به نام خدا



**دانشگاه صنعتی اصفهان**

دانشکده مهندسی برق و کامپیوتر

گروه مهندسی نرم افزار

Mini Project   2

**استاد**:

دکتر ناصر قدیری مدرس

**نویسنده**:

علی کتیرایی

تیرماه ۹۸

متد DESCR

```
1 import sklearn.datasets as datasets
2 data=datasets.load_breast_cancer()
3 print(data.DESCR)
```

```
.. _breast_cancer_dataset:

Breast cancer wisconsin (diagnostic) dataset
--------------------------------------------

**Data Set Characteristics:**

    :Number of Instances: 569

    :Number of Attributes: 30 numeric, predictive attributes and the class

    :Attribute Information:
        - radius (mean of distances from center to points on the perimeter)
        - texture (standard deviation of gray-scale values)
        - perimeter
        - area
        - smoothness (local variation in radius lengths)
        - compactness (perimeter^2 / area - 1.0)
        - concavity (severity of concave portions of the contour)
        - concave points (number of concave portions of the contour)
        - symmetry
        - fractal dimension ("coastline approximation" - 1)

        The mean, standard error, and "worst" or largest (mean of the three
        largest values) of these features were computed for each image,
        resulting in 30 features.  For instance, field 3 is Mean Radius, field
        13 is Radius SE, field 23 is Worst Radius.

        - class:
                - WDBC-Malignant
                - WDBC-Benign
```

:Summary Statistics:

| | Min | Max |
|---|---|---|
| radius (mean): | 6.981 | 28.11 |
| texture (mean): | 9.71 | 39.28 |
| perimeter (mean): | 43.79 | 188.5 |
| area (mean): | 143.5 | 2501.0 |
| smoothness (mean): | 0.053 | 0.163 |
| compactness (mean): | 0.019 | 0.345 |
| concavity (mean): | 0.0 | 0.427 |
| concave points (mean): | 0.0 | 0.201 |
| symmetry (mean): | 0.106 | 0.304 |
| fractal dimension (mean): | 0.05 | 0.097 |
| radius (standard error): | 0.112 | 2.873 |
| texture (standard error): | 0.36 | 4.885 |
| perimeter (standard error): | 0.757 | 21.98 |
| area (standard error): | 6.802 | 542.2 |
| smoothness (standard error): | 0.002 | 0.031 |
| compactness (standard error): | 0.002 | 0.135 |
| concavity (standard error): | 0.0 | 0.396 |
| concave points (standard error): | 0.0 | 0.053 |
| symmetry (standard error): | 0.008 | 0.079 |
| fractal dimension (standard error): | 0.001 | 0.03 |
| radius (worst): | 7.93 | 36.04 |
| texture (worst): | 12.02 | 49.54 |
| perimeter (worst): | 50.41 | 251.2 |
| area (worst): | 185.2 | 4254.0 |
| smoothness (worst): | 0.071 | 0.223 |
| compactness (worst): | 0.027 | 1.058 |
| concavity (worst): | 0.0 | 1.252 |
| concave points (worst): | 0.0 | 0.291 |
| symmetry (worst): | 0.156 | 0.664 |
| fractal dimension (worst): | 0.055 | 0.208 |

:Missing Attribute Values: None

:Class Distribution: 212 - Malignant, 357 - Benign

:Creator:  Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian

:Donor: Nick Street

:Date: November, 1995

This is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.
https://goo.gl/U2Uwz2

Features are computed from a digitized image of a fine needle
aspirate (FNA) of a breast mass.  They describe
characteristics of the cell nuclei present in the image.

Separating plane described above was obtained using
Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree
Construction Via Linear Programming." Proceedings of the 4th
Midwest Artificial Intelligence and Cognitive Science Society,
pp. 97-101, 1992], a classification method which uses linear
programming to construct a decision tree.  Relevant features
were selected using an exhaustive search in the space of 1-4
features and 1-3 separating planes.

The actual linear program used to obtain the separating plane
in the 3-dimensional space is that described in:
[K. P. Bennett and O. L. Mangasarian: "Robust Linear
Programming Discrimination of Two Linearly Inseparable Sets",
Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server:

ftp ftp.cs.wisc.edu
cd math-prog/cpo-dataset/machine-learn/WDBC/

.. topic:: References

- W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction
  for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on
  Electronic Imaging: Science and Technology, volume 1905, pages 861-870,
  San Jose, CA, 1993.
- O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and
  prognosis via linear programming. Operations Research, 43(4), pages 570-577,
  July-August 1995.
- W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques
  to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994)
  163-171.

سوال 1.2

متد های ()keys :feature_names: data

```
1 print(data.keys())
2 print(data.feature_names)
3 print(data.data)
```

```
dict_keys(['data', 'target', 'target_names', 'DESCR', 'feature_names', 'filename'])
['mean radius' 'mean texture' 'mean perimeter' 'mean area'
 'mean smoothness' 'mean compactness' 'mean concavity'
 'mean concave points' 'mean symmetry' 'mean fractal dimension'
 'radius error' 'texture error' 'perimeter error' 'area error'
 'smoothness error' 'compactness error' 'concavity error'
 'concave points error' 'symmetry error' 'fractal dimension error'
 'worst radius' 'worst texture' 'worst perimeter' 'worst area'
 'worst smoothness' 'worst compactness' 'worst concavity'
 'worst concave points' 'worst symmetry' 'worst fractal dimension']
[[1.799e+01 1.038e+01 1.228e+02 ... 2.654e-01 4.601e-01 1.189e-01]
 [2.057e+01 1.777e+01 1.329e+02 ... 1.860e-01 2.750e-01 8.902e-02]
 [1.969e+01 2.125e+01 1.300e+02 ... 2.430e-01 3.613e-01 8.758e-02]
 ...
 [1.660e+01 2.808e+01 1.083e+02 ... 1.418e-01 2.218e-01 7.820e-02]
 [2.060e+01 2.933e+01 1.401e+02 ... 2.650e-01 4.087e-01 1.240e-01]
 [7.760e+00 2.454e+01 4.792e+01 ... 0.000e+00 2.871e-01 7.039e-02]]
```

