

$$\text{density}(\mathbf{x}_{O6}, 1) = \left(\frac{1}{1} \cdot 0.331\right)^{-1} = 3.0211$$

$$\begin{aligned} \text{a.r.d.}(\mathbf{x}_{O4}, 1) &= \frac{\text{density}(\mathbf{x}_{O4}, 1)}{\frac{1}{1}(\text{density}(\mathbf{x}_{O6}, 1))} \\ &= \frac{1.8519}{\frac{1}{1} \cdot 3.0211} = 0.61 \end{aligned}$$

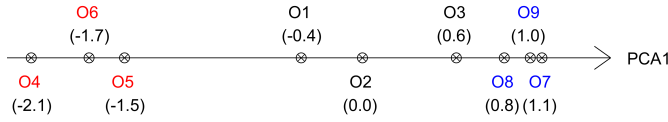


Figure 6: The nine observations considered in Table 2 projected onto the first principal component (the location of the projection is given in parenthesis).

Question 11. We will consider the nine observations projected onto the first principal component given in Figure 6. We will cluster this data using k-means with Euclidean distance into three clusters (i.e., $k=3$) and initialize the k-means algorithm with centroids located at observation O4, O6, and O5. Which one of the following statements is *correct*?

- A. The converged solution will be {O4}, {O6}, {O1, O2, O3, O5, O6, O7, O8, O9}.
- B. The converged solution will be {O4, O5, O6}, {O1, O2, O3}, {O7, O8, O9}.
- C. The converged solution will be {O4, O5, O6}, {O1, O2}, {O3, O7, O8, O9}.
- D. The converged solution will be {O4}, {O5, O6}, {O1, O2, O3, O7, O8, O9}.**
- E. Don't know.

Solution 11. With the described initialization, observation O4 will be assigned to the cluster located at O4, observation O6 will be assigned to the cluster located at O6, and the remaining observations {O1, O2, O3, O5, O7, O8, O9} assigned to the cluster located at O5. Thus, only cluster located at O5 will change location and the location updated to $\frac{-1.5 + -0.4 + 0.0 + 0.6 + 0.8 + 1.0 + 1.1}{7} = 0.2286$. For this new location O5 is closer to cluster located at O6 than the cluster located at 0.2286, resulting in the updated clustering {O4}, {O5, O6}, {O1, O2, O3, O7, O8, O9}. Thus the second cluster will change location to $\frac{-1.7 + -1.5}{2} = -1.6$ whereas the third cluster will change location to $\frac{-0.4 + 0.0 + 0.6 + 0.8 + 1.0 + 1.1}{6} = 0.5167$. As O1 is still closest to cluster 3 there is no change of assignment and the k-means procedure has converged.

Feature(s)	Training	Test
	Error Rate	Error Rate
No features	0.6667	0.6667
x_1	0.1143	0.1524
x_2	0.1143	0.1143
x_3	0.2190	0.1714
x_4	0.1524	0.1714
x_1 and x_2	0.0952	0.1619
x_1 and x_3	0.1143	0.1619
x_1 and x_4	0.1143	0.1619
x_2 and x_3	0.1238	0.1333
x_2 and x_4	0.1048	0.1429
x_3 and x_4	0.1143	0.1619
x_1 and x_2 and x_3	0.0571	0.1714
x_1 and x_2 and x_4	0.1048	0.1619
x_1 and x_3 and x_4	0.0857	0.1619
x_2 and x_3 and x_4	0.0762	0.1524
x_1 and x_2 and x_3 and x_4	0.0667	0.1810

Table 3: Error rate for the training and test set when using multinomial regression to predict the type of seed using different combinations of the four attributes ($x_1 - x_4$) based on the hold-out method with 50 % of the observations hold-out for testing.

Question 12. A multinomial regression classifier is trained using different combinations of the four attributes x_1, x_2, x_3 , and x_4 . Table 3 gives the training and test performance of the multinomial regression classifier when trained using different combinations of the four attributes. Which one of the following statements is *correct*?

- A. Forward selection will result in a better model being selected than backward selection.
- B. Neither forward nor backward selection will identify the optimal feature combination for this problem.
- C. Backward selection will use a model that includes three features.
- D. Forward selection will select only one feature.**
- E. Don't know.

