PROJECT TITLE

HEART FAILURE PREDICTION

GROUP NO. 10

GROUP MEMBERS

MEMBER-01: MUHAMMAD BAZIL CHOUHAN    (EP#20101027)

MEMBER-02: AZIZ ALI                                  (EP#19101009)

MEMBER-03: SYED AMMAR ALI                    (EP#20101052)

MEMBER-04: YUMAN SAEED                         (EP#20101060)

DEGREE NAME

MASTERS OF COMPUTER SCIENCE(MCS)

YEAR

MCS FINAL YEAR (2020-22)

UNIVERSITY NAME

UNIVERSITY OF KARACHI(UOK)

DEPARTMENT NAME

UMAER BASHA INSTITUTE OF INFORMATION TECHNOLOGY(UBIT)

# Table of Contents

# CHAPTER 1

# INTRODUCTION OF DATA MINING AND WAREHOUSING

We can define data mining and warehousing as an automatic or semi-automatic technical process which gathered the scattered or unreliable data and make that data into knowledgeable.

## APPLICATION OF DATA MINING AND WAREHOUSING

### HEALTH CARE SYSTEM

We can use data mining in health care to reduce cost. It is use in improvement of entire system. Researchers use data mining in database. It can use to predict the volume of patients in hospitals.

### MARKET BASE ANALYSIS

It can use in market for costs reduction. Through this retailer will know what group of items they should keep by model.

### EDUCATION

It is called EDM which purpose is to analysis student`s performance and predict student future learning. Its goal is to predict student`s future behavior and to predict student effectiveness on study.

### MANUFACTORING ENGINEERING

In this its works to solve pattern the complex manufacturing process. It can be very use for to predict the product development span time.

### CRM

Customer relationship management. It is use to maintain the relationship between business and customer for better customer service.

### FRAUD DETECTION

Its use to save billion dollars to from fraud with help of data mining we can convert data into information.

## LIE DETECTION

With the help of this we can detect lies

We collect data of past and present and make unstructured data into information and predict that either criminal speaking a truth or not.

# DIFFERNCE BETWEEN DATA AND INFORMATION

## DATA

It is unorganized form which is useless because it is not in structured form.

Data is not very useful until it is converted into information.

## INFORMATION

It's a form of structured data.

It's very useful because it is arranged and easy so understand.

# ADVANTAGES OF DATA MINING AND WAREHOUSING

It is use in health care, education, marketing, CRM. It is also use in businesses and banking. Information mining have a great deal of benefits when utilizing in a particular industry. Other than those benefits, information mining additionally has its own disservices e.g., protection, security, and abuse of data.

# CHAPTER 2
# LITERATURE REVIEW

Enormous information is arising as a very famous topic among experts and researchers. For models, many organizations utilize advanced advances to follow web-based media consistently, along these lines making longitudinal designs of millions of posts, tweets, or surveys George et al. (2014). Contemporaneously, an enormous number of examination papers are devoted on the progressive of large information in various industry and utilize 3Vs to portray huge information, the outrageous volume of information, the wide assortment of information types and the speed at which the information should be processed. Definitely, individuals from very different backgrounds these days are living in the period of enormous information.

In 2014, Bank of America designated another high level position, the Chief Analytics Officer(CAO), to Douglas Hague who is the previous senior VP of merchant analytics(Ferguson, 2014). He unequivocally communicated an acknowledgment of the significance of investigation in the utilization of enormous information, rather than simply writing about data(Ferguson, 2014). Exactly, organizations need to pick up taking investigation past the information and data which could come into the genuine bits of knowledge that sway direction. Sanders (2016 p.28) precisely gave the contrast between enormous information and examination, "Large information without investigation is only a huge measure of information. Examination without enormous information are basically numerical and factual apparatuses and applications". Henceforth, arranged by enormous information, we should place more endeavors on huge information analytics(BDA), which could be characterized as " another age of advances and structures, intended to monetarily remove esteem from exceptionally huge volumes of a wide assortment of information, by empowering high speed catch, disclosure and/or examination." (Carter, 2011).

Alongside the innovation of enormous information investigation being developed bit by bit, the idea of business knowledge emerged through the entire financial globe. Endeavors will more often than not investigate a wide range of profoundly concealed highlights from mass data

organizations, with pervasive and wise examination abilities, for example, information mining, process mining, web mining or text mining.

# CHAPTER 3
# RANDOM FOREST
## INTRODUCTION

Arbitrary timberlands or irregular choice woodlands are a group learning strategy for characterization, relapse and different undertakings that works by building a huge number of choice trees at preparing time. For relapse errands, the mean or normal expectation of the singular trees is returned. Random forest is bagged decision trees that built tree ang make large collection of decorrelated trees. its use to improve prediction of data mining. It's a very enjoyable algorithm to use.

It is a supervised machine learning Algorithm used in classification and regression problems. It takes data randomly from the original dataset and create decision trees and based on average of the output generated by decision trees makes final decision.

Random forest works on the bagging principle.

## BAGGING

It is also known as bootstrap aggregation. We create bootstraps by taking information randomly from the original data. The process of combining the output of the Decision tree which we created from the bootstrap and average the majority votes from the decision tree and make decision known as aggregation.

For Example:

A student wants to take admission in a specific course and he takes the advice about the course from the professional of that field and students of that particular course and gets the different voles, now he averages the votes and take decision accordingly.

## STEPS INVOLVE IN RANDOM FOREST ALGORITHM:

1. In random forest n number of random records taken from the k number of records from the dataset.
2. Individual decision trees are constructed from each sample.
3. Each decision tree will generate an output.

4. The final out will be considered by majority votes or averaging for classification and regression problems

# DECISION TREE

## INTRODUCTION

If we can say that decision tree is process of learning a function that maps data item into one of several pre-defined classes. Every classification based on inductive-learning algorithm.in this we create a model known as classifier in which it predicts the available values of their input attributes. Decision tree is a powerful method to solve the real-life mythologies.

## EXAMPLES OF CLASSIFICATION

PROBLEMS:

Classifying trends in financial market

Identified objects in large images

## ADVANTAGES OF DECISION TREE

One of its advantages is their output is to read

It does not require normalization of data.

It does not require scaling of data

It's very easy to explain someone like technical teams or stakeholders

## DISADVANTAGES OF DECISION TREE

Sometimes its algorithm goes in complex solution

It's a high time consume model

Its take higher time to train the mode.

## PRUNING

Pruning is an information pressure procedure in AI and search calculations that lessens the size of choice trees by eliminating areas of the tree that are non-basic and excess to order cases.

## CLASSIFICATION MODEL

In a Classification model, the association among classes and different properties of the examples can be characterized by a flowchart or as complicated and unstructured manual. Data-mining systems limit conversation to formalized, "executable" models of grouping. There are two altogether different manners by which they can be built. - Interviewing the pertinent master or specialists, and most information-based frameworks have been fabricated thusly (regardless of

the notable troubles orderly on adopting this strategy.) - Alternatively, various recorded orders may be analyzed and a model developed inductively by summing up from explicit models that are of essential interest for information mining applications.

# EXTENDED FOREST

Random forest uses same principle as decision tree. Random forest introduces us bagging. Random forest teaches us how to make leave trees bigger. Random forest use aggregates to improve prediction.

# OUT OF THE BOX

It's very popular because it provides out of the box performance.

Random forest has become least variability in their prediction.

# HYPERPARAMETER

We should consider hyperparameter when we train models.

There are main hyperparameters.

1. Trees in forest
2. Complexity of tree
3. How to do sampling scheme

# TREE IN FOREST

First, we should consider number of trees in random forest.

No of trees should be large to control error.

# COMPLEXITY IN TREE

In random forest there are one or more hyperparameter to control the complexity of tree.

# SAMPLING SCHEME

In random forest sampling scheme is bootstrapping where total 100 percent data is observed.

# TUNNING STRATEGIES

As we introduce more complex algorithm, we should me strategies.

We should know how to proceed our grid search.

# SPLITE RULE

We should recall the default split rule which we use in decision tree to build tree.

# CONCLUSION

1. Random forest provides out of shell prediction performance.

2. Random forest has more accuracy

3. Random forest also faster than bagging

4. It creates stability between correlation trees

# CHAPTER 4

# DBSCAN ALGORITHM

We can define DBSCAN as dense region covered or surrounded by the low-density region

It selects the no of clusters automatically to select core there should be 3 points in circle

The center point of circle is core. The neighbor points of core in circle are boundary points. The noise lies anywhere in circle.

How it works?

You should select a point p randomly from arbitrarily. Find epsilon neighborhood of point p.

If they are core points then the cluster is formed.

Then search for other points in the data set if they turn out to be a border points then move on and search for other points in the database.