Pandas DataFrame

.

```
1 import pandas as pd
2 df=pd.DataFrame(data.data,columns=data.feature_names)
3 df.head()
```

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | mean fractal dimension | ... | worst radius | worst texture | worst perimeter | worst area | worst smoothness |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.3001 | 0.14710 | 0.2419 | 0.07871 | ... | 25.38 | 17.33 | 184.60 | 2019.0 | 0.162 |
| 1 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.0869 | 0.07017 | 0.1812 | 0.05667 | ... | 24.99 | 23.41 | 158.80 | 1956.0 | 0.123 |
| 2 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.1974 | 0.12790 | 0.2069 | 0.05999 | ... | 23.57 | 25.53 | 152.50 | 1709.0 | 0.144 |
| 3 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.2414 | 0.10520 | 0.2597 | 0.09744 | ... | 14.91 | 26.50 | 98.87 | 567.7 | 0.209 |
| 4 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.1980 | 0.10430 | 0.1809 | 0.05883 | ... | 22.54 | 16.67 | 152.20 | 1575.0 | 0.137 |

5 rows × 30 columns

`df.describe()`

```
1 df.describe()
```

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | mean fractal dimension | ... | worst radius | wor: textu |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | ... | 569.000000 | 569.00000 |
| mean | 14.127292 | 19.289649 | 91.969033 | 654.889104 | 0.096360 | 0.104341 | 0.088799 | 0.048919 | 0.181162 | 0.062798 | ... | 16.269190 | 25.67722 |
| std | 3.524049 | 4.301036 | 24.298981 | 351.914129 | 0.014064 | 0.052813 | 0.079720 | 0.038803 | 0.027414 | 0.007060 | ... | 4.833242 | 6.14625 |
| min | 6.981000 | 9.710000 | 43.790000 | 143.500000 | 0.052630 | 0.019380 | 0.000000 | 0.000000 | 0.106000 | 0.049960 | ... | 7.930000 | 12.02000 |
| 25% | 11.700000 | 16.170000 | 75.170000 | 420.300000 | 0.086370 | 0.064920 | 0.029560 | 0.020310 | 0.161900 | 0.057700 | ... | 13.010000 | 21.08000 |
| 50% | 13.370000 | 18.840000 | 86.240000 | 551.100000 | 0.095870 | 0.092630 | 0.061540 | 0.033500 | 0.179200 | 0.061540 | ... | 14.970000 | 25.41000 |
| 75% | 15.780000 | 21.800000 | 104.100000 | 782.700000 | 0.105300 | 0.130400 | 0.130700 | 0.074000 | 0.195700 | 0.066120 | ... | 18.790000 | 29.72000 |
| max | 28.110000 | 39.280000 | 188.500000 | 2501.000000 | 0.163400 | 0.345400 | 0.426800 | 0.201200 | 0.304000 | 0.097440 | ... | 36.040000 | 49.54000 |

8 rows × 30 columns

کلید Target را برای دیتافریم تعریف می کنیم

```
1 df["target"]=data.get("target")
2 df.set_index('target',inplace=True)
3 df.head()
```

| target | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.3001 | 0.14710 | 0.2419 | 0.07871 | ... | 25.38 | 17.33 | 184.60 | 2019.0 |
| 0 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.0869 | 0.07017 | 0.1812 | 0.05667 | ... | 24.99 | 23.41 | 158.80 | 1956.0 |
| 0 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.1974 | 0.12790 | 0.2069 | 0.05999 | ... | 23.57 | 25.53 | 152.50 | 1709.0 |
| 0 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.2414 | 0.10520 | 0.2597 | 0.09744 | ... | 14.91 | 26.50 | 98.87 | 567.7 |
| 0 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.1980 | 0.10430 | 0.1809 | 0.05883 | ... | 22.54 | 16.67 | 152.20 | 1575.0 |

5 rows × 30 columns

این دو دستور یکی تعداد داده ها را در هر کدام از target های malignant و bengin نشان می دهد و دیگری نام هر کدام از تارگت ها را نمایان می کند.

```
1 print(df.index.value_counts())
2 print(data.target_names)
```

```
1     357
0     212
Name: target, dtype: int64
['malignant' 'benign']
```

داده های train , test را با X و Y می سازیم

```
1 import sklearn.model_selection as model
2 X=data.data
3 Y=data.target
4 x_train,x_test,y_train,y_test=model.train_test_split(X,Y)
```

مشاهده می کنیم که داده های train و test با نسبت دلخواه تقسیم شده اند.

```
1 print("Data Size: " + str(df.shape))
2 print("X_Train Size: " + str(x_train.shape))
3 print("X_Test Size: " + str(x_test.shape))
```

```
4 print("Y_Train Size: " + str(y_train.shape))
5 print("Y_Test Size: " + str(y_test.shape))
```

```
Data Size: (569, 30)
X_Train Size: (426, 30)
X_Test Size: (143, 30)
Y_Train Size: (426,)
Y_Test Size: (143,)
```

سوال 1.9

```
1 from sklearn.neighbors import KNeighborsClassifier
2 neighb=KNeighborsClassifier(n_neighbors=6)
3 neighb.fit(x_train,y_train)
4 neighb.score(x_test,y_test)
```

دقت دسته بندی برابر است با :

**0.9370629370629371**

سوال 1.10

با استفاده از متد predict داده ها را پیش بینی می کنیم و داخل Y_pred می ریزیم

```
1 y_pred=neighb.predict(x_test)
```

سوال1.11

توضیح Predict :

در حقیقت این تابع با استفاده از داده های train شده در مدل داده های test را پیش بینی می کند.

