

Mini project 2

Data Set Information: Seven different types of dry beans were used in this project, taking into account the features such as form, shape, type, and structure by the market situation. Use best machine learning algorithm to classify the seven types of beans in Turkey; Barbunya, Bombay, Cali, Dermason, Horoz, Seker and Sira, depending only on dimension and shape features of bean varieties with no external discriminatory features

Features Information

1. Area (A): The area of a bean zone and the number of pixels within its boundaries.
2. Perimeter (P): Bean circumference is defined as the length of its border.
3. Major axis length (L): The distance between the ends of the longest line that can be drawn from a bean.
4. Minor axis length (l): The longest line that can be drawn from the bean while standing perpendicular to the main axis.
5. Aspect ratio (r): Defines the relationship between L and l.
6. Eccentricity (Ec): Eccentricity of the ellipse having the same moments as the region.
7. Convex area (C): Number of pixels in the smallest convex polygon that can contain the area of a bean seed.
8. Equivalent diameter (Ed): The diameter of a circle having the same area as a bean seed area.
9. Extent (Ex): The ratio of the pixels in the bounding box to the pixels in the convex hull to those found in beans.
10. Solidity (S): Also known as convexity, the ratio of the pixels in the convex hull to those found in beans.
11. Roundness (R): Calculated with the following formula: $(4\pi A)/(P^2)$
12. Compactness (Co): Measures the roundness of an object. E_d/L .
13. ShapeFactor1 (SF1)
14. ShapeFactor2 (SF2)
15. ShapeFactor3 (SF3)
16. ShapeFactor4 (SF4)
17. Class (Seker, Barbunya, Bombay, Cali, Dermason, Horoz and Sira)

import libraries

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
```

Data Loading

```
In [5]: # read the data
df = pd.read_csv('data.csv')
```

```
In [6]: # checking whether data is loaded or not (it shows overview of the data)
df.head()
```

```
Out[6]:
```

	Area	Perimeter	MajorAxisLength	MinorAxisLength	AspectRatio	Eccentricity	ConvexArea	EquivalentDiameter	Extent	Solidity	roundness	Compactness	ShapeFactor1	ShapeFactor2
0	28395	610.291	208.178117	173.888747	1.197191	0.549812	28715	190.141097	0.763923	0.988856	0.958027	0.91358	0.007332	0.00314
1	28734	638.018	200.524796	182.734419	1.097356	0.411785	29172	191.272751	0.783968	0.984986	0.887034	0.953861	0.006979	0.00356
2	29380	624.110	212.826130	175.931143	1.209713	0.562727	29690	193.410904	0.778113	0.989559	0.947849	0.908774	0.007244	0.00304
3	30008	645.884	210.557999	182.516516	1.153638	0.498616	30724	195.467062	0.782681	0.976966	0.903936	0.928329	0.007017	0.00321
4	30140	620.134	201.847882	190.279279	1.060798	0.333880	30417	195.896503	0.773098	0.990893	0.984877	0.970516	0.006697	0.00366

```
In [15]: # checking shape
print(df.shape)
```

(13611, 17)

```
In [16]: df.columns
```

```
Out[16]:
```

Index(['Area', 'Perimeter', 'MajorAxisLength', 'MinorAxisLength', 'AspectRatio', 'Eccentricity', 'ConvexArea', 'EquivalentDiameter', 'Extent', 'Solidity', 'roundness', 'Compactness', 'ShapeFactor1', 'ShapeFactor2', 'ShapeFactor3', 'ShapeFactor4', 'class'], dtype='object')

```
In [17]: # checking size
print(df.size)
```

231387

```
In [18]: # describe the data
df.describe()
```

```
Out[18]:
```

	Area	Perimeter	MajorAxisLength	MinorAxisLength	AspectRatio	Eccentricity	ConvexArea	EquivalentDiameter	Extent	Solidity	roundness	Compactness	ShapeFactor1	ShapeFactor2
count	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000	13611.000000
mean	53408.284549	855.283459	320.141867	202.270714	1.583242	0.750805	53768.200206	253.064220	0.747933	0.967143	0.873282	0.7998		
std	29324.95717	214.896956	185.894186	44.700951	0.246678	0.092002	29774.915817	59.177120	0.049086	0.044660	0.059520	0.6167		
min	204.20.000000	524.236000	183.011615	122.512663	1.024868	0.218951	20684.000000	161.243764	0.553513	0.919246	0.496918	0.6405		
25%	3638.000000	703.232500	253.030333	175.931143	1.472582	0.432307	3714.500000	215.068003	0.718363	0.985670	0.832098	0.6241		
50%	44652.000000	794.941000	296.083367	192.317333	1.51124	0.674441	45174.000000	238.439266	0.759859	0.989283	0.883157	0.8012		
75%	61332.000000	977.213000	376.495012	217.031741	1.707109	0.810466	62294.000000	279.446467	0.786851	0.990013	0.916869	0.8342		
max	254616.000000	1985.270000	738.060154	460.186497	2.430306	0.911423	263261.000000	559.374358	0.866195	0.994677	0.990685	0.9873		