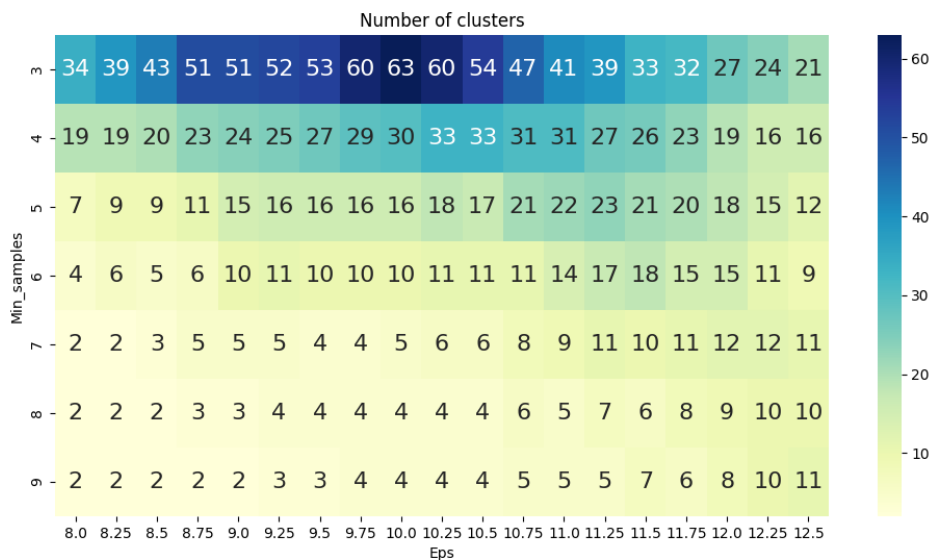Repeat these actions until all the points in the database are processed.
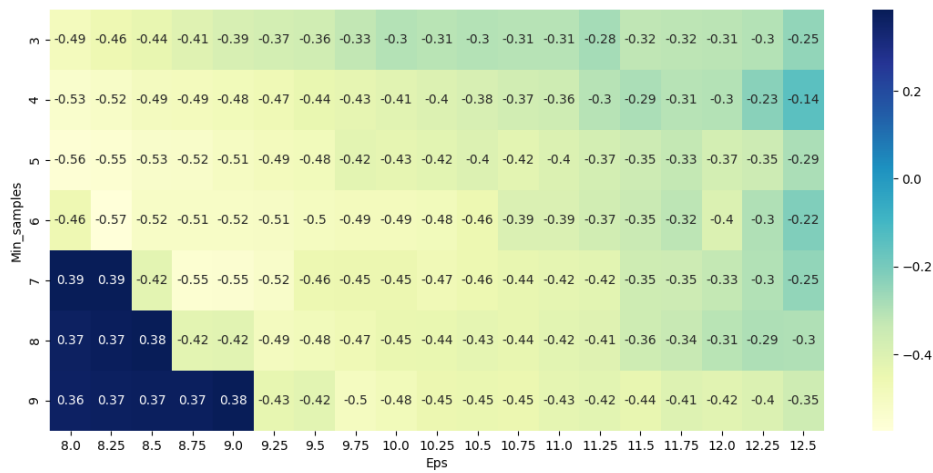


Figure 1

Figure 2
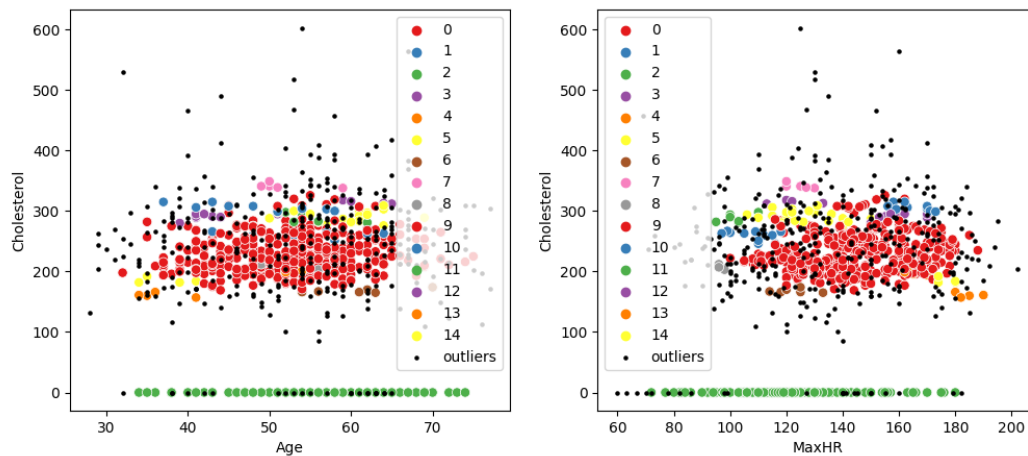


Figure 3

## TIME COMPLEXITY

The time complexity is O(m) and in worst case the time complexity will be O(m)^2.

## FUZZY K MEANS

Fuzzy k means we can say that it is an observation of soft clustering that can be assigned to multiple clusters.

## ADVANTAGES:

- It's very easy for implement
- Its scale very large datasets
- Guarantee of convergence

## DISADVANTAGES:

- Choose k manually it's a big disadvantage
- Always depends on initial values
- Clustering outliers

## CENTER BASED CLUSTERING

We define that goal of center-based clustering is to minimize the objective function. It's very efficient and work good on large or high dimensional data sets. We assume that cluster should be in convex shaped. We can represent cluster center as cluster. The example of this cluster is gaussian mixtures which is widely use in this clustering. We use Euclidean distance in this clustering.

Error generated in clustering which we can define as the sum of distance from the clusters Centre. K –means clustering use in in center clustering. If we talk about optimization in k-means its procedure is very greedy. In every iteration the of objective function decreases. It's working depends only on numeric values. If we see the clusters or identified so, we found ball-shaped clusters. The initial values make or depend on result.
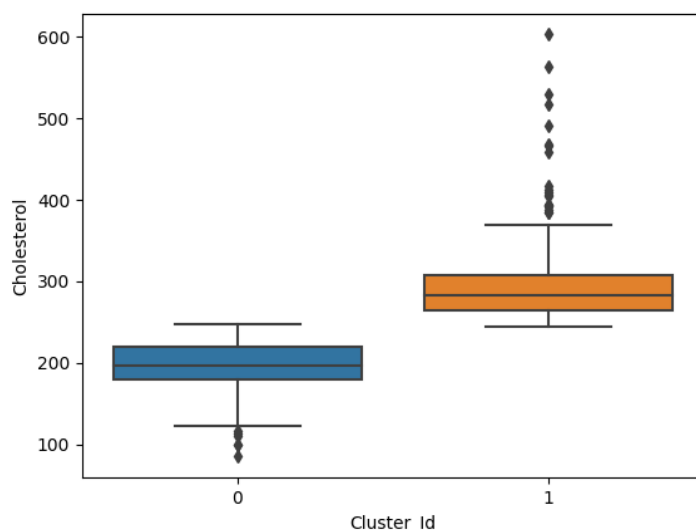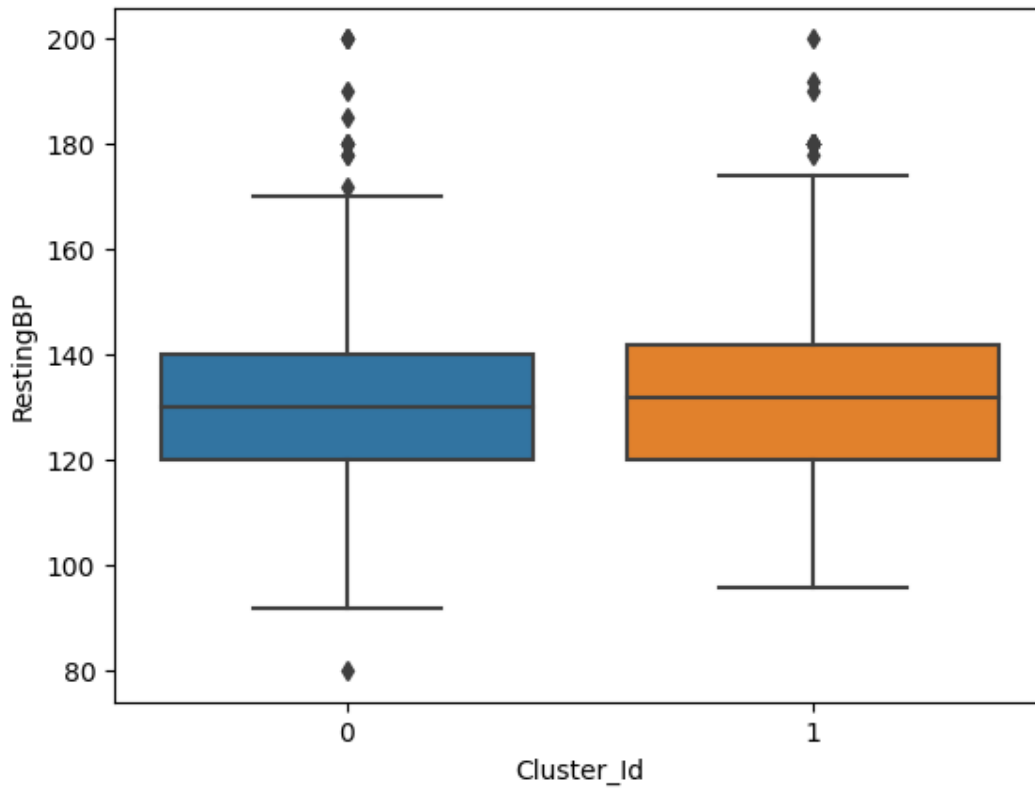
Figure 4

Figure 5

## CONCLUSION

In this chapter we learn about DBSCAN, center base clustering and fuzzy k-means technique

# CHAPTER 5

# DATA PREPROCESSING

It is done to improve the quality of data in data warehouse. Data preprocessing in used for increasing efficiency and mining process is easy. And it is also removed noisy data, incomplete data and inconsistent data. Data in the real world is dirty.