MinMaxScaler

```
1
2 from sklearn.preprocessing import MinMaxScaler
3 MinMax=MinMaxScaler()
4 MinMax.fit(X)
5 norm_x_train=MinMax.transform(x_train)
6 norm_x_test=MinMax.transform(x_test)
7
```

```
1 neighb=KNeighborsClassifier(n_neighbors=6)
2 neighb.fit(norm_x_train,y_train)
```

```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
          metric_params=None, n_jobs=None, n_neighbors=6, p=2,
          weights='uniform')
```

```
1 print("test :" + str(neighb.score(norm_x_test,y_test)))
2 print("train :" + str(neighb.score(norm_x_train,y_train)))
```

دقت دیتاست های trainو test به شرح زیر می باشد :
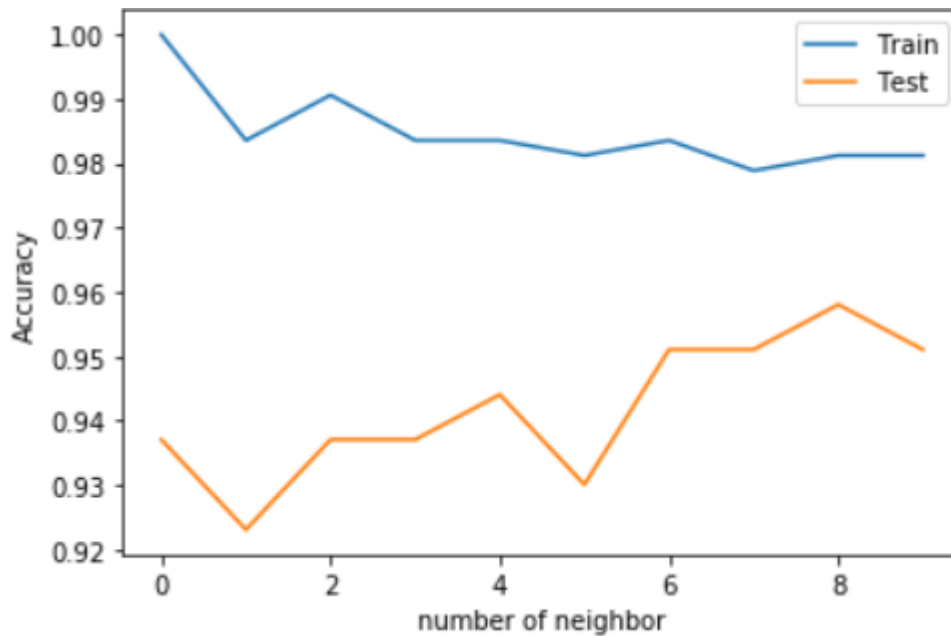
```
test :0.9300699300699301
train :0.9812206572769953
```

```
 1 train_accuracy=[]
 2 test_accuracy=[]
 3 for i in range(1,11):
 4     neighb=KNeighborsClassifier(n_neighbors=i)
 5     neighb.fit(norm_x_train,y_train)
 6     train_accuracy.insert(i,neighb.score(norm_x_train,y_train))
 7     test_accuracy.insert(i,neighb.score(norm_x_test,y_test))
 8
 9 print(train_accuracy)
10 print(test_accuracy)
```

```
[1.0, 0.9835680751173709, 0.9906103286384976, 0.9835680751173709, 0.9835680751173709, 0.9812206572769953, 0.9835680751173709,
0.9788732394366197, 0.9812206572769953, 0.9812206572769953]
[0.9370629370629371, 0.9230769230769231, 0.9370629370629371, 0.9370629370629371, 0.9440559440559441, 0.9300699300699301, 0.951
048951048951, 0.951048951048951, 0.958041958041958, 0.951048951048951]
```

```
1 import matplotlib.pyplot as plt
2 plt.plot(train_accuracy,label="Train")
3 plt.plot(test_accuracy, label="Test")
4 plt.ylabel('Accuracy')
5 plt.xlabel('number of neighbor')
6 plt.legend(loc='upper right')
7 plt.show()
```

همانطور که مشاهده می کنید در داده های train با افزایش تعداد همسایگی دقت کاهش پیدا می کند چرا که مدل دید جامع تری دارد و مجبور است همسایگان بیشتری را در بر بگیرد. ولی در مدل تست تقریبا برعکس است و با افزایش همسایگی دقت مدل افزایش چشم گیری را خواهد داشت. طبیعی است چرا که داده ای که در مدل train نشده است با دیدن همسایگان بیشتر با احتمال بیشتری دقیق پیشبینی می شود ولی داده ای که قبلا در مدل وجود داشته با افزایش همسایگی جوانب را نیز افزایش دادیم به همین جهت دقت کمتری را خواهیم داشت.

```
1 import numpy as np
2 import pandas as pd
3 df=pd.read_csv('vehicle.csv')
4 df.head()
```