Solution 12. Using forward and backward selection we would like to minimize the test error rate. Thus, the forward selection would first select x_2 having lowest test error rate. As no combination of the feature

x_2 lead to improvement in the test error rate the forward selection method terminates. Backward selection starts with all features and an improvement is found removing feature x_1 providing a test error rate of 0.1524 for x_2 and x_3 and x_4 . Subsequently, x_4 is removed providing a test error rate of 0.1333 for the features x_2 and x_3 . Finally x_3 is removed as having only feature x_2 provides an error rate of 0.1143 and the method terminates.

Question 13. We would like to investigate if we can predict the width of a seed kernel (x_4) based on the area (x_1), perimeter (x_2), and length of kernel (x_3). For this purpose regularized least squares regression is applied based on minimizing with respect to \mathbf{w} the cost function:

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

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 will consider the following four different values of λ : $\lambda_1 = 1$, $\lambda_2 = 10$, $\lambda_3 = 100$, and $\lambda_4 = 1000$. We obtain the following four different solutions for \mathbf{w} here given in random order of the values of λ considered:

$$\begin{aligned} \mathbf{w}_a &= \begin{bmatrix} 0.0538 \\ 0.0558 \\ 0.1861 \\ -0.0596 \end{bmatrix}, & \mathbf{w}_b &= \begin{bmatrix} 0.0089 \\ 0.0931 \\ 0.1093 \\ 0.0417 \end{bmatrix}, \\ \mathbf{w}_c &= \begin{bmatrix} 0.2811 \\ 0.0445 \\ 0.3379 \\ -0.4626 \end{bmatrix}, & \mathbf{w}_d &= \begin{bmatrix} 0.0167 \\ 0.0698 \\ 0.1354 \\ 0.0403 \end{bmatrix}. \end{aligned}$$

Which one of the following solutions to \mathbf{w} corresponds to the correct value of λ ?

- A. \mathbf{w}_a corresponds to λ_2 .**
- B. \mathbf{w}_b corresponds to λ_2 .
- C. \mathbf{w}_c corresponds to λ_2 .
- D. \mathbf{w}_d corresponds to λ_2 .
- E. Don't know.

Solution 13. As we increase λ we put more and more emphasis on the regularization penalization term $\|\mathbf{w}\|_2^2$. Thus, by evaluating this term for each solution vector we obtain: $\|\mathbf{w}_a\|_2^2 = 0.0442$, $\|\mathbf{w}_b\|_2^2 = 0.0224$, $\|\mathbf{w}_c\|_2^2 = 0.4092$, $\|\mathbf{w}_d\|_2^2 = 0.0251$. Thus sorting by these norm values we find that \mathbf{w}_c corresponds to λ_1 , \mathbf{w}_a corresponds to λ_2 , \mathbf{w}_d corresponds to λ_3 , \mathbf{w}_b corresponds to λ_4 .

No.	Attribute description
x_1	Occurrence of nausea
x_2	Lumbar pain
x_3	Urine pushing
x_4	Micturition pains
x_5	Burn/itch/swell urethra outlet
y	Inflammation of urinary bladder

Table 4: The attributes considered from the study on acute inflammation (taken from <https://archive.ics.uci.edu/ml/datasets/Acute+Inflammations>). The attributes x_1 - x_5 and y are binary where we use 1 for true and 0 for false.

Question 14. In a study of acute inflammation we would like to predict urinary bladder inflammation (the data is taken from <https://archive.ics.uci.edu/ml/datasets/Acute+Inflammations>). We will consider a subset of the attributes, these attributes are given in Table 4. From the study we have

- 49.17 pct. of the persons have inflammation of urinary bladder.
- 32.20 pct. of the persons that have inflammation of urinary bladder have occurrence of nausea.
- 16.39 pct. of the persons that do not have inflammation of urinary bladder have occurrence of nausea.

What is the probability that a person that has occurrence of nausea, i.e. $x_1 = 1$, has inflammation of the urinary bladder, i.e. $y = 1$, according to this study?

- A. 15.83 %
- B. 32.20 %
- C. 65.52 %**
- D. 98.82%
- E. Don't know.