Inconsistent: containing discrepancies in order or names. Data cleaning: It clean the data by filling in the missing values. Sometime noisy data resolving the inconsistency and removing the outliers. Data transformation: It is Data processing techniques that transfer or consolidate the data into alternate appropriate for mining. Data integration: Its is preprocessing method that involves merging of data form different sources in order to form a data solve like data warehouse.
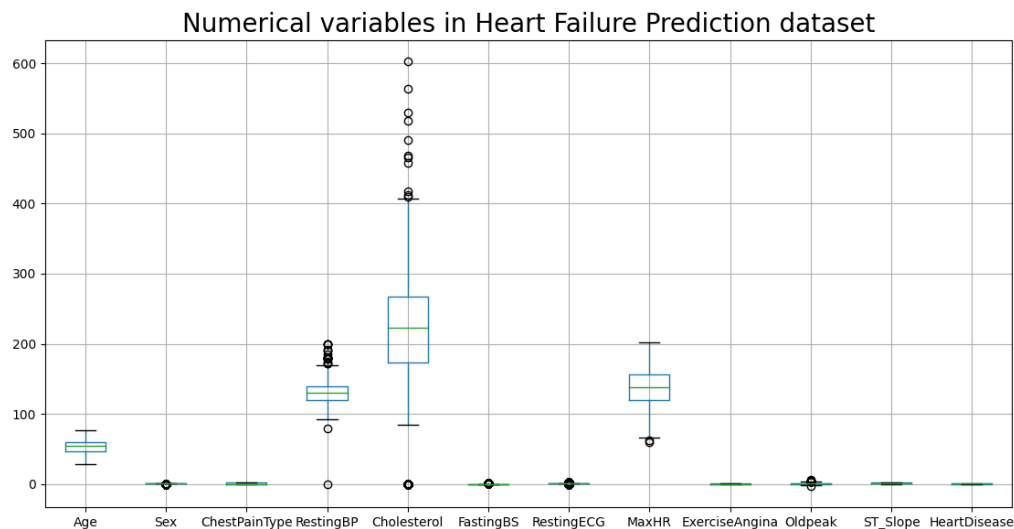
## OUTLIER DETECTION



Figure 6

## OUTLIER CORRECTION

Numerical variables in Heart Failure Prediction dataset

Figure 7

Issues in data integration

## INTEGRATION AND OBJECT MATCHING

For example, we have two companies A and B. Company A store employee's data in name with "employee ID" while company B store employee data in the name with "employee number". In this scenario both are same attributes but how the system understand that both are same. So, scheme integration and object matching are major problem in such a situation.

In redundancy some unwanted attributes are included in our data that are not needed so we have to remove such data which Are not necessary.

## DETECTION AND RESOLUTION OF DATA VALUE CONFLICT

Convert and modify the values. For example if we have if company show the amount in "RUPEES" while other company show in "DOLLAR" such in such a situation we have to correctly modify such a values.

Data reduction: reprocessing techniques that helps in obtaining related representation of data set from the available dataset.

- Integrity of original data should own after non reduction in data value
- It should produce same analytic result as an original data.
- Dimensionality Red: remove redundant attribute.
- Data compression: encoding of data.

# DATA REDUCTION STRATEGIES

## DATA CUBE AGGREGATION

The lowest level of data cube base cuboid. The aggregation data for an individual entity of interest. Multiple levels of aggregation in data cube.

## DIMENSIONALITY REDUCTION

It eliminates redundant attribute which are weakly important across the data.

- Step forward selection.
- Stepwise backward elimination
- Decision tree induction.

## NUMEROSITY REDUCTION

Replace the original data with small form of data representation. There are two methods.

## PARAMETRIC

Parametric reduction: used to eliminate the data so that only parameter of data is required store instead of actual data.

## NON-PARAMETRIC REDUCTION

Non-Parametric reduction: used to store thousand representations of data. It includes histogram, clustering, sampling and data aggregation.

# DATA DISCRETIZATION

It divides the range of attribute into intervals so as to reduce number of values for a given continuous attribute.

- Splitting
- Merging
- Supervised
- Un supervised

# SPLITTING

It's a top-down approach where attribute is split in to range of values.

# MERGING

It's a top-down approach, initially we consider all values later remove some during merging.

# SUPERVISED

If we know the class which they are lying the is called supervised

# UNSUPERVISED

If we have no idea of class which they are lying it is called unsupervised.

# CONCLUSION

In this chapter we learn about data preprocessing.

# CHAPTER 6
# LOGISTIC REGRESSION:

Logistic regression is a method used to recognized a data value based on observations of a data set. Logistic regression now a days become an important tool to use.  The approach allows an algorithm being used in a machine learning application. A logistic regression model now a days predicts a dependent data variable by analyzing and make the relationship between one or more existing independent variable. For example, a logistic regression could be used to predict whether a team will win or lose a match or whether a college student will be admitted to a particular university.

## PURPOSE AND EXAMPLE OF LOGISTIC REGRESSION

Logistic regression is one of the most commonly used machine learning for binary classification problems, which are problem with two class values, including prediction such as this or that yes or no or A and B.

## NON-LINEAR REGRESSION

Nonlinear regression is a form of regression analysis in which data is fit to a model and then expressed as a mathematical function simple linear regression relates two variables (X and Y) with a straight line (y= mx+b), while nonlinear regression relates the two variable is a nonlinear (curved) relationship. The goal of the model is to make the sum of the squares as Small as possible. The sum of squares is a measure that attacks how far the Y observations vary from the nonlinear (curved) function that is used to predict Y. Nonlinear regression modelling is similar to linear regression modelling is that both seek to track a particular response from a set of variables graphically. Nonlinear model is more complicated than linear model to develop because the function is created through a series of approximations that may stem from trial and error.

Logistic regression with one predictor

Example Rizatriptan

Rizatriptan is a headache medicine that narrow blood vessel around the brain it is also used reduce the substance in the body make the headache, nausea and also many symptoms. Rizatriptan is used to treat migraine headache. It is not used to treat the tension or any other headache. Use Rizatriptan at that movement if doctor confirm of migraine headache.

Multiple logistic regression

Multinomial logistic regression is used to predict in probability of categorically membership and categorical placement on dependent variable on a multiple independent variable. The type of dependent variable   may be binary or continuous or ratio scale. Multinomial logistic regression is simple logistic regression that allow more than two categories of the dependent of outcome variable. Like binary logistic regression and multinomial logistic regression.

Examples of multiple logistic regression

Example 1 : entering high school student make choices to among in general program like vocational and academic program. Their choice might be modeled using writing score and their social economic status.

Nonlinear regression

It's a form of regression data is fix and fit to a model and then make or expressed as mathematical function if for example of we have a linear regression have two variables (X and Y) with a  straight line (y=mx+b), while in nonlinear regression has two variables in a nonlinear relationship.

Examples p24 antigen

These are sometimes called fourth generation tests

One distinctive HIV antigen is a viral protein called p24, making of p24 antigen assays useful in diagnosing primary HIV infection.

AZT therapy

AZT therapy were tested by means of. Newly developed plaque assay in cd4+ heal cells. Fifty percent inhibitory Dose values of meter 0.05 nano meter.

Mk-639 in HIV patients

Mk-639 is a potent and specific HIV protease to have good oral bioavailability. target HIV protease mk-639 is a protease inhibitor used a component of high active antiretroviral therapy to treat HIV infection and AIDS MK-639 appear to have significant dose related antiviral activity and is well tolerated.

## CONCLUSION:

In this chapter we talk about logistic regression and find how to determine logistic regression.

# CHAPTER 7

# FUZZY C-MEAN

## INTRODUCTION:

We call it extension of k-means. In fuzzy c-mean it allows you to assigned data set into more than one cluster. In fuzzy c-mean data point has probability to belonging to each cluster. Fuzzy c-mean is better than k-mean.

Fuzzy-c-means clustering can be viewed as a superior calculation contrasted with the k-Means calculation. Not at all like the k-Means calculation where the information focuses only on a place with one group, on account of the Fuzzy-c-means calculation, the data point can have a place with more than one cluster with a probability

## OBJECTIVES OF FUZZY-C-MEANS

Mohsen Ghanea. Iranian National Institute for Oceanography. Minimizing objective function means increasing similarity among all the components within an object and reducing similarity between components of one object with others.

## MEMBERSHIP VALUES IN THE FUZZY C-MEANS CLUSTERING

Membership grades are allotted to every one of the main items (labels). These participation grades demonstrate how much information focuses have a place with each group. Accordingly, focuses on the edge of a cluster, with lower participation grades, might be in the group less significantly than focuses in the focal point of cluster.

