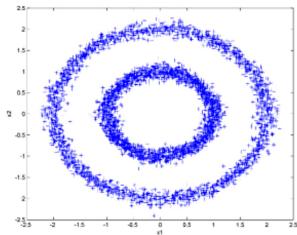


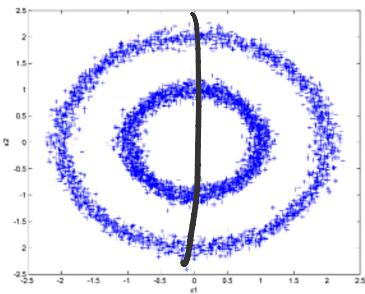
1. Perform clustering on the two-dimensional data illustrated below.
- Explain how the following four algorithms would split the data into two clusters: **K-means**, **hierarchical clustering** (separately **single linkage** and **complete linkage**), **DBSCAN**. Draw the clusters and give brief explanations to your answers!
 - Let us assume that the goal is to find the two annular natural clusters. Give an $R^2 \rightarrow R^2$ **coordinate transformation** that assists the bad-performing clustering algorithms to find the natural clusters. Plot the coordinate system transformation!

(20%)



16%

a.

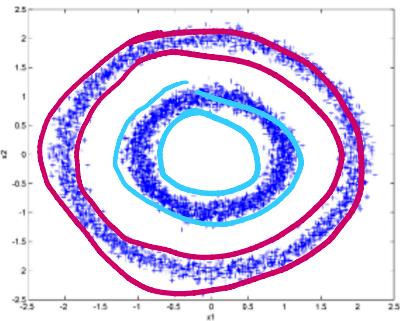


↳ K-mean

k-means would probably split the data in half vertically or horizontally. k-means will always form spherical clusters around its centroids. hence if it will not be able to find the natural clusters here.



↳ DBSCAN & singly linkage + complete linkage



DBSCAN would perform well here as it can handle clusters of different shapes and sizes like the ones we have here. Especially as the natural clusters are well separated.

Singly linkage will also perform well as the points of each natural cluster are very close to each other so when the algorithm needs to decide which clusters to merge with, it will be able to find the natural clusters. Complete linkage will do the same by the same process.

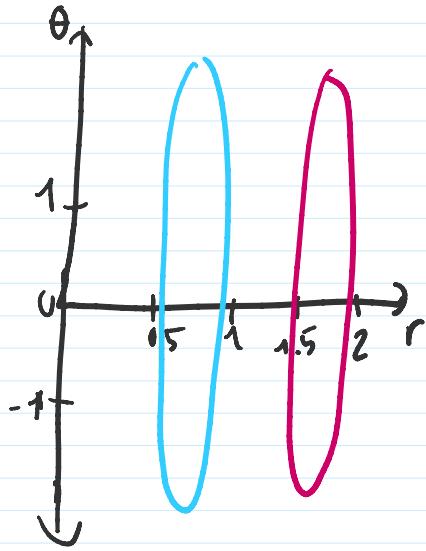
Are you sure about complete linkage?
There exist points from different natural clusters that are closer to each other than some points within the same natural clusters.

b. We could use something like a polar coordinates



Then, we are going to have the well separated natural clusters.





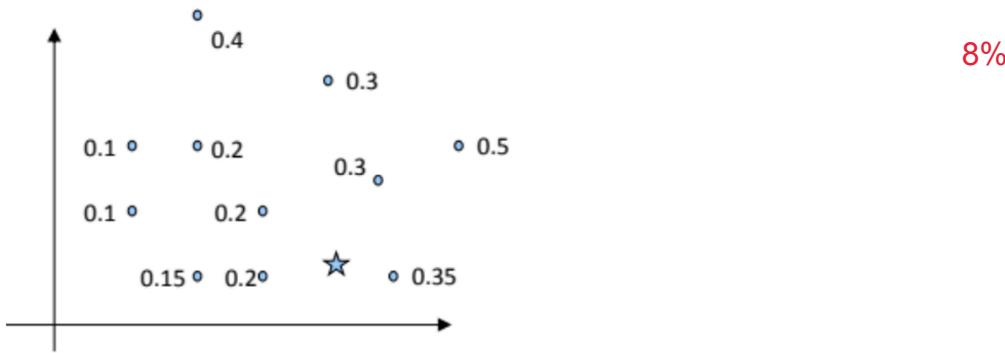
Then, we are going to have two well separated natural clusters.



2. kNN regression

- Determine the target value of the record marked by a star using **kNN regression** with the choice of $k = 4$ (without using distance-weights)!
- Name at least **three further algorithms** that can solve **regression** problems!

(10%)



a. $\frac{0.2 + 0.2 + 0.3 + 0.35}{4} = 0.2625$



b. Decision tree, linear regression, logistic regression.

Logistic regression is not a regression algorithm despite its name! It is a binary classification algorithm as you also note later!

3. Are the following statements true or false? Explain your answer!

(15%)

- a. Using AdaBoost algorithm, if the j th classifier correctly classifies the i th record then the weight of the i th record will be certainly reduced in the $(j+1)$ th step.