Solution 14. According to Bayes' theorem we have:

$$\begin{aligned}
 P(y = 1 | x_1 = 1) &= \frac{P(x_1=1|y=1)P(y=1)}{P(x_1=1)} \\
 &= \frac{P(x_1=1|y=1)P(y=1)}{P(x_1=1|y=1)P(y=1) + P(x_1=1|y=0)P(y=0)} \\
 &= \frac{0.3220 \cdot 0.4917}{0.3220 \cdot 0.4917 + 0.1639 \cdot (1 - 0.4917)} \\
 &= 0.6552
 \end{aligned}$$

	x_1	x_2	x_3	x_4	x_5	y
P1	1	1	1	1	0	1
P2	0	0	0	0	0	0
P3	1	1	0	1	0	0
P4	0	1	1	0	1	0
P5	1	1	1	1	1	1
P6	0	0	0	0	0	0
P7	1	1	0	1	0	0
P8	0	1	1	0	1	0
P9	1	1	1	1	0	1
P10	0	1	1	0	1	0
P11	0	0	0	0	0	0
P12	1	1	0	1	0	0
P13	0	1	1	0	1	0
P14	0	1	1	0	1	0

Table 5: Provided in the above table are the last 14 observations of the acute inflammation data.

Question 15. We will consider a subset of the acute inflammation data given by the last 14 observations provided in Table 5. We will consider this dataset a market basket with 14 persons (P1-P14) denoting the customers and six items denoted $x_1 - x_5$ and y corresponding to the five input attributes and output variable respectively of the features described in Table 4. What are all frequent itemsets with support greater than 40%?

- A. $\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\}, \{x_1, x_2\}, \{x_2, x_3\}, \{x_2, x_4\}, \{x_2, x_5\}, \{x_3, x_5\}$.
- B. $\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\}, \{x_1, x_2\}, \{x_1, x_4\}, \{x_2, x_3\}, \{x_2, x_4\}, \{x_2, x_5\}, \{x_3, x_5\}$.
- C. $\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\}, \{x_1, x_2\}, \{x_1, x_4\}, \{x_2, x_3\}, \{x_2, x_4\}, \{x_2, x_5\}, \{x_3, x_5\}, \{x_2, x_3, x_5\}$.
- D. $\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\}, \{x_1, x_2\}, \{x_1, x_4\}, \{x_2, x_3\}, \{x_2, x_4\}, \{x_2, x_5\}, \{x_3, x_5\}, \{x_1, x_2, x_4\}, \{x_2, x_3, x_5\}$.**
- E. Don't know.

Solution 15. For a set to have support more than 40% the set must occur at least $0.4 \cdot 14 = 5.6$, i.e. 6 out of the 14 times. All the itemsets that have this property are $\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\}, \{x_1, x_2\}, \{x_1, x_4\}, \{x_2, x_3\}, \{x_2, x_4\}, \{x_2, x_5\}, \{x_3, x_5\}, \{x_1, x_2, x_4\}, \{x_2, x_3, x_5\}$.

Question 16. What is the confidence of the association rule $\{x_1, x_2, x_3, x_4, x_5\} \rightarrow \{y\}$?

- A. 0.0%
- B. 7.1 %
- C. 21.4%
- D. 100.0 %**
- E. Don't know.

Solution 16. The confidence is given as

$$\begin{aligned} P(y = 1|x_1 = 1, x_2 = 1, x_3 = 1, x_4 = 1, x_5 = 1) &= \\ \frac{P(y = 1, x_1 = 1, x_2 = 1, x_3 = 1, x_4 = 1, x_5 = 1)}{P(x_1 = 1, x_2 = 1, x_3 = 1, x_4 = 1, x_5 = 1)} &= \\ = \frac{1/14}{1/14} = 1 = 100.0\% \end{aligned}$$

Question 17. We would like to predict whether a subject has inflammation of urinary bladder ($y = 1$) or not ($y = 0$) using the data in Table 5 and the attributes x_1 , and x_2 only. We will apply a Naïve Bayes classifier that assumes independence between the two attributes given y . Given that a person has $x_1 = 1$, and $x_2 = 1$ what is the probability that the person has an inflammation of urinary bladder ($y = 1$) according to the Naïve Bayes classifier?

- A. 1/14
- B. 3/14
- C. 1/2
- D. 11/19**
- E. Don't know.

