

CP 3252- Machine learning Lab

1. Use the k-means algorithm and Euclidean distance to cluster the following 8 examples into 2 clusters: $A_1=(2,4)$, $A_2=(2,9)$, $A_3=(8,3)$, $A_4=(5,7)$, $A_5=(7,5)$, $A_6=(6,4)$, $A_7=(1,2)$, $A_8=(4,9)$. Write the distance matrix based on the Euclidean distance.

Suppose that the initial seeds (centers of each cluster) are A_2 , A_3 and A_7 . Run the k-means algorithm for 2 epochs only. At the end of this epoch show:

- a) The new clusters (i.e. the examples belonging to each cluster)
- b) The centers of the new clusters
- c) Draw a 10 by 10 space with all the 8 points and show the clusters after the second epoch and the new centroids.
- d) Implement and show how many more iterations are needed to converge? Draw the result for each epoch.(only implementation for Qn.d)

2. For the given dataset,

Sepal Length	Sepal Width	Species
5.3	3.7	Setosa
5.1	3.8	Setosa
7.2	3.0	Virginica
5.4	3.4	Setosa
5.1	3.3	Setosa
5.4	3.9	Setosa
7.4	2.8	Virginica
6.1	2.8	Versicolor
7.3	2.9	Virginica
6.0	2.7	Versicolor
5.8	2.8	Virginica
6.3	2.3	Versicolor
5.1	2.5	Versicolor
6.3	2.5	Versicolor
5.5	2.4	Versicolor

1. Implement KNN Classifier on the given dataset to predict the species for the given flower.

Sepal Length	Sepal Width	Species
5.2	3.1	?

Calculate the distance between the query instance and all the training samples using Euclidean distance to find the k nearest neighbour for the given flower.

- Use the simple majority of the category of nearest neighbors as the prediction value of the query instance. (Use k=3, k=4 and k=5).
- Analyse whether the predicted species is same if the distance metric used is Manhattan distance ($|x_1 - x_2| + |y_1 - y_2|$)
- Construct Decision Tree and find accuracy of the classifier for the given dataset.

TASK

The customer Dataset (loan.csv) consists of information about 381 customers and status of their loan application result (Y/N) as a binary (2-class) classification problem.

- Consider the columns: Gender, Married, Education, ApplicantIncome, CoapplicantIncome, LoanAmount, Loan_Amount_Term, Credit_History, Property_Area as feature Input (X) and Loan_Status as label output (Y).
- Perform the necessary conversion (label encoding and feature scaling) of appropriate features.
- Classify the data set (20% test, 80% training) by using both Decision Tree.
- For the Decision Tree Algorithm generate the Tree and show/save it as an image file.
- Show the accuracies and confusion matrices for the test set.