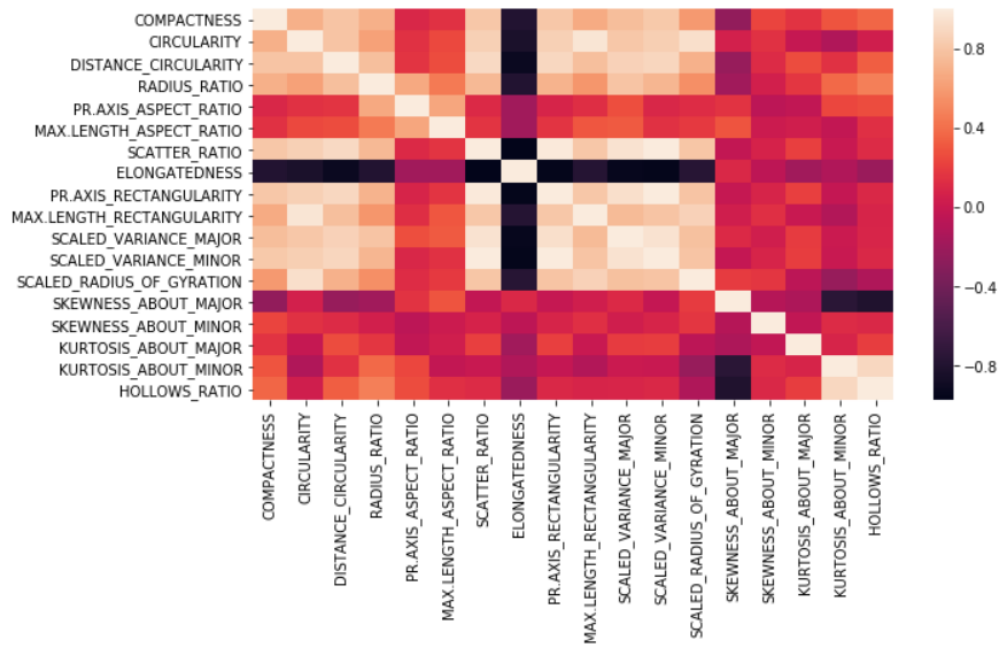| | COMPACTNESS | CIRCULARITY | DISTANCE_CIRCULARITY | RADIUS_RATIO | PR.AXIS_ASPECT_RATIO | MAX.LENGTH_ASPECT_RATIO | SCATTER_RATIO | ELONG. |
|---|---|---|---|---|---|---|---|---|
| 0 | 95 | 48 | 83 | 178 | 72 | 10 | 162 | |
| 1 | 91 | 41 | 84 | 141 | 57 | 9 | 149 | |
| 2 | 104 | 50 | 106 | 209 | 66 | 10 | 207 | |
| 3 | 93 | 41 | 82 | 159 | 63 | 9 | 144 | |
| 4 | 85 | 44 | 70 | 205 | 103 | 52 | 149 | |

## سوال 2.3

```
1 print(pd.unique(df["Class"]))
```

```
['van' 'saab' 'bus' 'opel']
```

## سوال2.4

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 plt.figure(figsize=(10,5))
4 sns.heatmap(df.corr())
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x28c347289e8>
```

```
1 X=df.loc[:,df.columns!="Class"]
2 Y=df["Class"]
```

```
1 import sklearn.model_selection as model
2 x_train,x_test,y_train,y_test=model.train_test_split(X,Y,test_size=0.2)
```

```
1 from sklearn.tree import DecisionTreeClassifier
2 decTree=DecisionTreeClassifier(max_depth = 5, max_features=4,
3 criterion='entropy')
```

```
decTree.fit(x_train,y_train)
```

.

```
1 import random
2 from scipy.stats import randint
3 Params = {"max_depth":[3 , None],
4           "max_feature": randint(1,9),
5           "min_samples_leaf": randint(1,9)}
6 print(Params)
```

```
{'max_depth': [3, None], 'max_feature': <scipy.stats._distn_infrastructure.rv_frozen object at 0x000002D0D0438A58>, 'min_sampl
es_leaf': <scipy.stats._distn_infrastructure.rv_frozen object at 0x000002D0D0438B70>}
```

```
1 from sklearn.model_selection import RandomizedSearchCV
2 tree = DecisionTreeClassifier()
3
4 tree_cv = RandomizedSearchCV(tree, Params, cv=5)
5
6 tree_cv.fit(X, Y)
```

```
RandomizedSearchCV(cv=5, error_score='raise-deprecating',
          estimator=DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
            max_features=None, max_leaf_nodes=None,
            min_impurity_decrease=0.0, min_impurity_split=None,
            min_samples_leaf=1, min_samples_split=2,
            min_weight_fraction_leaf=0.0, presort=False, random_state=None,
            splitter='best'),
          fit_params=None, iid='warn', n_iter=10, n_jobs=None,
          param_distributions={'max_depth': [3, None], 'max_features': <scipy.stats._distn_infrastructure.rv_frozen object at 0
x000001DB90FC6470>, 'min_samples_leaf': <scipy.stats._distn_infrastructure.rv_frozen object at 0x000001DB90FC6588>},
          pre_dispatch='2*n_jobs', random_state=None, refit=True,
          return_train_score='warn', scoring=None, verbose=0)
```

```
1 print("Tuned Decision Tree Parameters: {}".format(tree_cv.best_params_))
2 print("Best score is {}".format(tree_cv.best_score_))
```

```
Tuned Decision Tree Parameters: {'max_depth': None, 'max_features': 4, 'min_samples_leaf': 7}
Best score is 0.6855791962174941
```

```
1 clf_entropy = DecisionTreeClassifier(criterion = "entropy", max_depth =
2 None, max_features=6, min_samples_leaf=8)
3 clf_entropy.fit(x_train, y_train)
4
5 y_pred = clf_entropy.predict(x_test)
6
7 from sklearn.metrics import accuracy_score
  print("Accuracy is ", accuracy_score(y_test,y_pred)*100)
```

دقت بدست آمده برابر است با : 62.35294117647059

با افزایش مقدار cv دقت مدل نیز افزایش پیدا میکند

این تابع جهت مشخص کردن اهمیت feature ها برای انتخاب جهت Decision Node را تعیین می کند.