## CALCULATION of FUZZY C- MEAN

- First you have to find initial cluster centers by using formula
- After you find centers now you find distance to dataset from each data set
- After that you find distance, you update the new membership matrix using formula
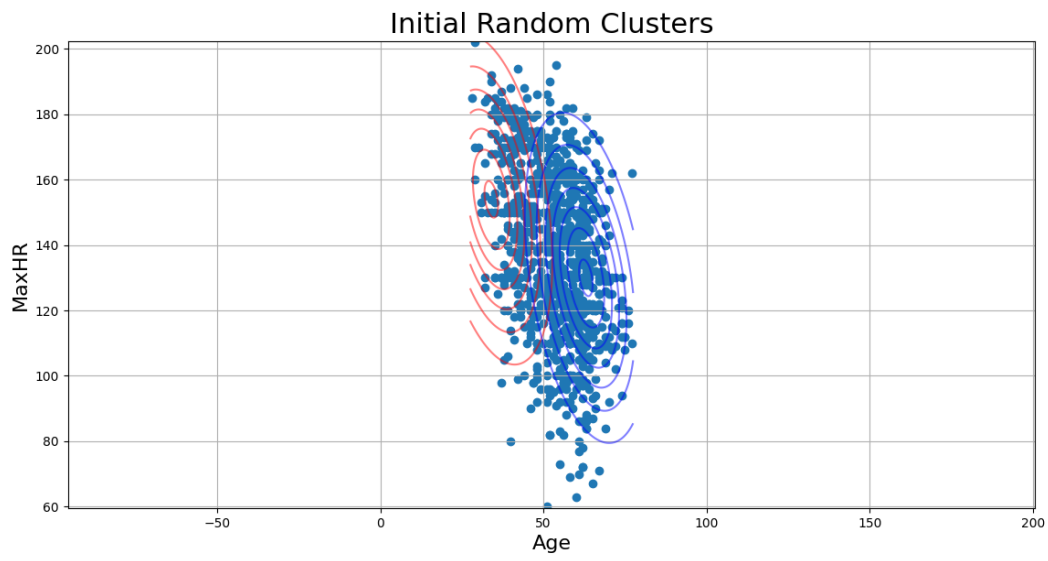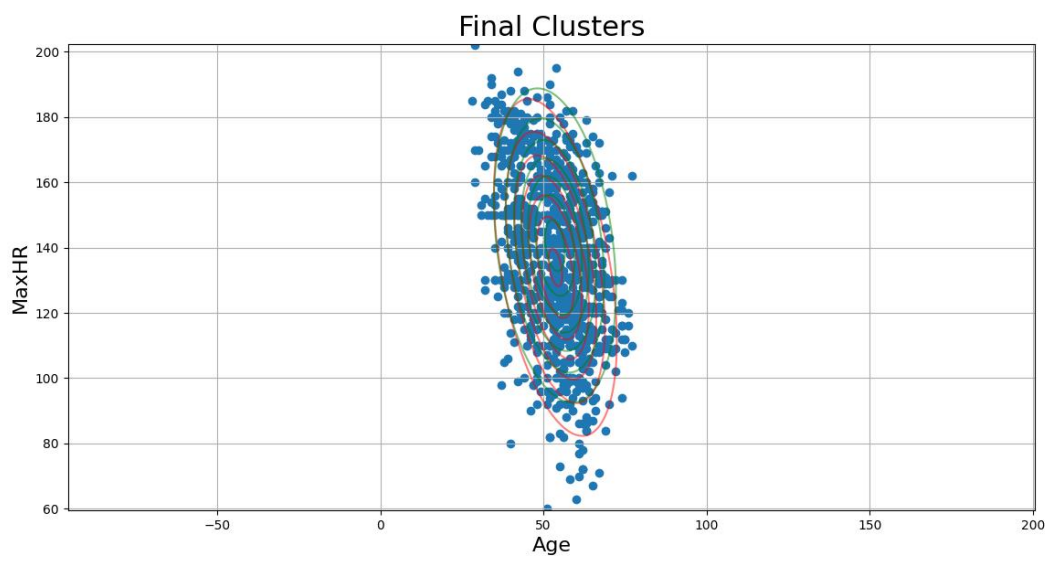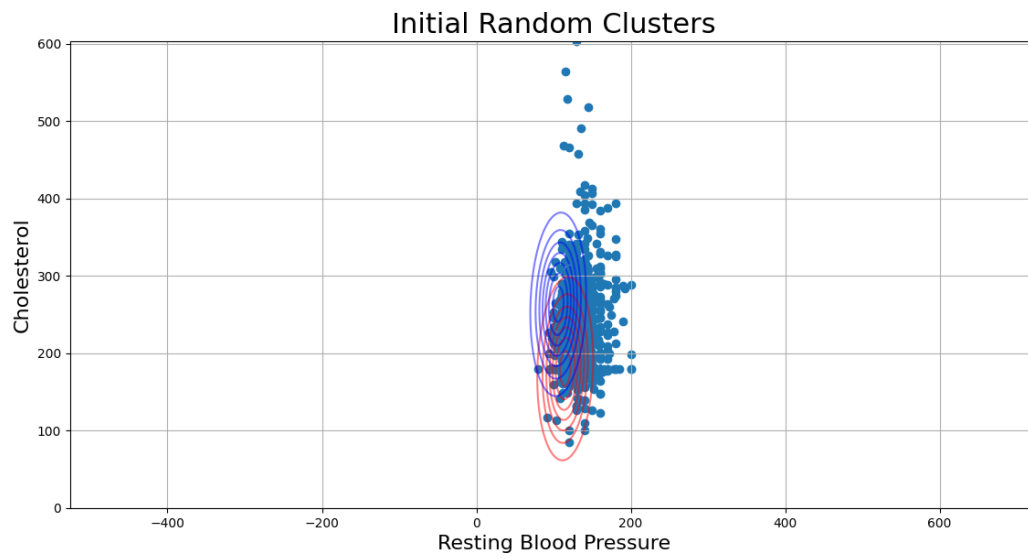- Do it after and after until you get same centroids.
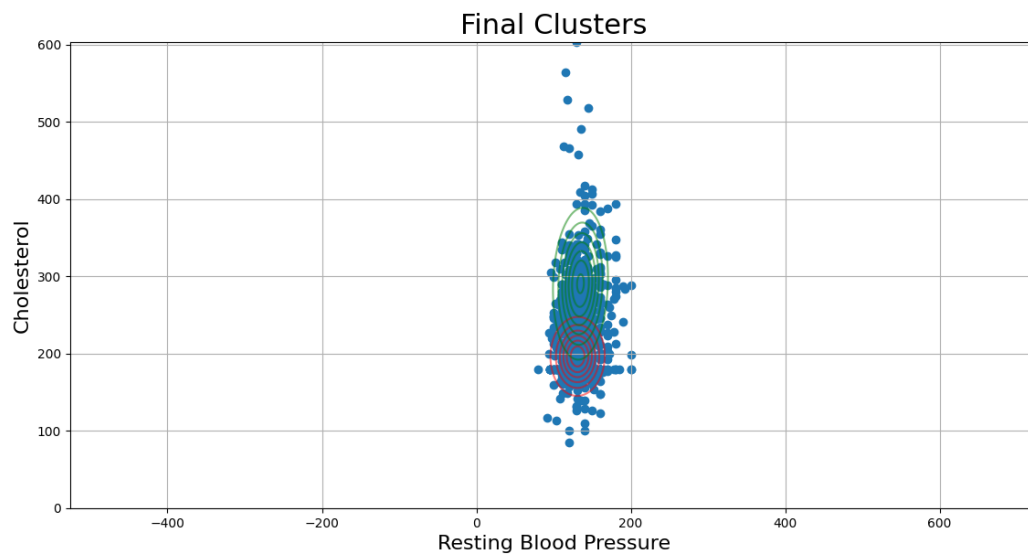
Figure 8



Figure 9

Figure 10



Figure 11

# CONCLUSION

Fuzzy c-mean much better than k-means

In fuzzy c-mean cluster belong to each other and there can be more one cluster assigned.

# CHAPTER 8

# ENTROPY AND CLUSTER

## ENTROPY FUNCTION ALGORITHM

Entropy measures the impurity of a collection of examples. It relies from the distribution of the random variable as we know

- S is a collection of training examples.
- P is Proportion of positive example in S
- P is Proportion of negative example in S
- Entropy and information gain

When Information gain measures the expected or outcoming reduction in entropy

Given entropy as a measure of impurity is a collection of training example, the information gains the expected reduction in entropy caused by partitioning the example according to a variable.

More precisely the information gain (S, A) of attribute A relative to a collection of example S in defined as

Gain (S, A) = entropy(S) - Sv\S (entropy S)

## PRINCIPAL COMPONENT ANALYSIS

A component analysis is concerned that explain the variance, covariance structure of a set through a few linear combinations of these variables.

Principal Analysis is an unsupervised learning algorithm that is used for the dimensionality reduce or make machine learning. It is a process that converts the observations of correlated into a set of linearly uncorrelated with the help of orthogonal transformation. These new converted features are called the Components. It is one of the famous tools that is used for explore and analysis data and predictive modeling. It is a technique to draw strong patterns from the given dataset by reducing the variances.

## FACTOR ANALYSIS

Factor analysis is a technique to combine question or variable to create new factor.

Purpose

Data reduction: to reduce the number of variables to a manageable set of factors

Substantive interpretation: to identify underlying constructs in the data. This make easier to understand the data.

## CLUSTER ANALYSIS

Cluster analysis is all about finding observations which are similar so that when you group them together, they should be homogeneous within group. And each group should be quite distinct or in other word any observations of one group should be quite dissimilar to any object of other group.

## TYPES OF CLUSTERING METHOD

As we know that for a successful grouping, we need to attain two major goals: one, a similarity between one data point with another and two, a distinction of those similar data points with others which most certainly, heuristically differ from those points. The basis of such divisions begins with our ability to scale large datasets and that's a major beginning point for us. Once we are in it, we are represented with a challenge that our data sustain different kinds of attributes – categorical, continuous data, etc., and we should be able to deal with them. Now, as we know that now a our data these days is not limited in terms of dimensions.

The clusters that we require, should not only be able to distinguish data points but also, they should be inclusive. We are Sure, a distance metric will help a lot but the cluster shape is often limited due to a geometric shape and that's a reason many important data points get excluded. This problem too needs to be taken care of.

In our progress, we seeing that our data is highly "noisy" in nature. Many unwanted features have been residing in the data which makes it rather Herculean task to bring about any similarity between the data I points – leading to the creation of improper groups. As we move towards the end of the line, we are faced with a challenge of business interpretation. The outputs from the clustering algorithm should be understandable and should fit the business criteria and address the business problem correctly.

The various types of clustering are:

Connectivity-based Clustering (Hierarchical clustering)

Centroids-based Clustering (Partitioning methods)

Distribution-based Clustering

Density-based Clustering (Model-based methods)

Fuzzy Clustering

Constraint-based (Supervised Clustering)

# CONNECTIVITY-BASED CLUSTERING (HIERARCHICAL CLUSTERING)

Hierarchical Clustering is a method of unsupervised machine learning clustering where it begins with a pre-defined top to bottom hierarchy of clusters. It then proceeds to perform a decomposition of the data I objects based on this hierarchy, hence obtaining the clusters. This method follows two approaches based on the I direction of progress, i.e., whether it is the top-down or bottom-up I flow of creating clusters. These are Divisive Approach and the Agglomerative Approach respectively.