Solution 17. According to the Naïve Bayes classifier we have

$$\begin{aligned} P(y = 1|x_1 = 1, x_2 = 1) &= \\ \frac{\begin{pmatrix} P(x_1 = 1|y = 1) \times \\ P(x_2 = 1|y = 1) \times \\ P(y = 1) \end{pmatrix}}{\begin{pmatrix} P(x_1 = 1|y = 1) \times \\ P(x_2 = 1|y = 1) \times \\ P(y = 1) \end{pmatrix} + \begin{pmatrix} P(x_1 = 1|y = 0) \times \\ P(x_2 = 1|y = 0) \times \\ P(y = 0) \end{pmatrix}} &= \\ = \frac{3/3 \cdot 3/3 \cdot 3/14}{3/3 \cdot 3/3 \cdot 3/14 + 3/11 \cdot 8/11 \cdot 11/14} = . \end{aligned}$$

Question 18. Considering the data in Table 5, we will use x_1 to classify whether a subject has inflammation of urinary bladder ($y = 1$) or not ($y = 0$). We will quantify how useful x_1 is for this purpose by calculating the area under curve (AUC) of the receiver operator characteristic (ROC). Which one of the ROC curves given in Figure 7 corresponds to using the feature x_1 to determine if a subject has inflammation of urinary bladder?

- A. The curve with AUC=0.636.
- B. The curve with AUC=0.864.**
- C. The curve with AUC=0.909.
- D. The curve with AUC=1.000.
- E. Don't know.

Solution 18. The ROC is defined by considering all conceivable thresholds and plotting the true positive

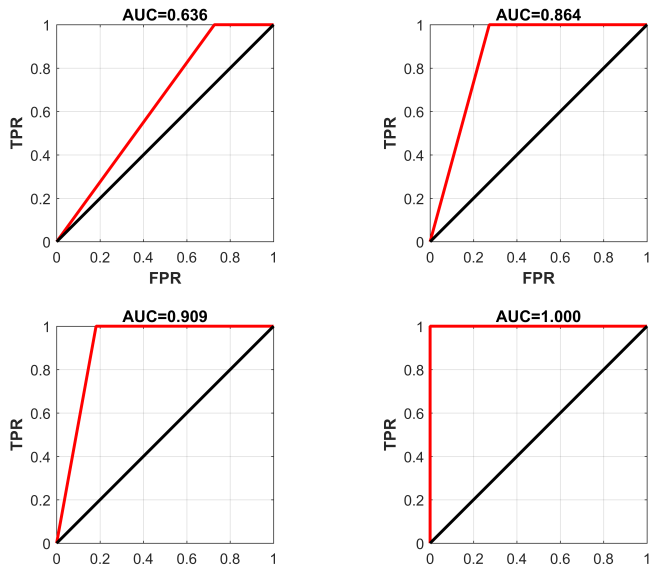


Figure 7: Four different receiver operating characteristic (ROC) curves and their corresponding area under the curve (AUC).

rate (TPR) against the false positive rate (FPR). For a threshold larger than 1 we have that $TPR=FPR=0$. For the threshold at 1 we have that 3 of the 11 observations having $y=0$ have $x_1 = 1$ thus $FPR=3/11$, whereas all three of the three observations with $y=1$ have $x_1 = 1$, thus $TPR=1$. lowering the threshold we get when thresholding above 0 that all 11 of the 11 observations or which $y=0$ are false positive, i.e. $FPR=1$, and all three of three observations where $y=1$ are true positive, thus $TPR=1$, giving an $AUC = 0.864$.

Question 19. Considering the data in Table 5, we will calculate the similarity between $P1$ given as the vector $\mathbf{r} = [1 \ 1 \ 1 \ 1 \ 0 \ 1]$ and $P3$ given by the vector $\mathbf{s} = [1 \ 1 \ 0 \ 1 \ 0 \ 0]$ using Jaccard, Simple Matching Coefficient, and Cosine similarity given respectively by:

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

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

$$\cos(\mathbf{r}, \mathbf{s}) = \frac{f_{11}}{\|\mathbf{r}\|_2 \|\mathbf{s}\|_2}.$$

Which one of the following statements regarding the similarity of \mathbf{r} and \mathbf{s} is correct?