True. The correctly classified record's weight has to decrease so the incorrectly classified ones' can increase.

- b. Perceptron algorithm is a universal function approximator.

False. A neural network with one hidden layer can do it according to the universal approximation theorem, not a perceptron.

- c. One-vs-one strategy is computationally more expensive than the one-vs-rest strategy.

True. For a problem with N classes, one-vs-one has to make $\binom{N}{2}$ binary classifications. One-vs-rest makes N classifiers.

- d. SMOTE is an efficient undersampling method.

False. It is an oversampling method.

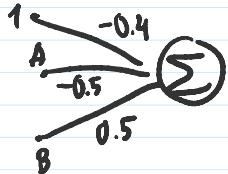
- e. Logistic regression is a binary classification algorithm.

True. It will predict the probability that a new record belongs to the labeled class. We can combine multiple logistic regression models to perform multiple class classification using OvO or OvR strategy.

4. Represent the "(NOT A) AND B" Boolean function with a perceptron or show that it is not possible to do so. In the latter case, construct a neural network with one hidden layer. Use the usual activation function!

10%

(10%)



$$y = -0.5A + 0.5B - 0.4 \text{ with the usual activation function}$$



5. Hierarchical clustering

- Using the following distance matrix draw the dendograms corresponding to the single linkage (MIN) and complete linkage (MAX) clustering algorithms.
- What advantages does hierarchical clustering have in general compared to K-means algorithm?
- What advantages and disadvantages do single and complete linkage techniques have compared to one another?

20%

(20%)

Item	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

a. Single linkage

Complete linkage

12%

Not necessarily! If the classifier has negative importance then the weight of the misclassified record will increase in the next step.

a. Single linkage

Item	A ∩ B	C	D	E
A ∩ B	0	2	2	3
C	2	0	1	5
D	2	1	3	7
E	3	5	0	0

Complete linkage

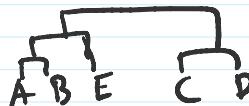
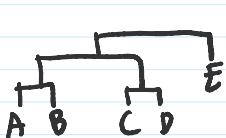
Item	A ∩ B	C	D	E
A ∩ B	0	2	4	3
C	2	0	1	5
D	4	1	0	3
E	3	5	0	0

Item	A ∩ B	C ∩ D	E
A ∩ B	0	2	3
C ∩ D	2	0	3
E	3	0	0

Item	A ∩ B	C ∩ D	E
A ∩ B	0	4	3
C ∩ D	4	0	5
E	3	5	0

Item	A ∩ B ∩ C ∩ D	
A ∩ B ∩ C ∩ D	0	3
E	3	0

Item	A ∩ B ∩ E	
A ∩ B ∩ E	0	5
C ∩ D	5	0



- b. We don't have to decide the number of clusters before running the algorithm.

The dendrogram it produces carries hierarchical information.
Can possibly deal better with non-spherical clusters.

Deal better with clusters of different sizes.



- c. Single linkage handles elliptical shapes better, but it is more sensitive for noise and outliers.

Complete linkage is less sensitive for noise and outliers, but it tends to divide large spherical clusters.



6. We aim to minimize the squared error with regularization term regarding a linear regression model:

$$((w_0 + w^T x) - y)^2 + \lambda \|w\|_2^2$$

where $w = (w_1, w_2, \dots, w_p)$: vector of weights (parameters, coefficients), x : feature vector, λ : regularization parameter, $\|w\|_2$ denotes the Euclidean length of the vector w . Note that the constant term (w_0) is not regularized.

18%

- What is the role of regularization?
- Give the formula for one update step of stochastic gradient descent method for w_0 and for w_i ($i \neq 0$), with η learning rate for a training record (x_1, x_2, \dots, x_p) with target variable y .
- What does the learning rate refer to? What are the advantages and disadvantages of using a small/large learning rate?
- What is the difference between gradient descent method and stochastic gradient descent method?

(25%)

a. To avoid overfitting.

b. Let $\vec{x} = (x_0, \dots, x_p)$. $\text{Err}^2 = \|w_0 + \vec{w}^T \vec{x}\|^2 - y^2 + \lambda \|w\|^2$



$$\vec{w}_0 = w_0 - \eta \left(\begin{array}{c} \frac{\partial \text{Err}^2}{\partial w_0} \\ \vdots \\ \frac{\partial \text{Err}^2}{\partial w_p} \end{array} \right) \parallel \vec{w}_0$$

w_0 and w_i are scalars not vectors.
What are the derivatives?

$$\vec{w}_i = w_i - \eta \left(\begin{array}{c} \frac{\partial \text{Err}^2}{\partial w_0} \\ \vdots \\ \frac{\partial \text{Err}^2}{\partial w_p} \end{array} \right) \parallel \vec{w}_i$$

c. Learning rate refers to the amount by which the weights will change at every step.



Using a large learning rate can make our algorithm faster, but it might miss / skip over the minimal minima we are trying to reach. On the other hand, using a small learning rate can fix the problem, but it is going to make the algorithm slower.

d. The gradient used in stochastic gradient descent is calculated for one random record, not the entire dataset like the gradient descent.