## CLUSTERING DENDOGRAM
## SINGAL LINKAGE

Figure 12

**COMPLETE LINKAGE**



Figure 13

## AVERAGE LINKAGE



Figure 14

## DIVISIVE APPROACH

This approach of hierarchical I clustering follows a top-down approach where we consider that all the data points belong to one large cluster and try I to divide the data into smaller groups based on a termination logic or, a point beyond which there will be no further division of data points. This termination logic can be based on the minimum sum of squares of error inside a cluster.

## CONCLUSION

In this chapter we learn about Entropy,Clusters and their types.

# CHAPTER 9
# RESULTS AND CONCLUSIONS

# DATASET

We took this dataset from www.kaggle.com.The dataset is about heart disease. This dataset will allow us to predict that someone has heart disease based on these variables ('Age', 'Sex',' ChestPainType', 'RestingBP', 'Cholesterol', 'FastingBS', 'RestingECG', 'MaxHR', 'ExerciseAngina', 'Oldpeak', 'ST_Slope', 'HeartDisease'). The data are organized in rows and columns and each row contains a mixture of rows and columns.



Figure 15

According to our dataset the figure 1 shows that heart rate is increased by age, As we can see the age group of between 50-60 have highest maximum heart rate, which lies under 80-180.

# PRE PROCESSING

Our dataset contains columns: Age, Sex, ChestPainType, RestingBP, Cholesterol, FastingBS, RestingECG, MaxHR, ExerciseAngina, Oldpeak, ST_Slope, HeartDisease.

We used boxplot to detect the outliers, the outliers were present in RestingBP, Cholesterol and MaxHR attribute. We corrected the outliers in RestingBP by assigning the average BP of a person according to their age and gender. We assigned Average cholesterol 180 to deal with outliers in Cholesterol. The heart rate outliers were not actually outliers as we know a person heart rate can go down to 80, 81. It lies under the normal range of a person heart rate.

We used boxplot again to show all the removed outliers. We also used different types of plots to understand our data more specifically.

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.graphics.gofplots import qqplot
import seaborn as sns
import warnings
from sklearn.preprocessing import LabelEncoder
import random
import operator
import math
from scipy.stats import multivariate_normal     # for generating pdf

warnings.filterwarnings('ignore')

df=pd.read_csv('heart.csv')




#to covert categorical data into numeric
columns=['Sex','ChestPainType','RestingECG','ExerciseAngina','ST_Slope']
encoders = {}
for col in columns:
    le = LabelEncoder().fit(df[col])
    df[col] = le.transform(df[col])
    encoders[col] = le
```

```python
33    #OUTLIER DETECTION
34    num_cols = ['Age','Sex','ChestPainType','RestingBP','Cholesterol','FastingBS','RestingECG','MaxHR',
35                'ExerciseAngina','Oldpeak','ST_Slope','HeartDisease']
36    plt.figure(figsize=(18,9))
37    df[num_cols].boxplot()
38    plt.title("Numerical variables in Heart Failure Prediction dataset", fontsize=20)
39    plt.show()
40
41
42
43
44
45
46
47    #OUTLIER CORRECTION
48    print(df[df.RestingBP == 0])
49
50    if (df.at[449, 'Age']>=18) and (df.at[449, 'Age']<=39):
51        if df.at[449, 'Sex']=='M':
52            df.at[449,'RestingBP'] = 119
53        else:
54            df.at[449,'RestingBP'] = 110
55    elif (df.at[449, 'Age']>=40) and (df.at[449, 'Age']<=59):
56        if df.at[449, 'Sex']=='M':
57            df.at[449,'RestingBP'] = 124
58        else:
59            df.at[449,'RestingBP'] = 122
60    else:
61        if df.at[449, 'Sex']=='M':
62            df.at[449,'RestingBP'] = 133
63        else:
64            df.at[449,'RestingBP'] = 139
```

```python
65
66    print(df.iloc[449])
67
68    print(df[df.Cholesterol == 0])
69    list1=[]
70  ∨ for i in range(len(df)):
71  ∨     if df.at[i,'Cholesterol'] ==0:
72            list1.append(i)
73
74  ∨ for j in list1:
75        df.at[j,'Cholesterol'] = 180
76    |
77  ∨ num_cols = ['Age','Sex','ChestPainType','RestingBP','Cholesterol','FastingBS','RestingECG','MaxHR',
78              'ExerciseAngina','Oldpeak','ST_Slope','HeartDisease']
79    plt.figure(figsize=(18,9))
80    df[num_cols].boxplot()
81    plt.title("Numerical variables in Heart Failure Prediction dataset", fontsize=20)
82    plt.show()
83
84
85
86  ∨ def plot_dist(col, ax):
87  ∨     if col != 'height':
88            df[col].value_counts().plot(kind='bar', facecolor='y', ax=ax)
89  ∨     else:
90            df[col].plot('density', ax=ax, bw_method = 0.15, color='y')
91            ax.set_xlim(130,200)
92            ax.set_ylim(0, 0.07)
93        ax.set_xlabel('{}'.format(col), fontsize=18)
94        ax.set_title("{} on HFP".format(col), fontsize= 18)
95        return ax
```

```python
97    f, ax = plt.subplots(2,4, figsize = (18,9))
98    f.tight_layout(h_pad=9, w_pad=2, rect=[0, 0.03, 1, 0.93])
99    cols = ['Age','Sex','ChestPainType','RestingBP','Cholesterol','FastingBS','RestingECG','MaxHR',
100             'ExerciseAngina','Oldpeak','ST_Slope','HeartDisease']
101   k = 0
102   for i in range(2):
103       for j in range(4):
104           plot_dist(cols[k], ax[i][j])
105           k += 1
106   __ = plt.suptitle("Final Distributions of different features", fontsize= 23)
107
108
109   plt.figure(figsize=(18,8))
110   plt.xlabel("Age", fontsize=18)
111   plt.ylabel("Max Heart Rate", fontsize=18)
112   plt.suptitle("Joint distribution of Age vs Max Heart Rate", fontsize= 20)
113   plt.plot(df['Age'], df['MaxHR'], 'bo', alpha=0.2)
114   plt.show()
115
```

# DECISION TREE

The above dataset is used in Decision tree.

We have split the data into dependent and independent variables. We have used the data in columns to predict one variable (Heart Disease).

We have used One Hot Encoding (A *one hot encoding* is a representation of categorical variables as binary vectors).

We split the data into training and testing sets and we plot a tree and it was huge so, then we use cost complexity pruning to simplify the tree which has more accuracy and less complicated to understand.

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.graphics.gofplots import qqplot
import seaborn as sns
from sklearn.tree import DecisionTreeClassifier
from sklearn.tree import plot_tree
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import plot_confusion_matrix


df=pd.read_csv('D:\\M.CS Studies\\4th Semester\\Data Mining\\Project\\Decision Tree\\heart.csv',header=None)

df.columns = ['Age','Sex','ChestPainType','RestingBP','Cholesterol','FastingBS',
              'RestingECG','MaxHR','ExerciseAngina'
              'Oldpeak','ST_Slope','HeartDisease',]


X=df.drop('HeartDisease', axis=1).copy()
y=df['HeartDisease'].copy()
y.head()
X.dtypes
X['ChestPainType'].unique()
pd.get_dummies(X, columns=['ChestPainType']).head()
X_encoded = pd.get_dummies(X, columns=['Sex','ChestPainType','RestingECG','ST_Slope','ExerciseAngina'])
X_encoded.head()

y_zero_index = y > 0
y[y_zero_index] = 1
```

```python
33    X_train, X_test , y_train, y_test = train_test_split (X_encoded ,y, random_state=45)
34    clf_dt = DecisionTreeClassifier(random_state=45)
35    clf_dt = clf_dt.fit(X_train, y_train)
36
37    plt.figure(figsize=(40,20))
38    plot_tree (clf_dt,filled=True,rounded=True,class_names=["No HD", "Yes HD"] ,feature_names=X_encoded.columns);
39    plot_confusion_matrix(clf_dt, X_test, y_test, display_labels=["Does not have HD", "Has HD"])
40    path =clf_dt.cost_complexity_pruning_path(X_train , y_train)
41
42    ccp_alphas = path.ccp_alphas
43    ccp_alphas = ccp_alphas[:-1]
44    clf_dts = []
45
46    for ccp_alpha in ccp_alphas:
47        clf_dt = DecisionTreeClassifier(random_state=0 , ccp_alpha=ccp_alpha)
48        clf_dt.fit(X_train, y_train)
49        clf_dts.append(clf_dt)
50
51
52    train_scores = [clf_dt.score(X_train , y_train) for clf_dt in clf_dts]
53    test_scores = [clf_dt.score(X_test , y_test) for clf_dt in clf_dts]
54
55    fig, ax = plt.subplots()
56    ax.set_xlabel("Alpha")
57    ax.set_ylabel("Accuracy")
58    ax.set_title("Accuracy vs Alpha for testing and training sets")
59    ax.plot(ccp_alphas, train_scores, marker='o', label="train", drawstyle="steps-post")
60    ax.plot(ccp_alphas, test_scores, marker='o', label="test", drawstyle="steps-post")
61    ax.legend()
62    plt.show()
```

```python
63
64
65    clf_dt = DecisionTreeClassifier(random_state=45 , ccp_alpha=0.016)
66    scores = cross_val_score (clf_dt, X_train , y_train, cv=5 )
67    df = pd.DataFrame(data={'tree': range(5), 'accuracy': scores})
68
69    df.plot(x='tree', y='accuracy', marker='o', linestyle='--')
70
71    alpha_loop_values= []
72
73    for ccp_alpha in ccp_alphas:
74        clf_dt = DecisionTreeClassifier(random_state=0 , ccp_alpha=ccp_alpha)
75        scores = cross_val_score (clf_dt, X_train , y_train, cv=5 )
76        alpha_loop_values.append([ccp_alpha, np.mean(scores), np.std(scores)])
77
78    alpha_results=pd.DataFrame(alpha_loop_values,columns=['alpha', 'mean_accuracy','std'])
79    alpha_results.plot(x='alpha', y='mean_accuracy', yerr='std',marker='o', linestyle='--')
80
81    alpha_results[(alpha_results['alpha'] > 0.015)
82                    &
83                    (alpha_results['alpha'] < 0.016)]
84
85    ideal_ccp_alpha = alpha_results[(alpha_results['alpha'] > 0.015)
86                                        &
87                                        (alpha_results['alpha'] < 0.016)]['alpha']
88
89
90    ideal_ccp_alpha = float(ideal_ccp_alpha)
91    ideal_ccp_alpha
92
```

```
 93    clf_dt_pruned= DecisionTreeClassifier(random_state=45 , ccp_alpha = ideal_ccp_alpha )
 94    clf_dt_pruned = clf_dt_pruned.fit(X_train , y_train)
 95
 96    plot_confusion_matrix(clf_dt_pruned, X_test , y_test , display_labels=["Does not have HD" , "Has HD"])
 97    plt.figure(figsize=(15,7.5))
 98    plot_tree(clf_dt_pruned, filled=True, rounded=True,class_names=["No HD", "Yes HD"] ,
 99            feature_names=X_encoded.columns);
100
```