- A. $J(\mathbf{r}, \mathbf{s}) < SMC(\mathbf{r}, \mathbf{s})$
- B. $J(\mathbf{r}, \mathbf{s}) > \cos(\mathbf{r}, \mathbf{s})$
- C. $SMC(\mathbf{r}, \mathbf{s}) > \cos(\mathbf{r}, \mathbf{s})$
- D. $\cos(\mathbf{r}, \mathbf{s}) = 3/15$
- E. Don't know.

Solution 19. For \mathbf{r} and \mathbf{s} we have:

$$J(\mathbf{r}, \mathbf{s}) = \frac{f_{11}}{M - f_{00}} = 3/5 = 0.6000,$$

$$SMC(\mathbf{r}, \mathbf{s}) = \frac{f_{11} + f_{00}}{M} = 4/6 = 0.6667,$$

$$\cos(\mathbf{r}, \mathbf{s}) = \frac{f_{11}}{\|\mathbf{r}\|_2 \|\mathbf{s}\|_2} = 3/(\sqrt{5}\sqrt{3}) = 0.7746.$$

Thus, $J(P1, P3) < SMC(P1, P3)$ is the only correct statement.

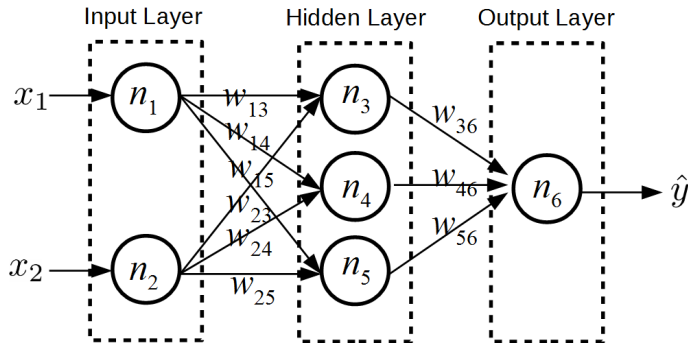


Figure 8: The architecture of the considered neural network having one hidden layer.

Question 20. A neural network is trained to separate persons with urinary inflammation ($y = 1$) from persons not having urinary inflammation based on the features x_1 and x_2 . The structure of the neural network is outlined in Figure 8. The activation function used for all six neurons n_1, n_2, n_3, n_4, n_5 , and n_6 is the rectified linear unit

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise.} \end{cases}$$

The neural network has no biases, i.e. all the biases of all units are zero. The weights of the network are:

$$\begin{aligned} w_{13} &= 0.5, & w_{14} &= 0.5, & w_{15} &= -0.5, \\ w_{23} &= 0.5, & w_{24} &= -0.5, & w_{25} &= 0.25, \\ w_{36} &= 0.25, & w_{46} &= -0.25, & w_{56} &= 0.25. \end{aligned}$$

What will be the output (\hat{y}) of the neural network for an observation having $x_1 = 1$ and $x_2 = 1$?

- A. 0
- B. 0.25**
- C. 0.75
- D. 1
- E. Don't know.

Solution 20. The output of the neurons in the hidden layer will be:

$$\begin{aligned} n_3 : f(0.5 \cdot 1 + 0.5 \cdot 1) &= f(1) = 1, \\ n_4 : f(0.5 \cdot 1 - 0.5 \cdot 1) &= f(0) = 0, \\ n_5 : f(-0.5 \cdot 1 + 0.25 \cdot 1) &= f(-0.25) = 0. \end{aligned}$$

