



**L**OVELY  
**P**ROFESSIONAL  
**U**NIVERSITY

# **MAJOR PROJECT REPORT**

**On dataset: Wholesale Customers Data Set**

**Name: Krishna Mohan Mishra**

**Reg no: 11605503**

**Roll no: A10**

**Submitted to: Mr. GokulNath K**

# OBJECTIVE:

1. Preparing data (Feature Identification) .
2. Choosing a model (Comparative study of minimum 3 models).
3. Training (Reiterating for optimization) .
4. Evaluation (Highest accuracy with confusion matrix)
5. Hyperparameter tuning.

## DataSet:

### Wholesale Customers Data Set

<http://archive.ics.uci.edu/ml/machine-learning-databases/00292/Wholesale%20customers%20data.csv>

#### Data Set Information:

Wholesale Customer dataset refers to clients of a wholesale distributor. It includes the annual spending in monetary units (m.u.) on diverse product categories. Number of Instances: 440, Number of Attributes: 8, Data Set Characteristics: Multivariate, Attribute Characteristics: Integer

Data Set Characteristics:	Multivariate	Number of Instances:	440	Area:	Business
Attribute Characteristics:	Integer	Number of Attributes:	8	Date Donated	2014-03-31
Associated Tasks:	Classification, Clustering	Missing Values?	N/A	Number of Web Hits:	270339

Channel	Region	Fresh	Milk	Grocery	Frozen	Detergents	Delicassen
2	3	12669	9656	7561	214	2674	1338
2	3	7057	9810	9568	1762	3293	1776
2	3	6353	8808	7684	2405	3516	7844
1	3	13265	1196	4221	6404	507	1788
2	3	22615	5410	7198	3915	1777	5185
2	3	9413	8259	5126	666	1795	1451
2	3	12126	3199	6975	480	3140	545
2	3	7579	4956	9426	1669	3321	2566
1	3	5963	3648	6192	425	1716	750
2	3	6006	11093	18881	1159	7425	2098
2	3	3366	5403	12974	4400	5977	1744
2	3	13146	1124	4523	1420	549	497
2	3	31714	12319	11757	287	3881	2931
2	3	21217	6208	14982	3095	6707	602

### Attribute Information:

1. FRESH: annual spending (m.u.) on fresh products (Continuous)
2. MILK: annual spending (m.u.) on milk products (Continuous)
3. GROCERY: annual spending (m.u.) on grocery products (Continuous)
4. FROZEN: annual spending (m.u.) on frozen products (Continuous)
5. DETERGENTS\_PAPER: annual spending (m.u.) on detergents and paper products (Continuous)
6. DELICATESSEN: annual spending (m.u.) on and delicatessen products (Continuous)
7. CHANNEL: customersale Channel - Horeca (Hotel/Restaurant/Cafe) or Retail channel (Nominal)
8. REGION: customersale Region - Lisbon, Oporto or Other (Nominal)

# Objective 1: Preparing data

Machine learning helps us find patterns in data—patterns we then use to make predictions about new data points. To get those predictions right, we must

*construct* the [data set](#) and *transform* the data correctly.

The quality of the data and the amount of useful information that it contains are key factors that determine how well a machine learning algorithm can learn. Therefore, it is absolutely critical that we make sure to examine and preprocess a dataset before we feed it to a learning algorithm.

#import dataset

In this we import libraries and dataset and check if the dataset missing?

```
8 import pandas as pd
9 import numpy as np
10 # Load the wholesale customers dataset
11 #check whether any dataset missing ?
12 try:
13     data = pd.read_csv("customers.csv")
14     data.drop(['Region', 'Channel'], axis = 1, inplace = True)
15     print( "Wholesale customers dataset has {} samples with {} features each.".format(*data.shape))
16 except:
17     print ("Dataset could not be loaded. Is the dataset missing?")
18
19
20 stats = data.describe()
21 print(stats)
22 # Using data.loc to filter a pandas DataFrame
23 data.loc[[100, 200, 300],:]
24 data.columns
25 # Select three indices of your choice you wish to sample from the dataset
26 indices = [43, 12, 39]
27
28 # Create a DataFrame of the chosen samples
29 # .reset_index(drop = True) resets the index from 0, 1 and 2 instead of 100, 200 and 300
30 samples = pd.DataFrame(data.loc[indices], columns = data.columns).reset_index(drop = True)
31 print ("Chosen samples of wholesale customers dataset:")
```

Wholesale customers dataset has 440 samples with 6 features each.