# DBSCAN

DBSCAN stands for Density-Based Spatial Clustering of Applications with Noise. In this type of clustering method, clusters are formed by locating all the points in a dense region surrounded low density points. This algorithm uses epsilon ( ε ) behaves as a radius of a circle, and minimum points which tells that minimum how many points are required to form a cluster. When a cluster is formed by selecting the core point, then all the points having shorter distance from the epsilon are added to the cluster as boundary points. Continue the process until all the points are processed. Euclidean distance formula is used to find the distance between points and epsilon. It contains noise points also which behave as outliers.

We implemented DBSCAN on our dataset. We created matrix of investigated combinations to select the best then generated the clusters. We used heatplot to show how many clusters are formed. Then we found the size of each cluster to have a better idea. In last 15 clusters were formed. We used Scatter Plot to show the cluster.

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.graphics.gofplots import qqplot
import seaborn as sns
import sklearn
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
import warnings
from sklearn.preprocessing import LabelEncoder
from sklearn.cluster import DBSCAN
from itertools import product


warnings.filterwarnings('ignore')

dataset=pd.read_csv('heart.csv')




#to covert categorical data into numeric
columns=['Sex','ChestPainType','RestingECG','ExerciseAngina','ST_Slope']
encoders = {}
for col in columns:
    le = LabelEncoder().fit(dataset[col])
    dataset[col] = le.transform(dataset[col])
    encoders[col] = le

X_numerics=dataset[['Age','Sex','ChestPainType','RestingBP','Cholesterol',
                    'FastingBS','RestingECG','MaxHR','ExerciseAngina','Oldpeak','ST_Slope','HeartDisease']]
```

```
35    eps_values = np.arange(8,12.75,0.25) # eps values to be investigated
36    min_samples = np.arange(3,10) # min_samples values to be investigated
37    DBSCAN_params = list(product(eps_values, min_samples))
38
39
40    no_of_clusters = []
41    sil_score = []
42
43    for p in DBSCAN_params:
44        DBS_clustering = DBSCAN(eps=p[0], min_samples=p[1]).fit(X_numerics)
45        no_of_clusters.append(len(np.unique(DBS_clustering.labels_)))
46        sil_score.append(silhouette_score(X_numerics, DBS_clustering.labels_))
47
48
49    tmp = pd.DataFrame.from_records(DBSCAN_params, columns =['Eps', 'Min_samples'])
50    tmp['No_of_clusters'] = no_of_clusters
51
52    pivot_1 = pd.pivot_table(tmp, values='No_of_clusters', index='Min_samples', columns='Eps')
53
54    fig, ax = plt.subplots(figsize=(12,6))
55    sns.heatmap(pivot_1, annot=True,annot_kws={"size": 16}, cmap="YlGnBu", ax=ax)
56    ax.set_title('Number of clusters')
57    plt.show()
58
59
60
61    tmp = pd.DataFrame.from_records(DBSCAN_params, columns =['Eps', 'Min_samples'])
62    tmp['Sil_score'] = sil_score
63
64    pivot_1 = pd.pivot_table(tmp, values='Sil_score', index='Min_samples', columns='Eps')
```

```
65
66    fig, ax = plt.subplots(figsize=(18,6))
67    sns.heatmap(pivot_1, annot=True, annot_kws={"size": 10}, cmap="YlGnBu", ax=ax)
68    plt.show()
69
70
71    DBS_clustering = DBSCAN(eps=12.5, min_samples=4).fit(X_numerics)
72
73    DBSCAN_clustered = X_numerics.copy()
74    DBSCAN_clustered.loc[:,'Cluster'] = DBS_clustering.labels_ # append labels to points
75
76    DBSCAN_clust_sizes = DBSCAN_clustered.groupby('Cluster').size().to_frame()
77    DBSCAN_clust_sizes.columns = ["DBSCAN_size"]
78    DBSCAN_clust_sizes
79
80
81    outliers = DBSCAN_clustered[DBSCAN_clustered['Cluster']==-1]
82
83    fig2, (axes) = plt.subplots(1,2,figsize=(12,5))
84
85
86    sns.scatterplot('Age', 'Cholesterol',
87                    data=DBSCAN_clustered[DBSCAN_clustered['Cluster']!=-1],
88                    hue='Cluster', ax=axes[0], palette='Set1', legend='full', s=45)
89
90    sns.scatterplot('MaxHR', 'Cholesterol',
91                    data=DBSCAN_clustered[DBSCAN_clustered['Cluster']!=-1],
92                    hue='Cluster', palette='Set1', ax=axes[1], legend='full', s=45)
93
94    axes[0].scatter(outliers['Age'], outliers['Cholesterol'], s=5, label='outliers', c="k")
95    axes[1].scatter(outliers['MaxHR'], outliers['Cholesterol'], s=5, label='outliers', c="k")
96    axes[0].legend()
```

```
96    axes[0].legend()
97    axes[1].legend()
98
99    plt.setp(axes[0].get_legend().get_texts(), fontsize='10')
100   plt.setp(axes[1].get_legend().get_texts(), fontsize='10')
101
102   plt.show()
103
```

# K-MEANS CLUSTERING

K-means clustering is one of the simplest unsupervised machine learning algorithms.

First we initialize clusters by selecting the random points from the data.

Then calculate distance to categorize each point to the cluster having shorter distance from it, then we recalculate the clusters by finding average between the cluster point and the point that have recently been added to the cluster.

Continue this process until all the points have been added to the clusters.

K-mean clustering had implemented to heart failure prediction dataset. No of clusters were 2, Heart Disease (1) and No Heart Disease (0). The result also showed that considering 2 clusters have higher silhouette score. We used boxplot to show the clusters.

```python
1   import numpy as np
2   import pandas as pd
3   import matplotlib.pyplot as plt
4   from statsmodels.graphics.gofplots import qqplot
5   import seaborn as sns
6   import sklearn
7   from sklearn.preprocessing import StandardScaler
8   from sklearn.cluster import KMeans
9   from sklearn.metrics import silhouette_score
10  import warnings
11  from sklearn.preprocessing import LabelEncoder
12
13  warnings.filterwarnings('ignore')
14
15  dataset=pd.read_csv('heart.csv')
16
17
18
19
20  #to covert categorical data into numeric
21  columns=['Sex','ChestPainType','RestingECG','ExerciseAngina','ST_Slope']
22  encoders = {}
23  for col in columns:
24      le = LabelEncoder().fit(dataset[col])
25      dataset[col] = le.transform(dataset[col])
26      encoders[col] = le
27
28
29
30  #K MEAN CLUSTERING
31
32  ssd = []
```

```python
33    range_n_clusters = [2, 3, 4, 5, 6, 7, 8]
34    for num_clusters in range_n_clusters:
35        kmeans = KMeans(n_clusters=num_clusters, max_iter=50)
36        kmeans.fit(dataset)
37
38        ssd.append(kmeans.inertia_)
39
40    # plot the SSDs for each n_clusters
41    plt.plot(ssd)
42    plt.show()
43
44    range_n_clusters = [2, 3, 4, 5, 6, 7, 8]
45
46    for num_clusters in range_n_clusters:
47
48        # intialise kmeans
49        kmeans = KMeans(n_clusters=num_clusters, max_iter=50)
50        kmeans.fit(dataset)
51
52        cluster_labels = kmeans.labels_
53
54        # silhouette score
55        silhouette_avg = silhouette_score(dataset, cluster_labels)
56        print("For n_clusters={0}, the silhouette score is {1}".format(num_clusters, silhouette_avg))
57
58    kmeans = KMeans(n_clusters=2, max_iter=50)
59    kmeans.fit(dataset)
60    kmeans.labels_
61
62    dataset['Cluster_Id'] = kmeans.labels_
63    dataset.head
```