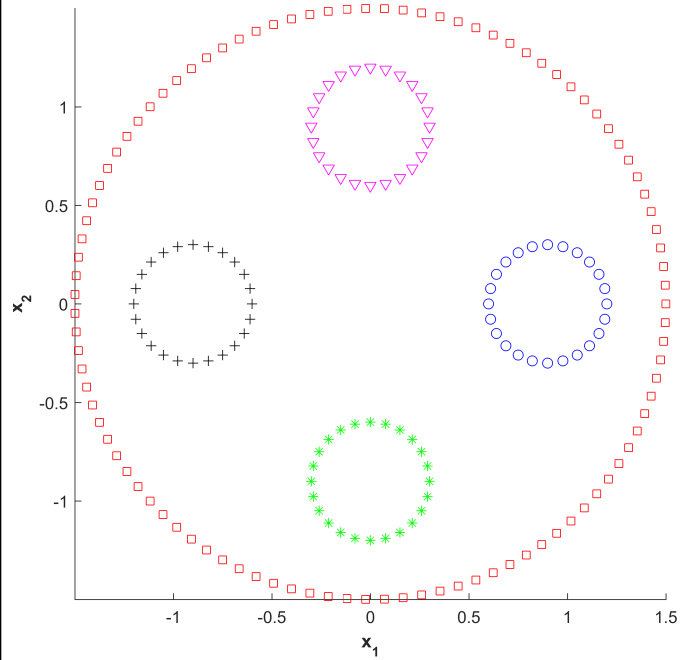


Figure 9: A dataset with five classes given respectively by a large circle and four smaller circles.

The output of the output neurons will therefore be:

$$n_6 : f(0.25 \cdot 1 - 0.25 \cdot 0 + 0.25 \cdot 0) = f(0.25) = 0.25.$$

Question 21. We will consider the dataset with five classes given in Figure 9 defined respectively by the four inner circles and the larger outer circle. We will cluster this dataset using hierarchical clustering. What would be a suitable measure of proximity and linkage in order to perfectly separate the five classes into five clusters?

- A. Average linkage using the 2-norm (i.e. $\|x - y\|_2$) as proximity measure.
- B. Single linkage using the 1-norm (i.e. $\|x - y\|_1$) as proximity measure.**
- C. Complete linkage using the 2-norm (i.e. $\|x - y\|_2$) as proximity measure.
- D. Complete linkage using the 1-norm (i.e. $\|x - y\|_1$) as proximity measure.
- E. Don't know.

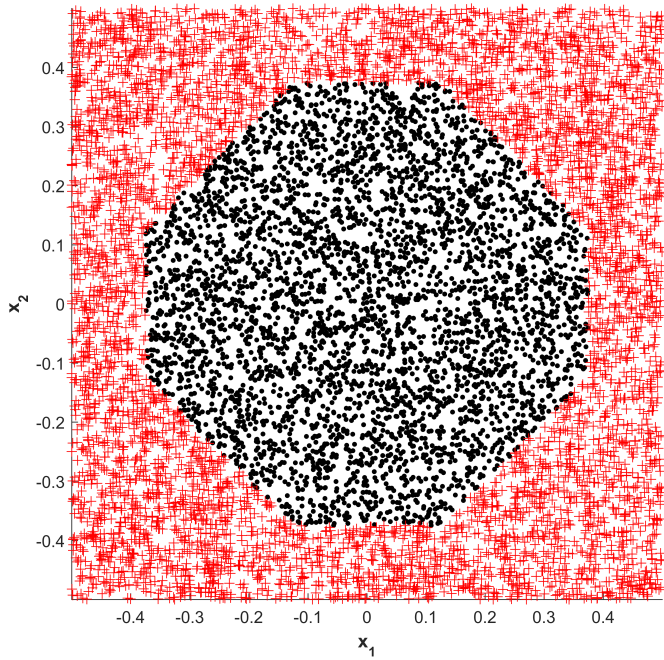


Figure 10: A two class classification problem.

Solution 21. The choice of norm will not generally influence the results much in this example as observations that are close in 1-norm will also be close in 2-norm. However, linkage function will heavily influence the results. As all clusters have the property that at least one observation is closer within the cluster than an observation in another cluster, thus, a contiguity based approach will be well-suited. Hence, single linkage clustering will perfectly separate the classes whereas the other approaches will fail.

Question 22. Consider the dataset with two classes given in Figure 10. Which one of the following decisions would lead to a perfect separation of the two classes?

- A. If $\|x\|_1 \leq \frac{1}{4}$ and $\|x\|_2 \leq \frac{3}{8}$ then black dot, otherwise red plus.
- B. If $\|x\|_2 \leq \frac{3}{8}$ and $\|x\|_\infty \leq \frac{1}{4}$ then black dot, otherwise red plus.
- C. If $\|x\|_1 \leq \frac{1}{2}$ and $\|x\|_\infty \leq \frac{1}{2}$ then black dot, otherwise red plus.
- D. If $\|x\|_1 \leq \frac{1}{2}$ and $\|x\|_\infty \leq \frac{3}{8}$ then black dot, otherwise red plus.**
- E. Don't know.