```
1 fimportant = dict(zip(df.columns, clf_entropy.feature_importances_))
2 for key,val in fimportant.items():
3     print(key, "=>", val)
```

```
COMPACTNESS => 0.057243604187455936
CIRCULARITY => 0.021646551795873187
DISTANCE_CIRCULARITY => 0.0331075292040589
RADIUS_RATIO => 0.0
PR.AXIS_ASPECT_RATIO => 0.06292942492213029
MAX.LENGTH_ASPECT_RATIO => 0.11923060972148906
SCATTER_RATIO => 0.0210759292214674386
ELONGATEDNESS => 0.19069324997836132
PR.AXIS_RECTANGULARITY => 0.0
MAX.LENGTH_RECTANGULARITY => 0.08826263261284722
SCALED_VARIANCE_MAJOR => 0.0
SCALED_VARIANCE_MINOR => 0.19412643509730335
SCALED_RADIUS_OF_GYRATION => 0.014039747171418277
SKEWNESS_ABOUT_MAJOR => 0.07560526810188511
SKEWNESS_ABOUT_MINOR => 0.0019039046049978466
KURTOSIS_ABOUT_MAJOR => 0.026054267154103403
KURTOSIS_ABOUT_MINOR => 0.07847339011861136
HOLLOWS_RATIO => 0.015607456114790403
```

```python
from sklearn import tree
import pydotplus
import collections
from sklearn.tree import export_graphviz
from IPython.display import Image
import matplotlib
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
matplotlib.style.use('ggplot')

dot_data = tree.export_graphviz(clf_entropy,
                                feature_names=df.iloc[0,0:-1],
                                out_file=None,
                                filled=True,
                                rounded=True)
graph = pydotplus.graph_from_dot_data(dot_data)

colors = ('turquoise', 'orange')
edges = collections.defaultdict(list)

for edge in graph.get_edge_list():
    edges[edge.get_source()].append(int(edge.get_destination()))

for edge in edges:
    edges[edge].sort()
    for i in range(2):
        dest = graph.get_node(str(edges[edge][i]))[0]
        dest.set_fillcolor(colors[i])
```
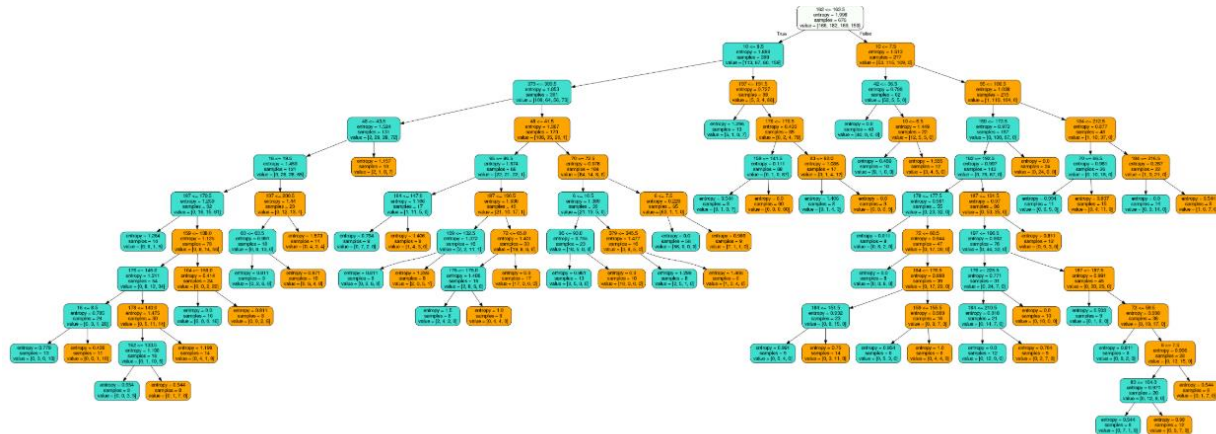
```
31 graph.write_png('dot_data.png')
32 Image('dot_data.png')
```

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn import datasets
5 iris=datasets.load_iris()
6 df=pd.DataFrame(iris.data,columns=iris.feature_names)
7 df.head()
```

|   | sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) |
|---|---|---|---|---|
| **0** | 5.1 | 3.5 | 1.4 | 0.2 |
| **1** | 4.9 | 3.0 | 1.4 | 0.2 |
| **2** | 4.7 | 3.2 | 1.3 | 0.2 |
| **3** | 4.6 | 3.1 | 1.5 | 0.2 |
| **4** | 5.0 | 3.6 | 1.4 | 0.2 |

```
1 from sklearn.cluster import KMeans
```

```
2 X = iris.data
3 y = iris.target
4 kmeans = KMeans(n_clusters=3)
5 kmeans.fit(X)
6 y_kmeans = kmeans.predict(X)
```

```
1 centroids = kmeans.cluster_centers_
```

```
  plt.scatter(X[:, 1], X[:, 3], c=y_kmeans, s=50, cmap='viridis')
1 plt.scatter(centroids[:, 1], centroids[:, 3], c='blue',marker="X", s=200,
2 alpha=0.9);
```

```
1 kmeans.inertia_
```

**78.94084142614602**

```
1 interia_list=[]
2 for i in range(5):
3     kmeans = KMeans(n_clusters=i+1)
4     kmeans.fit(X)
5     y_kmeans = kmeans.predict(X)
6     interia_list.append(kmeans.inertia_)
7
8 interia_list
```

**[680.8244,
152.36870647733906,
78.94084142614602,
57.31787321428571,
46.53558205128205]**

```
1 plt.bar(list(range(1,6)),interia_list, label="Number Cluster")
2 plt.legend()
3 plt.xlabel('number cluster')
4 plt.ylabel('interia score')
5 plt.title('inertia')
6 plt.show()
```

با افزایش تعداد کلاستر مقدار اینرسی نیز کاهش می یابد چرا که مقادیر درون هر کلاستر منسجم تر می شود تا جایی که هر کلاستر خود یک نود شود. جهت فهمیدن تعداد کلاستر بهینه از دو روش معروف می توان استفاده کرد:

1. Elbow : در این روش با استفاده از معیار فاصله درون خوشه ای اقدام به پیدا کردن کمترین فاصله درون خوشه ای به ازای k می کنیم.