	Fresh	Milk	Grocery	Frozen \
count	440.000000	440.000000	440.000000	440.000000
mean	12000.297727	5796.265909	7951.277273	3071.931818
std	12647.328865	7380.377175	9503.162829	4854.673333
min	3.000000	55.000000	3.000000	25.000000
25%	3127.750000	1533.000000	2153.000000	742.250000
50%	8504.000000	3627.000000	4755.500000	1526.000000
75%	16933.750000	7190.250000	10655.750000	3554.250000
max	112151.000000	73498.000000	92780.000000	60869.000000

	Detergents_Paper	Delicassen
count	440.000000	440.000000
mean	2881.493182	1524.870455
std	4767.854448	2820.105937
min	3.000000	3.000000
25%	256.750000	408.250000
50%	816.500000	965.500000
75%	3922.000000	1820.250000
max	40827.000000	47943.000000

## #train\_test\_split:

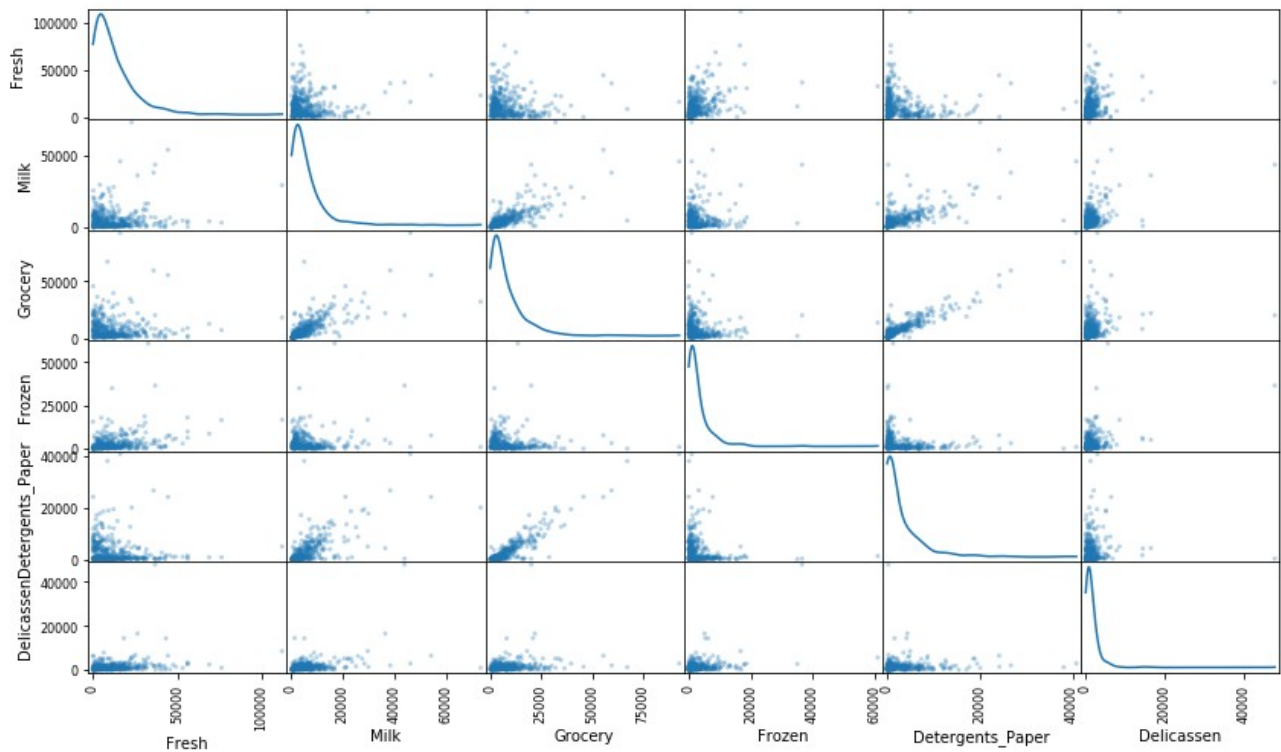
To evaluate how well a trained model performs on unseen data, we have to split the dataset into separate training and test datasets.

```
# Split the data into training and testing sets using the given feature as the target
X_train, X_test, y_train, y_test = train_test_split(new_data, new_feature, test_size=0.25, random_state=42)
```