Solution 22. The decision boundary is formed by a hexagonal shape. Inspecting the position of the edges it is seen that the decision boundary traverses the coordinates $(0, 3/8)$ and $(0.25, 0.25)$ corresponding to a point where $\|x\|_1 = \frac{1}{2}$ and $\|x\|_\infty = \frac{3}{8}$ respectively.

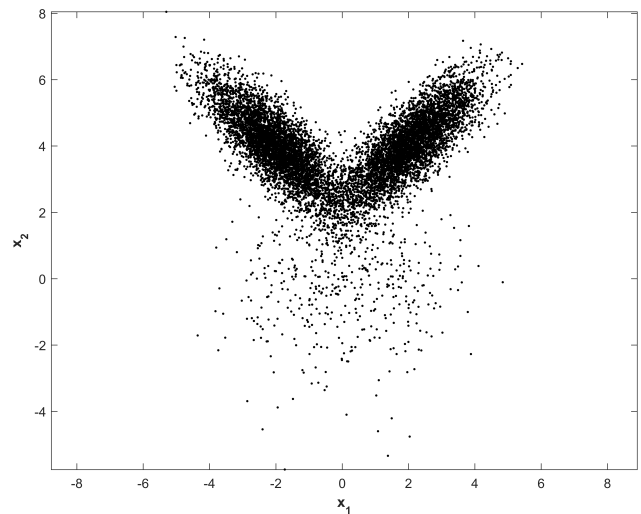


Figure 11: 10,000 data observations drawn from a Gaussian Mixture Model (GMM).

Question 23. Consider the 10.000 observations drawn from a Guassian Mixture Model (GMM) shown in Figure 11. 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 best characterize the data?

A.

$$p(\mathbf{x}) = 0.5 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}\right) + 0.5 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}\right)$$

B.

$$p(\mathbf{x}) = 0.05 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}\right) + 0.475 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}\right) + 0.475 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}\right)$$

C.

$$p(\mathbf{x}) = 0.5 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}\right) + 0.25 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}\right) + 0.25 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}\right)$$

D.

$$p(\mathbf{x}) = 0.1 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}\right) + 0.45 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}\right) + 0.45 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}\right)$$

E. Don't know.

Solution 23. There are three clusters and the centroids of these clusters are $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 2 \\ 4 \end{bmatrix}$, $\begin{bmatrix} -2 \\ 4 \end{bmatrix}$. The

cluster at $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$ is not very dense and should therefore have a very low weight, furthermore, the clusters at $\begin{bmatrix} 2 \\ 4 \end{bmatrix}$ has positive covariance whereas the cluster at

$\begin{bmatrix} -2 \\ 4 \end{bmatrix}$ has negative covariance. This property only holds for the answer option:

$$p(\mathbf{x}) = 0.05 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}\right) + 0.475 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}\right) + 0.475 \cdot \mathcal{N}\left(\mathbf{x} \middle| \begin{bmatrix} -2 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}\right)$$

Question 24. We will consider a very large dataset with 100 mio. observations and ten features, i.e. $N = 100.000.000$ and $M = 10$. We would like to perform two-level cross-validation in order to select between 3 different settings of the parameters of a model (inner fold) and estimate the generalization error (outer fold). We are only allowed to train maximally 65 models in total. Which one of the following procedures satisfies this constraint?

- A. Five fold cross-validation in both the outer and inner folds.
- B. Leave-one-out cross-validation for the outer fold and hold-out 50 % for the inner fold.
- C. Ten-fold cross-validation for the outer fold and two fold cross-validation for the inner fold.
- D. Two-fold cross-validation for the outer fold and ten fold cross-validation for the inner fold.**
- E. Don't know.