2. میانگین شاخص silhouette: در این روش نیز مانند قبل یک بازه ای از K ها را تست کرده و برای هر k یک میانگین silhouette بدست می آوریم و با رسم نمودار curve آنها بالا ترین یا بیشترین مقدار همان بهینه ترین است

سوال 4.1

```
1 import numpy as np
2 import pandas as pd
3 from sklearn import datasets
4 iris=datasets.load_iris()
5 df=pd.DataFrame(iris.data,columns=iris.feature_names)
6 df.head()
```

```
1 from scipy.cluster.hierarchy import dendrogram, linkage
2 from matplotlib import pyplot as plt
3 X = df
4 Z = linkage(X, 'complete')
```

|   | sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) |
|---|---|---|---|---|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 |

سوال 4.2

```
1 fig = plt.figure(figsize=(25, 10))
2 dn = dendrogram(Z)
```



سوال 4.3

```
1 from scipy.cluster.hierarchy import fcluster
```

```
2 fcluster(Z, t=6, criterion='distance')
```

```
array([2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
       2, 2, 2, 2, 2, 2, 1, 1, 1, 2, 1, 2, 1, 2, 1, 2, 2, 2, 2, 1, 2, 1,
       2, 2, 1, 2, 1, 2, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 1, 2, 1, 1, 1,
       2, 2, 2, 1, 2, 2, 2, 2, 2, 1, 2, 2, 1, 1, 1, 1, 1, 1, 2, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1], dtype=int32)
```

```python
import matplotlib.pyplot as plt
plt.figure(figsize=(10, 8))
plt.scatter(iris.data[:,0], iris.data[:,1], c=clusters, cmap='prism')
plt.show()
```

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn import datasets
5 boston=datasets.load_boston()
6 df=pd.DataFrame(boston.data,columns=boston.feature_names)
7 df.head()
```

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT |
|---|------|-----|-------|------|------|-------|------|--------|-----|-------|---------|--------|-------|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 | 4.98 |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 | 9.14 |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 | 4.03 |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 | 2.94 |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 | 5.33 |

```
1 df["Price"]=boston.target
2 df.head()
```

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | Price |
|---|------|-----|-------|------|------|-------|------|--------|-----|-------|---------|--------|-------|-------|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 | 4.98 | 24.0 |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 | 9.14 | 21.6 |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 | 4.03 | 34.7 |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 | 2.94 | 33.4 |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 | 5.33 | 36.2 |

```
1 X=df[["CRIM","ZN"]]
2 Y=df["Price"]
```

```
1 import sklearn.model_selection as model
2 x_train,x_test,y_train,y_test=model.train_test_split(X,Y,test_size=0.3)
```

```
1 from sklearn.linear_model import LinearRegression
2 lm = LinearRegression()
3 lm.fit(x_train, y_train)
4 y_pred = lm.predict(x_test)
5 plt.scatter(y_test, y_pred)
6 plt.xlabel("Prices:")
7 plt.ylabel("Predicted prices")
8 plt.title("Prices vs Predicted prices")
```

Text(0.5,1,'Prices vs Predicted prices')

## Prices vs Predicted prices

```
1 from sklearn.metrics import mean_squared_error
2 print("MSE : " + str(mean_squared_error(y_test,y_pred)))
```

MSE : 65.20918707439772

```
 1  X=df.loc[:,df.columns != "Price"]
 2  Y=df["Price"]
 3
 4  x_train,x_test,y_train,y_test=model.train_test_split(X,Y,test_size=0.3)
 5
 6
 7  lm = LinearRegression()
 8  lm.fit(x_train, y_train)
 9
10  y_pred = lm.predict(x_test)
11
12  plt.scatter(y_test, y_pred)
13  plt.xlabel("Prices:")
14  plt.ylabel("Predicted prices")
15  plt.title("Prices vs Predicted prices")ed)))
```

Text(0.5,1,'Prices vs Predicted prices')

```
1 print("MSE : " + str(mean_squared_error(y_test,y_pred)))
```

**MSE : 20.381895735190685**

مقدار MSE کاهش چشم گیری را داشته و طبیعی است که با افزایش تعدادfeature ها دقت افزایش و خطا نیز کاهش می یابد. چرا که
جوانب بیشتری را در مدل کردن مد نظر دارد.