## Scatter matrix

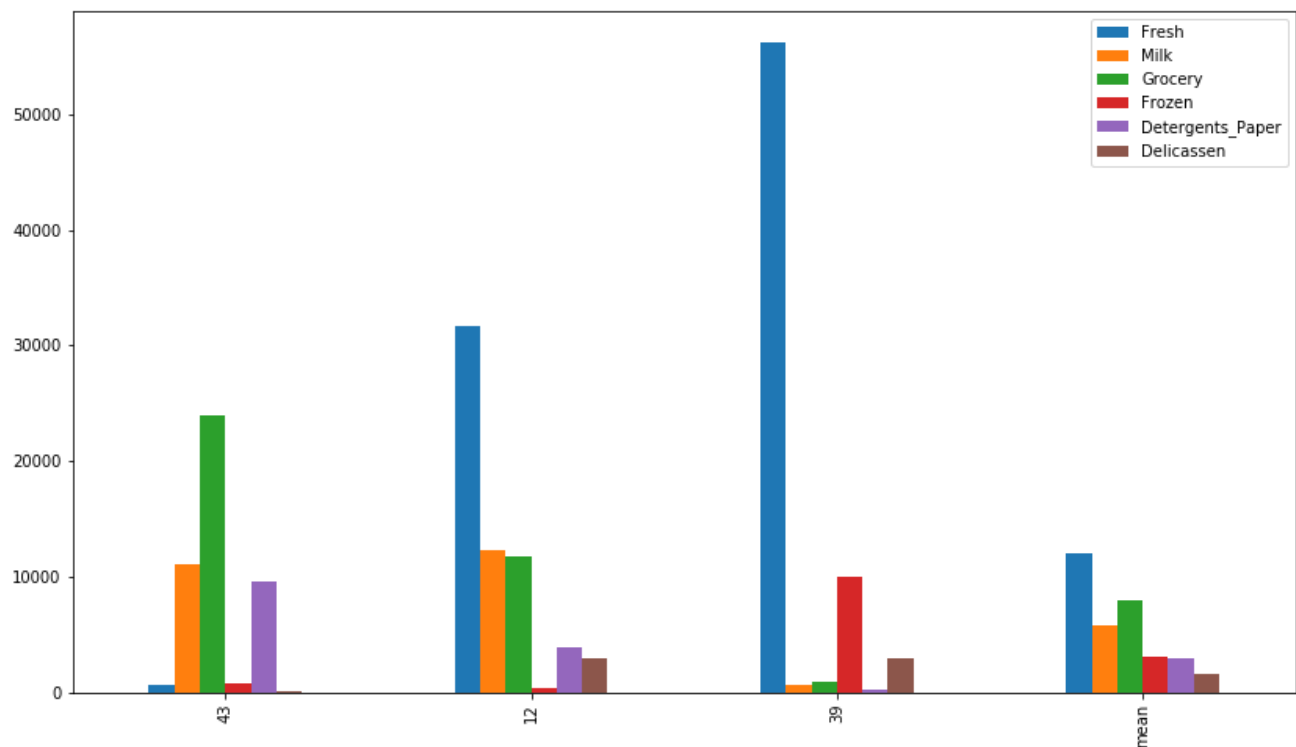
A scatter matrix is a pair-wise scatter plot of several variables presented in a matrix format. It can be used to determine whether the variables are correlated and whether the correlation is positive or negative.

```
70  
71 pd.scatter_matrix(data, alpha = 0.3, figsize = (14,8), diagonal = 'kde'); #scatter matrix
```