```python
66    sns.boxplot(x='Cluster_Id', y='Age', data=dataset)
67    plt.show()
68    sns.boxplot(x='Cluster_Id', y='Sex', data=dataset)
69    plt.show()
70    sns.boxplot(x='Cluster_Id', y='ChestPainType', data=dataset)
71    plt.show()
72    sns.boxplot(x='Cluster_Id', y='RestingBP', data=dataset)
73    plt.show()
74    sns.boxplot(x='Cluster_Id', y='FastingBS', data=dataset)
75    plt.show()
76    sns.boxplot(x='Cluster_Id', y='RestingECG', data=dataset)
77    plt.show()
78    sns.boxplot(x='Cluster_Id', y='MaxHR', data=dataset)
79    plt.show()
80    sns.boxplot(x='Cluster_Id', y='ExerciseAngina', data=dataset)
81    plt.show()
82    sns.boxplot(x='Cluster_Id', y='Oldpeak', data=dataset)
83    plt.show()
84    sns.boxplot(x='Cluster_Id', y='ST_Slope', data=dataset)
85    plt.show()
86
```

# HIERARCHICAL CLUSTERING

Hierarchical clustering is a type of clustering method in which the clusters are formed by ordering from top to bottom, in a hierarchical way. Two types of hierarchical clustering are used. One is divisive and the other is agglomerative. We are using agglomerative approach in this report.

Three types of methods are used to perform agglomerative clustering

1. Single linkage 2. Complete linkage 3. Average linkage

We used all these 3 methods to have a better idea and result. In single linkage we focus on considering the shortest distance between two points. Unlike the complete linkage hierarchical clustering that focuses on considering the longest distance between two points. Whereas in Average linkage we consider the average distance. In these hierarchical clustering methods we make dendrogram based on k clusters. We implemented single, complete and average linkage hierarchical clustering on our dataset to plot the dendogram based on 2 clusters. And showed all the results by using box plot

```
1   import numpy as np
2   import pandas as pd
3   import matplotlib.pyplot as plt
4   from statsmodels.graphics.gofplots import qqplot
5   import seaborn as sns
6   from scipy.cluster.hierarchy import linkage
7   from scipy.cluster.hierarchy import dendrogram
8   from scipy.cluster.hierarchy import cut_tree
9   import warnings
10  from sklearn.preprocessing import LabelEncoder
11
12  warnings.filterwarnings('ignore')
13
14  dataset=pd.read_csv('heart.csv')
15
16
17
18
19  #to covert categorical data into numeric
20  columns=['Sex','ChestPainType','RestingECG','ExerciseAngina','ST_Slope']
21  encoders = {}
22  for col in columns:
23      le = LabelEncoder().fit(dataset[col])
24      dataset[col] = le.transform(dataset[col])
25      encoders[col] = le
26
27
28
29
30  #HIERARICHAL CLUSTERING
31  mergings = linkage(dataset, method="single", metric='euclidean')
32  dendrogram(mergings)
```

```
33  plt.show()
34
35  mergings = linkage(dataset, method="complete", metric='euclidean')
36  dendrogram(mergings)
37  plt.show()
38
39  mergings = linkage(dataset, method="average", metric='euclidean')
40  dendrogram(mergings)
41  plt.show()
42
43
44  cluster_labels = cut_tree(mergings, n_clusters=2).reshape(-1, )
45  cluster_labels
46
47
48  dataset['Cluster_Labels'] = cluster_labels
49  dataset.head()
50
51  sns.boxplot(x='Cluster_Labels', y='Age', data=dataset)
52  plt.show()
53  sns.boxplot(x='Cluster_Labels', y='Sex', data=dataset)
54  plt.show()
55  sns.boxplot(x='Cluster_Labels', y='ChestPainType', data=dataset)
56  plt.show()
57  sns.boxplot(x='Cluster_Labels', y='RestingBP', data=dataset)
58  plt.show()
59  sns.boxplot(x='Cluster_Labels', y='FastingBS', data=dataset)
60  plt.show()
61  sns.boxplot(x='Cluster_Labels', y='RestingECG', data=dataset)
62  plt.show()
63  sns.boxplot(x='Cluster_Labels', y='MaxHR', data=dataset)
64  plt.show()
```

```
65    sns.boxplot(x='Cluster_Labels', y='ExerciseAngina', data=dataset)
66    plt.show()
67    sns.boxplot(x='Cluster_Labels', y='Oldpeak', data=dataset)
68    plt.show()
69    sns.boxplot(x='Cluster_Labels', y='ST_Slope', data=dataset)
70    plt.show()
71
```

# FUZZY C MEAN CLUSTERING

Fuzzy c mean is a type of clustering method in which a data point can belong to two or more clusters by initializing membership matrix then calculating cluster centers based on initial membership values. Update the membership value to attain the accuracy of clusters.

We implemented Fuzyy c mean on our dataset related heart failure prediction to cluster attributes. We used only few attributes to show the results. We performed clustering on Age and Max Heart Rate, and on Resting Blood Pressure and Cholesterol. We initialized the clusters at the origin, at random locations within a multi-variate Gaussian distribution with zero-mean and unit-variance and at random vectors chosen from the data. Scatter plot was used to plot the clusters.

```
1   import numpy as np
2   import pandas as pd
3   import matplotlib.pyplot as plt
4   from statsmodels.graphics.gofplots import qqplot
5   import seaborn as sns
6   import warnings
7   from sklearn.preprocessing import LabelEncoder
8   import random
9   import operator
10  import math
11  from scipy.stats import multivariate_normal     # for generating pdf
12
13  warnings.filterwarnings('ignore')
14
15  df=pd.read_csv('heart.csv')
16
17
18
19
20  #to covert categorical data into numeric
21  columns=['Sex','ChestPainType','RestingECG','ExerciseAngina','ST_Slope']
22  encoders = {}
23  for col in columns:
24      le = LabelEncoder().fit(df[col])
25      df[col] = le.transform(df[col])
26      encoders[col] = le
27
28
29  #FUZZY C MEAN
30  # Number of Clusters
31  k = 2
32  # Maximum number of iterations
33  MAX_ITER = 100
34  # Number of data points
35  n = len(df)
36  # Fuzzy parameter
37  m = 1.7 #Select a value greater than 1 else it will be knn
38
39
40  def initializeMembershipMatrix(): # initializing the membership matrix
41      membership_mat = []
42      for i in range(n):
43          random_num_list = [random.random() for i in range(k)]
44          summation = sum(random_num_list)
45          temp_list = [x/summation for x in random_num_list]
46
47          flag = temp_list.index(max(temp_list))
48          for j in range(0,len(temp_list)):
49              if(j == flag):
50                  temp_list[j] = 1
51              else:
52                  temp_list[j] = 0
53
54          membership_mat.append(temp_list)
55      return membership_mat
56  membership_mat = initializeMembershipMatrix()
57
58
59
60  def calculateClusterCenter(membership_mat): # calculating the cluster center
61      cluster_mem_val = list(zip(*membership_mat))
62      cluster_centers = []
63      for j in range(k):
64          x = list(cluster_mem_val[j])
```

```python
65              xraised = [p ** m for p in x]
66              denominator = sum(xraised)
67              temp_num = []
68              for i in range(n):
69                  data_point = list(df.iloc[i])
70                  prod = [xraised[i] * val for val in data_point]
71                  temp_num.append(prod)
72              numerator = map(sum, list(zip(*temp_num)))
73              center = [z/denominator for z in numerator]
74              cluster_centers.append(center)
75          return cluster_centers
76      #cluster_centers = calculateClusterCenter(membership_mat)
77      calculateClusterCenter(membership_mat)
78
79
80      def updateMembershipValue(membership_mat, cluster_centers): # Updating the membership value
81          p = float(2/(m-1))
82          for i in range(n):
83              x = list(df.iloc[i])
84              distances = [np.linalg.norm(np.array(list(map(operator.sub, x, cluster_centers[j])))) for j in range(k)]
85              for j in range(k):
86                  den = sum([math.pow(float(distances[j]/distances[c]), p) for c in range(k)])
87                  membership_mat[i][j] = float(1/den)
88          return membership_mat
89
90
91      def getClusters(membership_mat): # getting the clusters
92          cluster_labels = list()
93          for i in range(n):
94              max_val, idx = max((val, idx) for (idx, val) in enumerate(membership_mat[i]))
95              cluster_labels.append(idx)
96          return cluster_labels
97
98
99      def fuzzyCMeansClustering(): #First Iteration with centers at 0
100         # Membership Matrix
101         membership_mat = initializeMembershipMatrix()
102         curr = 0
103         acc=[]
104         cent_temp = [[0, 0, 0, 0],[0, 0, 0, 0],[0, 0, 0, 0]]
105         while curr < MAX_ITER:
106             if(curr == 0):
107                 cluster_centers = cent_temp
108                 print("Cluster Centers:")
109                 print(np.array(cluster_centers))
110             else:
111                 cluster_centers = calculateClusterCenter(membership_mat)
112             #cluster_centers = calculateClusterCenter(membership_mat)
113             membership_mat = updateMembershipValue(membership_mat, cluster_centers)
114             cluster_labels = getClusters(membership_mat)
115             acc.append(cluster_labels)
116             curr += 1
117         print("--------------------------")
118         print("Membership Matrix:")
119         print(np.array(membership_mat))
120         return cluster_labels, cluster_centers, acc
121
122
123
124     def fuzzyCMeansClustering(): #Second Iteration Multivariate Gaussian
125         # Membership Matrix
126         membership_mat = initializeMembershipMatrix()
127         curr = 0
128         acc=[]
```