```
1 from sklearn.model_selection import cross_val_score
2 import statistics
3 cvs=cross_val_score(lm, X, Y, cv=5)
4 print(cvs)
5 print("Mean cvs: "+ str(statistics.mean(cvs)))
```

```
[ 0.63861069  0.71334432  0.58645134  0.07842495 -0.26312455]
Mean cvs: 0.3507413509325258
```

```
1 import numpy as np
2 import pandas as pd
3 import sklearn.datasets as datasets
4 bcancer=datasets.load_breast_cancer()
5 df=pd.DataFrame(bcancer.data,columns=bcancer.feature_names)
6 df.head()
```

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | mean fractal dimension | ... | worst radius | worst texture | worst perimeter | worst area | worst smoothnes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.3001 | 0.14710 | 0.2419 | 0.07871 | ... | 25.38 | 17.33 | 184.60 | 2019.0 | 0.162 |
| 1 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.0869 | 0.07017 | 0.1812 | 0.05667 | ... | 24.99 | 23.41 | 158.80 | 1956.0 | 0.123 |
| 2 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.1974 | 0.12790 | 0.2069 | 0.05999 | ... | 23.57 | 25.53 | 152.50 | 1709.0 | 0.144 |
| 3 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.2414 | 0.10520 | 0.2597 | 0.09744 | ... | 14.91 | 26.50 | 98.87 | 567.7 | 0.209 |
| 4 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.1980 | 0.10430 | 0.1809 | 0.05883 | ... | 22.54 | 16.67 | 152.20 | 1575.0 | 0.137 |

5 rows × 30 columns

سوال 6.2

```
1 import sklearn.model_selection as model
2 X=bcancer.data
3 Y=bcancer.target
4 x_train,x_test,y_train,y_test=model.train_test_split(X,Y,test_size=0.2)
5
6 from sklearn.neighbors import KNeighborsClassifier
7 neighb=KNeighborsClassifier(n_neighbors=8)
8 neighb.fit(x_train,y_train)
9 y_pred=neighb.predict(x_test)
```

سوال 6.3

```
1 from sklearn.metrics import confusion_matrix,classification_report
```

سوال 6.4

```
1 print(confusion_matrix(y_test,y_pred))
```

```
[[41  3]
 [ 2 68]]
```

این چهار عدد نشان دهنده 4 حالت پیش بینی ما و درستی جواب می باشد که شامل :TN , TP , FN , FP

<div align="center">سوال 6.5</div>

```
1 print(classification_report(y_test,y_pred))
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.95 | 0.93 | 0.94 | 44 |
| 1 | 0.96 | 0.97 | 0.96 | 70 |
| micro avg | 0.96 | 0.96 | 0.96 | 114 |
| macro avg | 0.96 | 0.95 | 0.95 | 114 |
| weighted avg | 0.96 | 0.96 | 0.96 | 114 |

میزان دقت در دو دسته بندی را نشان می دهد که هر کدام از این 0 و 1 ها نمایان گر یک دسته می باشند.

<div align="center">سوال 6.6</div>

```
1 from sklearn.preprocessing import normalize
2 normalize(confusion_matrix(y_test,y_pred),norm="l1")
```

```
array([[0.93181818, 0.06818182],
       [0.02857143, 0.97142857]])
```

```
1  dff = pd.DataFrame(Normal,
2  index=['benign','malignant'],columns=['benign','malignant'])
   print (dff)
```

```
              benign   malignant
benign       0.958333   0.041667
malignant    0.106061   0.893939
```

هر چه در این نمونه از نمودار ها به محور x=y نزدیک تر باشد خطای بیشتری دارد به همین منظور می توان گفت هر چه حالت elbow داشته باشت خطای آن کم تر و در نتیجه جواب ما ایده آل تر خواهد بود

```
1 y_pred_prob=neighb.predict_proba(x_test)
2 print(y_pred_prob)
```

```
[[1.     0.    ]
 [1.     0.    ]
 [0.875 0.125]
 [0.     1.    ]
 [0.125 0.875]
 [0.5   0.5  ]
 [0.     1.    ]
 [0.     1.    ]
 [1.     0.    ]
 [0.     1.    ]
 [0.     1.    ]
 [0.625 0.375]
 [0.     1.    ]
 [0.875 0.125]
 [0.     1.    ]
 [1.     0.    ]
 [0.     1.    ]
 [1.     0.    ]
 [1.     0.    ]
 [0.25  0.75 ]
```

<div style="text-align:center">سوال 6.10</div>

```python
1 from sklearn.metrics import roc_curve
2 roc_curve(y_pred,y_pred_prob)
```

<div style="text-align:center">مشکل در ورودی تابع داشتیم و نتونستیم شکل ROC را رسم کنیم ☹</div>

<div style="text-align:center">سوال 7.1</div>

```python
1 import pandas as pd
```

```
2 import numpy as np
3 df=pd.read_excel("OnlineRetail.xlsx")
4 df.head()
```

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country |
|---|---|---|---|---|---|---|---|---|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 2010-12-01 08:26:00 | 2.55 | 17850.0 | United Kingdom |
| 1 | 536365 | 71053 | WHITE METAL LANTERN | 6 | 2010-12-01 08:26:00 | 3.39 | 17850.0 | United Kingdom |
| 2 | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 8 | 2010-12-01 08:26:00 | 2.75 | 17850.0 | United Kingdom |
| 3 | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 2010-12-01 08:26:00 | 3.39 | 17850.0 | United Kingdom |
| 4 | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 6 | 2010-12-01 08:26:00 | 3.39 | 17850.0 | United Kingdom |

سوال 7.2

```
1 ! pip install mlxtend --upgrade --no-deps
```

سوال 7.3

```
1 from mlxtend.frequent_patterns import apriori,association_rules
```