## Bar Plot

```
32 # Get the means
33 mean_data = data.describe().loc['mean', :]
34 # Append means to the samples' data
35 samples_bar = samples.append(mean_data)
36 # Construct indices
37 samples_bar.index = indices + ['mean']
38 # Plot bar plot
39 samples_bar.plot(kind='bar', figsize=(14,8))
40 # First, calculate the percentile ranks of the whole dataset.
41 percentiles = data.rank(pct=True)
42 # Then, round it up, and multiply by 100
43 percentiles = 100*percentiles.round(decimals=3)
44 # Select the indices you chose from the percentiles dataframe
45 percentiles = percentiles.iloc[indices]
46 data.columns
47
```



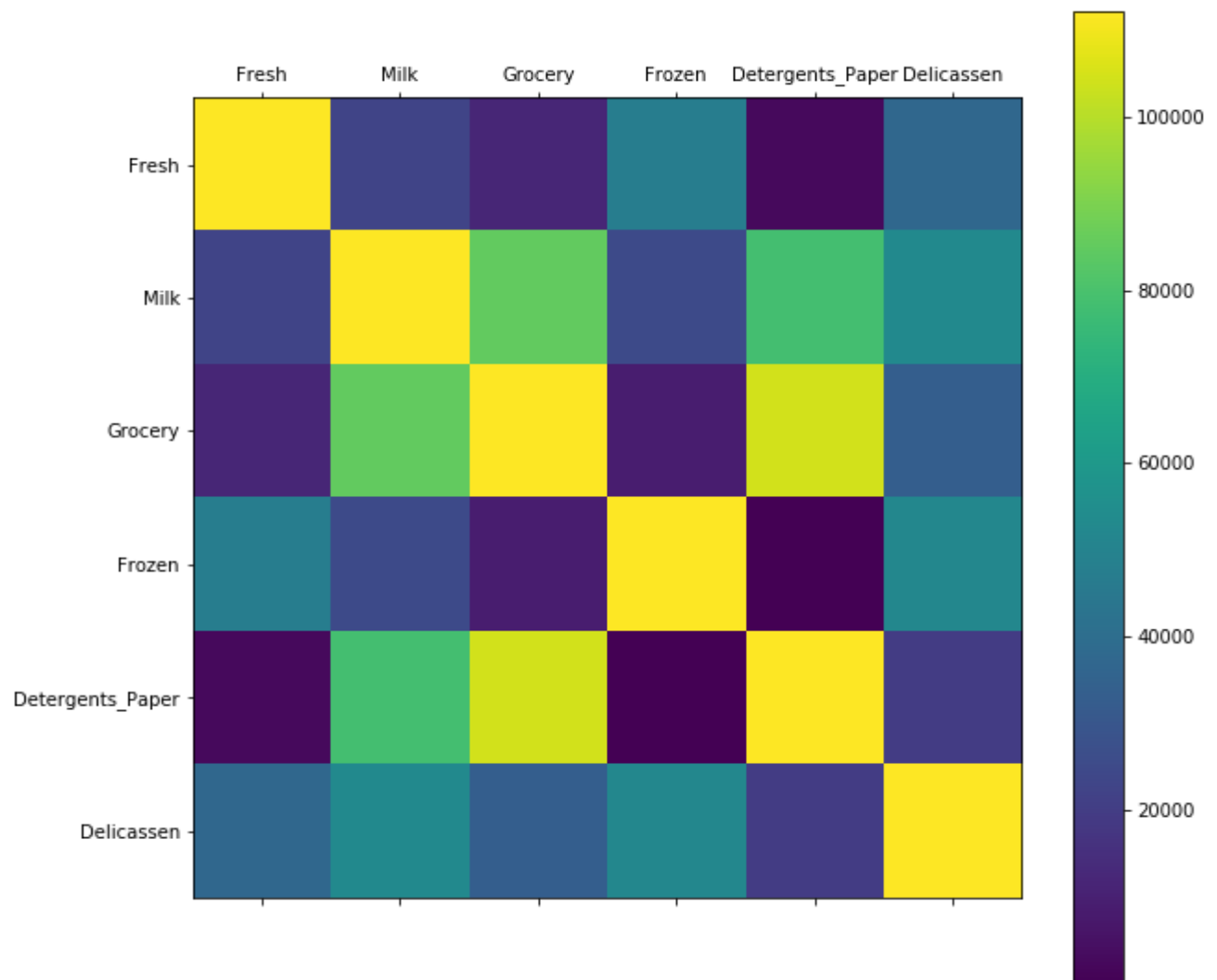
## Correlation Matrix

A correlation matrix is a table showing correlation coefficients between variables. Each cell in the table shows the correlation between two variables. A correlation matrix is used as a way to summarize data, as an input into a more advanced analysis, and as a diagnostic for advanced analyses.

```

73
74 def plot_corr(df, size=10):
75     corr = df.corr()
76     fig, ax = plt.subplots(figsize=(size, size))
77     cax = ax.matshow(df, interpolation='nearest')
78     ax.matshow(corr)
79     fig.colorbar(cax)
80     plt.xticks(range(len(corr.columns)), corr.columns);
81     plt.yticks(range(len(corr.columns)), corr.columns);
82
83 plot_corr(data)

```





## Dimensionality reduction using pca

Principal component analysis (PCA) is an unsupervised linear transformation technique that is widely used across different fields, most prominently for dimensionality reduction.

```
173 from sklearn.decomposition import PCA
174
175 #Apply PCA by fitting the good data with the same number of dimensions as features
176 pca = PCA(n_components=6)
177 # Fit
178 pca.fit(good_data)
179