```
In [19]: # checking information about data (gives count of data types)
df.info()
```

<class 'pandas.core.frame.DataFrame'>

```
Ram usage: 1.8* MB
RangeIndex: 13611 entries, 0 to 13610
Data columns (total 17 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   Area            13611 non-null   int64  
 1   Perimeter       13611 non-null   float64
 2   MajorAxisLength 13611 non-null   float64
 3   MinorAxisLength 13611 non-null   float64
 4   AspectRatio      13611 non-null   float64
 5   Eccentricity    13611 non-null   float64
 6   ConvexArea       13611 non-null   float64
 7   EquivalentDiameter 13611 non-null   float64
 8   Extent           13611 non-null   float64
 9   Solidity          13611 non-null   float64
 10  roundness        13611 non-null   float64
 11  compactness       13611 non-null   float64
 12  ShapeFactor1     13611 non-null   float64
 13  ShapeFactor2     13611 non-null   float64
 14  ShapeFactor3     13611 non-null   float64
 15  ShapeFactor4     13611 non-null   float64
 16  Class             13611 non-null   object 
dtypes: float64(14), int64(2), object(1)
memory usage: 1.8* MB
```

```
In [20]: # unique = it shows different data types present in the data
df.dtypes.unique()
```

array('dtype('float64'), dtype('float64'), dtype('O'), dtype('object'))

```
In [21]: # data cleaning step
# drop missing values
print(df.isnull().sum())
df = df.dropna(axis=0)
```

Area: 0

Perimeter: 0

MajorAxisLength: 0

MinorAxisLength: 0

AspectRatio: 0

Eccentricity: 0

ConvexArea: 0

EquivalentDiameter: 0

Extent: 0

Solidity: 0

roundness: 0

compactness: 0

ShapeFactor1: 0

ShapeFactor2: 0

ShapeFactor3: 0

ShapeFactor4: 0

Class: 0

dtype: int64

```
In [23]: df.duplicated().sum()
68
```

```
In [24]: # check for unique values in the target variable
# target variable is class and others are features
print(df['Class'].unique())
['SEKER', 'BARBUNYA', 'BOMBAY', 'CALI', 'HOROZ', 'SIRA', 'DERMASON']
```

```
In [25]: print(df['Class'].value_counts())
```

DERMASON: 2546

SIRA: 2636

SEKER: 2827

HOROZ: 1928

CALI: 1550

BARBUNYA: 1322

BOMBAY: 522

Name: Class, dtype: int64

```
In [26]: df.head(6)
```

Area: 28395

Perimeter: 610.291

MajorAxisLength: 208.178117

MinorAxisLength: 173.888747

AspectRatio: 1.197191

Eccentricity: 0.549812

ConvexArea: 28715

EquivalentDiameter: 0.747933

Extent: 0.190109

Solidity: 0.763923

roundness: 0.988856

Compactness: 0.958027

ShapeFactor1: 0.91358

ShapeFactor2: 0.006979

ShapeFactor3: 0.003035

ShapeFactor4: 0.003036

Class: 0

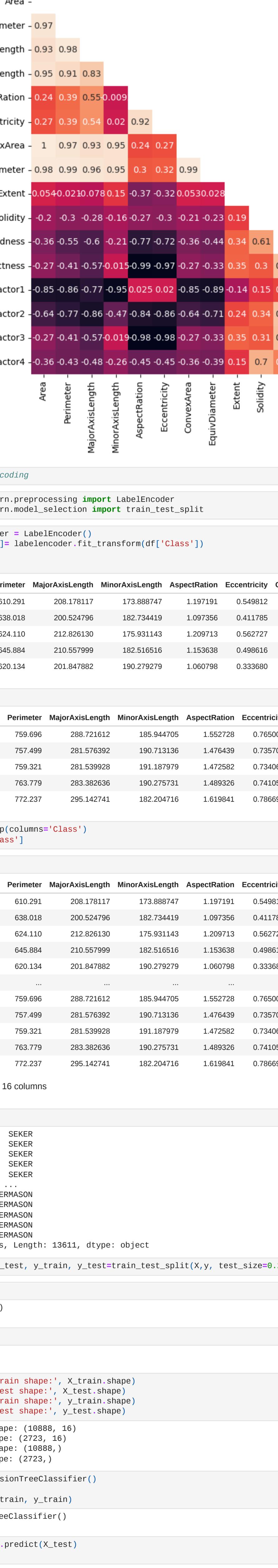
dtype: int64

```
In [27]: # visualize the distribution of the target variable
sns.countplot(x=df['Class'], data=df)
plt.title('Distribution of Bean Types')
plt.show()
```

Distribution of Bean Types


```
In [28]: # plot the class distribution
class_counts = df['Class'].value_counts()
plt.xlabel('Bean class')
plt.ylabel('Number of Instances')
plt.title('Class Distribution')
plt.show()
```

Class Distribution



```
In [29]: sns.pairplot(df)
```

```
Out[29]: 
```

pair plot



```
In [30]: df = df.drop(columns=['Class'])
y = df['Class']
plt.show()
```

X


```
In [31]: from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
```

```
Out[31]:
```