Solution 24. In the inner fold we have to train as many models as we have folds to identify the optimal parameters. We then use the optimal parameters to train a model on the entire training set in order to evaluate this model on the test set defined in the outer fold. Thus we need for each outer fold to train the number of inner folds + one model (i.e. the model trained on the entire training data) times number of outer folds, i.e. $K_1(K_2 \cdot S + 1)$. This gives for:

Five fold cross-validation in both the outer and inner folds: $5(5 \cdot 3 + 1) = 80$ models

Leave-one-out cross-validation for the outer fold and hold-out: $100.000.000(1 \cdot 3 + 1) = 400.000.000$ models

Ten-fold cross-validation for the outer fold and two fold cross-validation for the inner fold: $10(2 \cdot 3 + 1) = 70$

Two-fold cross-validation for the outer fold and ten fold cross-validation for the inner fold: $2(10 \cdot 3 + 1) = 62$.

Question 25. We recall that the AdaBoost algorithm is given by updating the weight to the i 'th data observation (w_i) based on the classifier f_t at round t according to:

$$w_i(t+1) = \frac{\tilde{w}_i(t+1)}{\sum_{j=1}^N \tilde{w}_j(t+1)}, \text{ where}$$

$$\tilde{w}_i(t+1) = \begin{cases} w_i(t)e^{-\alpha_t} & \text{if } f_t(\mathbf{x}_i) = y_i \\ w_i(t)e^{\alpha_t} & \text{if } f_t(\mathbf{x}_i) \neq y_i. \end{cases}$$

Here $\alpha_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$ (where \log is the natural logarithm) and $\epsilon_t = \sum_{i=1}^N w_i (1 - \delta_{f_t(\mathbf{x}_i), y_i})$, where $\delta_{f_t(\mathbf{x}_i), y_i} = 1$ if $f_t(\mathbf{x}_i) = y_i$ and zero otherwise. Initially the weights are uniform across samples, i.e. $w_1 = w_2 = \dots = w_N = 1/N$ where N is the number of observations.

A dataset is sampled with replacement from this uniform distribution and the classifier is trained on this sampled data. Using this trained classifier 5 of the original 25 observations are misclassified. What will the updated weights be for these misclassified observations according to the AdaBoost algorithm?

- A. 0.02
- B. 0.025
- C. 0.08
- D. 0.1**
- E. Don't know.

Solution 25. As we have 5 misclassified observations the weighted error rate will be $\epsilon_1 = 5/25 = 1/5$, thus $\alpha_t = \frac{1}{2} \log(\frac{1-1/5}{1/5}) = \frac{1}{2} \log 4 = 0.6931$. Thus for correctly classified observations we have: $\tilde{w}_j(t+1) = w_j(t)e^{-\alpha_t} = 1/25e^{-0.6931}$ and for incorrectly classified observations $\tilde{w}_j(t+1) = w_j(t)e^{\alpha_t} = 1/25e^{+0.6931}$. We thereby obtain for the updated weights for misclassified observations: $w_i(t+1) = \frac{1/25e^{+0.6931}}{20/25e^{-0.6931} + 5/25e^{+0.6931}} = 0.1$

Question 26. For which of the following purposes is cross-validation *the least* well suited?

- A. Select the number of hidden units in artificial neural networks (ANN).
- B. Select the width of the Gaussian kernel in kernel density estimation (KDE).
- C. Select the observations that minimize the training error.**
- D. Select the number of neighbors in KNN classification.
- E. Don't know.

Solution 26. Cross-validation can trivially be used to quantify the number of hidden units in ANN and nearest neighbors in KNN classification by evaluating the performance predicting the output in these supervised

learning problems. As we have seen in the course cross-validation can also be used to quantify the width of the kernel density estimator. Cross-validation is used to quantify models generalization through the use of the test sets and not to minimize the training error.

Question 27. Which of the following statements regarding ensemble methods is correct?

- A. In ensemble methods it is important that the different trained classifiers perform very similar.
- B. In Random Forest features are randomly sampled at each node of the tree.**
- C. Random Forest is the same as fitting several decision trees and classifying according to the tree for which the leaf has highest purity.
- D. In bagging the output class labels are randomly changed to introduce noise for robustness.
- E. Don't know.

Solution 27. Using ensemble methods it is important that the different methods are not the same but as independent as possible. In Random Forest $m < M$ features are indeed randomly sampled at each node of the tree. Majority voting is used for the classification and not the tree with highest purity of the leaf. In bagging observations are uniformly sampled with replacement and there is no additional emphasis to misclassified observations.