180 # Transform the sample log-data using the PCA fit above
181 pca_samples = pca.transform(log_samples)
182
183 # Generate PCA results
184 #Apply PCA by fitting the good data with only two dimensions
185
186 pca = PCA(n_components=2)
187 pca.fit(good_data)
188 #Transform the good data using the PCA fit above
189 reduced_data = pca.transform(good_data)
190 #Transform the sample log-data using the PCA fit above
191 pca_samples = pca.transform(log_samples)
192 # Create a DataFrame for the reduced data
193 reduced_data = pd.DataFrame(reduced_data, columns = ['Dimension 1', 'Dimension 2'])
194
195 # Transform the sample log-data using the PCA fit above
```

## Objective 2: Choosing a model

I have taken three models and using these three models we find out the accuracy.

### Support vector machine (SVM):

Support Vector Machine (**SVM**) is a supervised binary classification algorithm. Given a set of points of two types in N dimensional place SVM generates a (N-1) dimensional hyperplane to separate those points into two groups.

```
245 for var in dep_vars:
246
247     #Make a copy of the DataFrame, using the 'drop' function to drop the given feature
248     new_data = data.drop([var], axis = 1)
249
250     # Create feature Series (Vector)
251     new_feature = pd.DataFrame(data.loc[:, var])
252
253     # Split the data into training and testing sets using the given feature as the target
254     X_train, X_test, y_train, y_test = train_test_split(new_data, new_feature, test_size=0.25, random_state=42)
255     dtr = SVC(kernel='linear', C=1, random_state = 42)
256     dtr.fit(X_train, y_train)
257     y_pred = dtr.predict(X_test)
258     from sklearn.metrics import confusion_matrix
259     cm = confusion_matrix(y_test, y_pred)
260     print("confusion matrix of SVM\n", cm)
261
```