سوال 7.4

```
1 df["Description"]=df["Description"].str.strip()
```

```
1 print("Orginal Size : " + str(df.size))
2 df["InvoiceNo"].replace('', np.nan, inplace=True)
3 df.dropna(subset=['InvoiceNo'], inplace=True)
4 print("Reduced Size : " + str(df.size))
5
6 df["InvoiceNo"]=df["InvoiceNo"].astype("str")
```

```
Orginal Size : 4335272
Reduced Size : 4335272
```

```
1 df=df[~df.InvoiceNo.str.contains("C")]
```

```
1 basket = (df[df['Country'] =="France"]
2 .groupby(['InvoiceNo', 'Description'])['Quantity']
3 .sum().unstack().reset_index().fillna(0).set_index('InvoiceNo'))
4 basket.head()
```

| Description | 10 COLOUR SPACEBOY PEN | 12 COLOURED PARTY BALLOONS | 12 EGG HOUSE PAINTED WOOD | 12 MESSAGE CARDS WITH ENVELOPES | 12 PENCIL SMALL TUBE WOODLAND | 12 PENCILS SMALL TUBE RED RETROSPOT | 12 PENCILS SMALL TUBE SKULL | 12 PENCILS TALL TUBE POSY | 12 PENCILS TALL TUBE RED RETROSPOT | 12 PENCILS TALL TUBE WOODLAND | ... | WRAP VINTAGE PETALS DESIGN | YELLO CO RA PAF FASHIC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **InvoiceNo** | | | | | | | | | | | | | |
| **536370** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | |
| **536852** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | |
| **536974** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | |
| **537065** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | |
| **537463** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | |

5 rows × 1563 columns

```
1 basket=basket.applymap(lambda x: 1 if x > 0 else 0)
2 basket.head()
```

| Description | 10 COLOUR SPACEBOY PEN | 12 COLOURED PARTY BALLOONS | 12 EGG HOUSE PAINTED WOOD | 12 MESSAGE CARDS WITH ENVELOPES | 12 PENCIL SMALL TUBE WOODLAND | 12 PENCILS SMALL TUBE RED RETROSPOT | 12 PENCILS SMALL TUBE SKULL | 12 PENCILS TALL TUBE POSY | 12 PENCILS TALL TUBE RED RETROSPOT | 12 PENCILS TALL TUBE WOODLAND | ... | WRAP VINTAGE PETALS DESIGN | YELLO CO RA PAF FASHIC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **InvoiceNo** | | | | | | | | | | | | | |
| **536370** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | |
| **536852** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | |
| **536974** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | |
| **537065** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | |
| **537463** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | |

5 rows × 1563 columns

```
1 basket=basket.drop("POSTAGE",axis=1)
```

```
1 frequent_itemsets=apriori(basket,min_support=0.07)
2 frequent_itemsets.head()
```

| | support | itemsets |
|---|---|---|
| **0** | 0.071429 | (35) |
| **1** | 0.096939 | (61) |
| **2** | 0.102041 | (64) |
| **3** | 0.094388 | (65) |
| **4** | 0.081633 | (103) |

```
1 rules = association_rules(frequent_itemsets, metric="lift",
  min_threshold=1)
2 rules.head()
```

| | antecedents | consequents | antecedent support | consequent support | support | confidence | lift | leverage | conviction |
|---|---|---|---|---|---|---|---|---|---|
| **0** | (64) | (61) | 0.102041 | 0.096939 | 0.073980 | 0.725000 | 7.478947 | 0.064088 | 3.283859 |
| **1** | (61) | (64) | 0.096939 | 0.102041 | 0.073980 | 0.763158 | 7.478947 | 0.064088 | 3.791383 |
| **2** | (65) | (61) | 0.094388 | 0.096939 | 0.079082 | 0.837838 | 8.642959 | 0.069932 | 5.568878 |
| **3** | (61) | (65) | 0.096939 | 0.094388 | 0.079082 | 0.815789 | 8.642959 | 0.069932 | 4.916181 |
| **4** | (64) | (65) | 0.102041 | 0.094388 | 0.073980 | 0.725000 | 7.681081 | 0.064348 | 3.293135 |

```
1    rules[ (rules['lift'] >= 6) &
             (rules['confidence'] >= 0.8) ]
```

| | antecedents | consequents | antecedent support | consequent support | support | confidence | lift | leverage | conviction |
|---|---|---|---|---|---|---|---|---|---|
| 2 | (65) | (61) | 0.094388 | 0.096939 | 0.079082 | 0.837838 | 8.642959 | 0.069932 | 5.568878 |
| 3 | (61) | (65) | 0.096939 | 0.094388 | 0.079082 | 0.815789 | 8.642959 | 0.069932 | 4.916181 |
| 16 | (1267) | (1247) | 0.127551 | 0.132653 | 0.102041 | 0.800000 | 6.030769 | 0.085121 | 4.336735 |
| 18 | (1266) | (1267) | 0.137755 | 0.127551 | 0.122449 | 0.888889 | 6.968889 | 0.104878 | 7.852041 |
| 19 | (1267) | (1266) | 0.127551 | 0.137755 | 0.122449 | 0.960000 | 6.968889 | 0.104878 | 21.556122 |
| 20 | (1266, 1267) | (1247) | 0.122449 | 0.132653 | 0.099490 | 0.812500 | 6.125000 | 0.083247 | 4.625850 |
| 21 | (1266, 1247) | (1267) | 0.102041 | 0.127551 | 0.099490 | 0.975000 | 7.644000 | 0.086474 | 34.897959 |
| 22 | (1267, 1247) | (1266) | 0.102041 | 0.137755 | 0.099490 | 0.975000 | 7.077778 | 0.085433 | 34.489796 |

به عنوان مثال آیتم 65 را اگر موجود باشد با دقت 0.83 امکان دارد آیتم 61 نیز باشد. یا به صورت دیگر اگر آیتم ها را مربوط به سوپرمارکت بدانیم می توان گفت اگر آیتم 1266 و 1267 با هم انتخاب شوند با confidence 0.97 می توان گفت آیتم 1267 نیز خریداری می شود.