```python
129        mean = [0, 0]
130        cov = [[1, 0], [0, 1]]
131
132        lis1,cent_temp=[],[]
133
134        for i in range(0,k):
135            Z = list(np.random.multivariate_normal(mean, cov))
136            Z1 = list(np.random.multivariate_normal(mean, cov))
137            lis1 = Z+Z1
138            cent_temp.append(lis1)
139
140
141        while curr < MAX_ITER:
142            if(curr == 0):
143                cluster_centers = cent_temp
144                print("Cluster Centers:")
145                print(np.array(cluster_centers))
146            else:
147                cluster_centers = calculateClusterCenter(membership_mat)
148            #cluster_centers = calculateClusterCenter(membership_mat)
149            membership_mat = updateMembershipValue(membership_mat, cluster_centers)
150            cluster_labels = getClusters(membership_mat)
151            acc.append(cluster_labels)
152            curr += 1
153        print("---------------------------")
154        print("Membership Matrix:")
155        print(np.array(membership_mat))
156        return cluster_labels, cluster_centers, acc
157
158
159
160    def fuzzyCMeansClustering(): #Third iteration Random vectors from data
```

```python
161        # Membership Matrix
162        membership_mat = initializeMembershipMatrix()
163        curr = 0
164        acc=[]
165        while curr < MAX_ITER:
166            cluster_centers = calculateClusterCenter(membership_mat)
167            membership_mat = updateMembershipValue(membership_mat, cluster_centers)
168            cluster_labels = getClusters(membership_mat)
169
170            acc.append(cluster_labels)
171
172            if(curr == 0):
173                print("Cluster Centers:")
174                print(np.array(cluster_centers))
175            curr += 1
176        print("---------------------------")
177        print("Partition matrix:")
178        print(np.array(membership_mat))
179        #return cluster_labels, cluster_centers
180        return cluster_labels, cluster_centers, acc
181
182    labels, centers, acc = fuzzyCMeansClustering()
183    print("Cluster center vectors:") #final cluster centers
184    print(np.array(centers))
185
186
187    #FOR COLUMNS AGE AND MAX HEART RATE
188    df1 = df.iloc[:,[0,7]]
189    df1 = np.array(df1)
190    #First initialization
191    #m1 = [0,0]
192    #m2 = [0,0]
```

```python
193    #Second initialization
194    m1 = random.choice(df1)
195    m2 = random.choice(df1)
196
197    cov1 = np.cov(np.transpose(df1))
198    cov2 = np.cov(np.transpose(df1))
199    x1 = np.linspace(27.5,77.5,918)
200    x2 = np.linspace(59.5,202.5,918)
201    X, Y = np.meshgrid(x1,x2)
202
203    Z1 = multivariate_normal(m1, cov1)
204    Z2 = multivariate_normal(m2, cov2)
205
206    pos = np.empty(X.shape + (2,))                        # a new array of given shape and type, without initializing entrie
207    pos[:, :, 0] = X; pos[:, :, 1] = Y
208
209    plt.figure(figsize=(10,10))                                              # creating the figure and a
210    plt.scatter(df1[:,0], df1[:,1], marker='o')
211    plt.contour(X, Y, Z1.pdf(pos), colors="r" ,alpha = 0.5)
212    plt.contour(X, Y, Z2.pdf(pos), colors="b" ,alpha = 0.5)
213    plt.axis('equal')                                          # making both the axis equal
214    plt.xlabel('Age', fontsize=16)                          # X-Axis
215    plt.ylabel('MaxHR', fontsize=16)                        # Y-Axis
216    plt.title('Initial Random Clusters', fontsize=22)             # Title of the plot
217    plt.grid()                                             # displaying gridlines
218    plt.show()
219
220    #finding mode
221    seto = max(set(labels[0:459]), key=labels[0:459].count)
222    virg = max(set(labels[459:]), key=labels[459:].count)
223    seto=0
224    virg=1
```

```python
225    s_mean_clus1 = np.array([centers[seto][0],centers[seto][7]])
226    s_mean_clus2 = np.array([centers[virg][0],centers[virg][7]])
227    values = np.array(labels) #label
228
229    #search all 3 species
230    searchval_seto = seto
231    searchval_virg = virg
232
233    #index of all 3 species
234    ii_seto = np.where(values == searchval_seto)[0]
235    ii_virg = np.where(values == searchval_virg)[0]
236    ind_seto = list(ii_seto)
237    ind_virg = list(ii_virg)
238    df1 = df.iloc[:,[0,7]]
239    seto_df = df1[df1.index.isin(ind_seto)]
240    virg_df = df1[df1.index.isin(ind_virg)]
241    cov_seto = np.cov(np.transpose(np.array(seto_df)))
242    cov_virg = np.cov(np.transpose(np.array(virg_df)))
243    df1 = np.array(df1)
244    x1 = np.linspace(27.5,77.5,918)
245    x2 = np.linspace(59.5,202.5,918)
246    X, Y = np.meshgrid(x1,x2)
247
248    Z1 = multivariate_normal(s_mean_clus1, cov_seto)
249    Z2 = multivariate_normal(s_mean_clus2, cov_virg)
250
251    pos = np.empty(X.shape + (2,))                        # a new array of given shape and type, without initializing entri
252    pos[:, :, 0] = X; pos[:, :, 1] = Y
253
254    plt.figure(figsize=(10,10))                                              # creating the figure and
255    plt.scatter(df1[:,0], df1[:,1], marker='o')
256    plt.contour(X, Y, Z1.pdf(pos), colors="r" ,alpha = 0.5)
```

```python
257    plt.contour(X, Y, Z2.pdf(pos), colors="g" ,alpha = 0.5)
258    plt.axis('equal')                                                   # making both the axis equal
259    plt.xlabel('Age', fontsize=16)                                   # X-Axis
260    plt.ylabel('MaxHR', fontsize=16)                                 # Y-Axis
261    plt.title('Final Clusters', fontsize=22)                           # Title of the plot
262    plt.grid()                                                         # displaying gridlines
263    plt.show()
264
265
266    #FOR RESTING BLOOD PRESSURE AND CHOLESTEROL
267    df1 = df.iloc[:,3:5]
268    df1 = np.array(df1)
269    #first initialization
270    #m1 = [0,0]
271    #m2 = [0,0]
272    #Second initialization
273    m1 = random.choice(df1)
274    m2 = random.choice(df1)
275    cov1 = np.cov(np.transpose(df1))
276    cov2 = np.cov(np.transpose(df1))
277    x1 = np.linspace(-0.5,200.5,918)
278    x2 = np.linspace(-0.5,603.5,918)
279    X, Y = np.meshgrid(x1,x2)
280
281    Z1 = multivariate_normal(m1, cov1)
282    Z2 = multivariate_normal(m2, cov2)
283
284    pos = np.empty(X.shape + (2,))                    # a new array of given shape and type, without initializing entrie
285    pos[:, :, 0] = X; pos[:, :, 1] = Y
286
287    plt.figure(figsize=(10,10))                                        # creating the figure and a
288    plt.scatter(df1[:,0], df1[:,1], marker='o')
```

```python
289    plt.contour(X, Y, Z1.pdf(pos), colors="r" ,alpha = 0.5)
290    plt.contour(X, Y, Z2.pdf(pos), colors="b" ,alpha = 0.5)
291    plt.axis('equal')                                                   # making both the axis equal
292    plt.xlabel('Resting Blood Pressure', fontsize=16)                       # X-Axis
293    plt.ylabel('Cholesterol', fontsize=16)                           # Y-Axis
294    plt.title('Initial Random Clusters', fontsize=22)                  # Title of the plot
295    plt.grid()                                                       # displaying gridlines
296    plt.show()
297
298    #finding mode
299    seto = max(set(labels[0:459]), key=labels[0:459].count)
300    virg = max(set(labels[459:]), key=labels[459:].count)
301    seto=0
302    virg=1
303    s_mean_clus1 = np.array([centers[seto][3],centers[seto][4]])
304    s_mean_clus2 = np.array([centers[virg][3],centers[virg][4]])
305    values = np.array(labels) #label
306
307    #search all 3 species
308    searchval_seto = seto
309    searchval_virg = virg
310
311    #index of all 3 species
312    ii_seto = np.where(values == searchval_seto)[0]
313    ii_virg = np.where(values == searchval_virg)[0]
314    ind_seto = list(ii_seto)
315    ind_virg = list(ii_virg)
316    df1 = df.iloc[:,3:5]
317    seto_df = df1[df1.index.isin(ind_seto)]
318    virg_df = df1[df1.index.isin(ind_virg)]
319    cov_seto = np.cov(np.transpose(np.array(seto_df)))
320    cov_virg = np.cov(np.transpose(np.array(virg_df)))
```

```
321    df1 = np.array(df1)
322    x1 = np.linspace(-0.5,200.5,918)
323    x2 = np.linspace(-0.5,603.5,918)
324    X, Y = np.meshgrid(x1,x2)
325
326    Z1 = multivariate_normal(s_mean_clus1, cov_seto)
327    Z2 = multivariate_normal(s_mean_clus2, cov_virg)
328
329    pos = np.empty(X.shape + (2,))                      # a new array of given shape and type, without initializing entri
330    pos[:, :, 0] = X; pos[:, :, 1] = Y
331
332    plt.figure(figsize=(10,10))                                                      # creating the figure and
333    plt.scatter(df1[:,0], df1[:,1], marker='o')
334    plt.contour(X, Y, Z1.pdf(pos), colors="r" ,alpha = 0.5)
335    plt.contour(X, Y, Z2.pdf(pos), colors="g" ,alpha = 0.5)
336    plt.axis('equal')                                                        # making both the axis equal
337    plt.xlabel('Resting Blood Pressure', fontsize=16)                             # X-Axis
338    plt.ylabel('Cholesterol', fontsize=16)                                        # Y-Axis
339    plt.title('Final Clusters', fontsize=22)                          # Title of the plot
340    plt.grid()                                                        # displaying gridlines
341    plt.show()
342
```