### DecisionTreeRegressor:

```
48 # Importing the libraries
49 from sklearn.model_selection import train_test_split
50 from sklearn.tree import DecisionTreeRegressor
51
52 # Create list to loop through
53 dep_vars = list(data.columns)
54
55
56 # Create loop to test each feature as a dependent variable
57 for var in dep_vars:
58     #Make a copy of the DataFrame, using the 'drop' function to drop the given feature
59     new_data = data.drop([var], axis = 1)
60     new_feature = pd.DataFrame(data.loc[:, var]) # feature series
61     # Split the data into training and testing sets using the given feature as the target
62     X_train, X_test, y_train, y_test = train_test_split(new_data, new_feature, test_size=0.25, random_state=42)
63     # Creating a decision tree regressor and fitting it to the training set
64     dtr = DecisionTreeRegressor(random_state=42)
65     dtr.fit(X_train, y_train) # fitting the data
66     y_pred = dtr.predict(X_test)
67     from sklearn.metrics import confusion_matrix
68     cm = confusion_matrix(y_test, y_pred)
69     print("confusion matrix of decision regression\n", cm)
70
```

## KNeighborsClassifier:

```
195 # Importing libraries
196
197 from sklearn.cluster import KMeans
198 from sklearn.metrics import silhouette_score
199 # Create range of clusters
200 range_n_clusters = list(range(2,11))
201 print(range_n_clusters)
202
203 #KNN Implementation
204 for n_clusters in range_n_clusters:
205     # Apply your clustering algorithm of choice to the reduced data
206     clusterer = KMeans(n_clusters=n_clusters).fit(reduced_data)
207
208     # Predict the cluster for each data point
209     preds = clusterer.predict(reduced_data)
210
211     #Find the cluster centers
212     centers = clusterer.cluster_centers_
213
214     #Predict the cluster for each transformed sample data point
215     sample_preds = clusterer.predict(pca_samples)
216
217     #Calculate the mean silhouette coefficient for the number of clusters chosen
218     score = silhouette_score(reduced_data, preds, metric='euclidean')
219
220     print ("For n_clusters = {}. The average silhouette_score is : {}".format(n_clusters, score))
221
222
```

## Objective 3: Training

In this objective given to us we have to reiterate the models to get an optimum or optimized model.

```

48 # Importing the libraries
49 from sklearn.model_selection import train_test_split
50 from sklearn.tree import DecisionTreeRegressor
51
52 # Create list to loop through
53 dep_vars = list(data.columns)
54
55
56 # Create loop to test each feature as a dependent variable
57 for var in dep_vars:
58     #Make a copy of the DataFrame, using the 'drop' function to drop the given feature
59     new_data = data.drop([var], axis = 1)
60     new_feature = pd.DataFrame(data.loc[:, var]) # feature series
61     # Split the data into training and testing sets using the given feature as the target
62     X_train, X_test, y_train, y_test = train_test_split(new_data, new_feature, test_size=0.25, random_state=42)
63     # Creating a decision tree regressor and fitting it to the training set
64     dtr = DecisionTreeRegressor(random_state=42)
65     dtr.fit(X_train, y_train) # fitting the data
66     y_pred = dtr.predict(X_test)
67     from sklearn.metrics import confusion_matrix
68     cm = confusion_matrix(y_test, y_pred)
69     print("confusion matrix of decision regression\n",cm)
70

```

```

213 #KNN Implementation
214 for n_clusters in range_n_clusters:
215     # Apply your clustering algorithm of choice to the reduced data
216     clusterer = KMeans(n_clusters=n_clusters).fit(reduced_data)
217
218     # Predict the cluster for each data point
219     preds = clusterer.predict(reduced_data)
220
221     #Find the cluster centers
222     centers = clusterer.cluster_centers_
223
224     #Predict the cluster for each transformed sample data point
225     sample_preds = clusterer.predict(pca_samples)
226
227     #Calculate the mean silhouette coefficient for the number of clusters chosen
228     score = silhouette_score(reduced_data, preds, metric='euclidean')
229
230     print ("For n_clusters = {}. The average silhouette_score is : {}".format(n_clusters, score))
231

```

```

245 for var in dep_vars:
246
247     #Make a copy of the DataFrame, using the 'drop' function to drop the given feature
248     new_data = data.drop([var], axis = 1)
249
250     # Create feature Series (Vector)
251     new_feature = pd.DataFrame(data.loc[:, var])
252
253     # Split the data into training and testing sets using the given feature as the target
254     X_train, X_test, y_train, y_test = train_test_split(new_data, new_feature, test_size=0.25, random_state=42)
255     dtr = SVC(kernel='linear', C=1, random_state = 42)
256     dtr.fit(X_train, y_train)
257     y_pred = dtr.predict(X_test)
258     from sklearn.metrics import confusion_matrix
259     cm = confusion_matrix(y_test, y_pred)
260     print("confusion matrix of SVM\n",cm)
261

```

## Objective 4: Evaluation

Confusion matrix: a matrix that lays out the performance of a learning algorithm. The confusion matrix is simply a square matrix that reports the counts of the true positive, true negative, false positive, and false negative predictions of a classifier.

Silhouette Score – It is measure of how similar an object to its own cluster compared to other cluster.

### Confusion Matrix For SVM:

```
confusion matrix of SVM
[[0 0 1 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
confusion matrix of SVM
[[0 1 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
confusion matrix of SVM
[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 1 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
confusion matrix of SVM
[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
```

## Confusion Matrix For Decision Tree Regression:

```
confusion matrix of decision regression
[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
confusion matrix of decision regression
[[0 1 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
confusion matrix of decision regression
[[0 0 0 ... 0 0 0]
 [0 0 1 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 1 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
confusion matrix of decision regression
[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]]
```

## Silhouette Score for KneighborsClassifier:

```
For n_clusters = 2. The average silhouette_score is : 0.42628101546910835
For n_clusters = 3. The average silhouette_score is : 0.39418735404519234
For n_clusters = 4. The average silhouette_score is : 0.3329304196814193
For n_clusters = 5. The average silhouette_score is : 0.3496074606159664
For n_clusters = 6. The average silhouette_score is : 0.3650885924121462
For n_clusters = 7. The average silhouette_score is : 0.36490353248489016
For n_clusters = 8. The average silhouette_score is : 0.36405497488401767
For n_clusters = 9. The average silhouette_score is : 0.3572844712605567
For n_clusters = 10. The average silhouette_score is : 0.3542011421758249
```

## Objective 5: Hyperparameter tuning

The approach of grid search is quite simple, it's a brute-force exhaustive search paradigm where we specify a list of values for different hyperparameters, and the computer evaluates the model performance for each combination of those to obtain the optimal set.

Parameters which define the model architecture are referred to as hyperparameters and thus this process of searching for the ideal model architecture is referred to as hyperparameter tuning.

```
212 from sklearn.svm import SVC
213 from sklearn.model_selection import GridSearchCV
214 from sklearn.pipeline import Pipeline
215 from sklearn.preprocessing import StandardScaler
216 X_train, X_test, y_train, y_test = train_test_split(new_data, new_feature, test_size=0.25, random_state=42)
217 pipe_svc = Pipeline([('scl', StandardScaler()), ('clf', SVC(random_state=1))])
218 param_range = [0.0001, 0.001, 0.01, 0.1, 1.0, 10.0, 100.0, 1000.0]
219 param_grid = [{'clf_C': param_range, 'clf_kernel': ['linear']}, {'clf_C': param_range, 'clf_gamma': param_range, 'clf_kernel': ['rbf']}]
220 gs=GridSearchCV(estimator=pipe_svc,param_grid=param_grid,scoring='accuracy',cv=3,n_jobs=-1)
221 gs=gs.fit(X_train,y_train)
222 print(gs.best_score_)
223 print(gs.best_params_)
224
```

