

## 1. Visualizing data ( Iris dataset – Part of NB Iris program )

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
matplotlib.interactive(True)
sns.set(color_codes = True)

#ignore warning messages
import warnings
warnings.filterwarnings('ignore')

In [2]: iris = pd.read_csv('Iris.csv')

In [3]: # Printing top 5 rows
iris.head()

Out[4]:
```

	Id	SepalLengthCm	SepalWidthCm	Petal.LengthCm	Petal.WidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa

```


In [4]: # we can see the count of each column along with their mean value, standard deviation, minium and maxiuma values.
iris.describe()

Out[4]:
```

	Id	Sepal.LengthCm	Sepal.WidthCm	Petal.LengthCm	Petal.WidthCm	Species
count	150.00000	150.00000	150.00000	150.00000	150.00000	
mean	75.000000	5.843333	3.054000	3.758333	1.198667	
std	43.443368	0.828066	0.433594	1.764400	0.761641	
min	1.000000	4.300000	2.000000	1.000000	0.100000	
25%	38.250000	5.100000	2.800000	1.600000	0.300000	
50%	75.000000	5.800000	3.000000	4.350000	1.300000	
75%	112.750000	6.400000	3.300000	5.100000	1.800000	
max	150.000000	7.900000	4.400000	6.900000	2.500000	

```


In [5]: # we can see that only one column has categorical data and all the other columns are of the numeric type with non-null entries.
iris.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   Id          150 non-null    int64
 1   Sepal.LengthCm  150 non-null    float64
 2   Sepal.WidthCm   150 non-null    float64
 3   Petal.LengthCm  150 non-null    float64
 4   Petal.WidthCm   150 non-null    float64
 5   Species       150 non-null    object
dtypes: float64(4), int64(1), object(1)
memory usage: 7.2+ KB

In [6]: # we will use the shape parameter to get the shape of the dataset.
iris.shape
#as we can see, there are 6 column and 150 rows.

Out[6]: (150, 6)

In [7]: # we will check if our data contains any missing values or not.
iris.isnull().sum()
#as u can see there is 0 null values

Out[7]:
```

	Id	Sepal.LengthCm	Sepal.WidthCm	Petal.LengthCm	Petal.WidthCm	Species
0	Id	8	8	8	8	8
1	Sepal.LengthCm	8	8	8	8	8
2	Sepal.WidthCm	8	8	8	8	8
3	Petal.LengthCm	8	8	8	8	8
4	Petal.WidthCm	8	8	8	8	8
5	Species	8	8	8	8	8
6	dtype:	int64	float64	float64	float64	object

```


In [8]: # let's see if the dataset is balanced or not i.e. all the species contain equal amounts of rows or not.
iris.value_counts('Species')

Out[8]:
```

Species	count
Iris-setosa	50
Iris-versicolor	50
Iris-virginica	50

```
dtype: int64

In [ ]:

In [ ]:
```

### data visualization

```
In [9]: # importing packages
import seaborn as sns
#seaborn is the python library that is used to create the plot and graphical representation of the data,
#internally it maps our data and create an informative plot of the data, which is easy to understand the data.
import matplotlib.pyplot as plt
#matplotlib is a low-level graph plotting library in python that serves as a visualization utility.

sns.countplot(x='Species', data=iris, )
plt.show()

#Seaborn is much more functional and organized than Matplotlib.

Out[9]:
```



the above plot is showing the counts of each species, present in dataset.

```
In [10]: #to check the no. of each species in dataset by bar graph
sns.set(font_scale = 1.5)
sns.countplot(x = 'Species', data=iris, palette='Set2')
plt.xlabel('flower species')

Out[10]:
```



```


In [11]: #pie chart
plt.figure(figsize=(10,10))
pienew = iris['Species'].value_counts()
explode = (0.05, 0)
colors = ['deeppink', 'thistle', 'red']
labels = ('setosa', 'versicolor', 'virginica')
sns.set(font_scale = 2.5)
plt.pie(pienew, autopct = '%.2f%%', colors = colors)
plt.legend(labels, loc = 'best')

Out[11]:
```



```


In [12]: iris[iris['Sepal.LengthCm'] == 0 ] # no missing values

Out[12]:
```

Id	Sepal.LengthCm	Sepal.WidthCm	Petal.LengthCm	Petal.WidthCm	Species
----	----------------	---------------	----------------	---------------	---------

```


In [13]: iris[iris['Sepal.WidthCm'] == 0 ] # no missing values

Out[13]:
```

Id	Sepal.LengthCm	Sepal.WidthCm	Petal.LengthCm	Petal.WidthCm	Species
----	----------------	---------------	----------------	---------------	---------

```


In [14]: iris[iris['Petal.LengthCm'] == 0 ] # no missing values

Out[14]:
```

Id	Sepal.LengthCm	Sepal.WidthCm	Petal.LengthCm	Petal.WidthCm	Species
----	----------------	---------------	----------------	---------------	---------

```


In [15]: iris[iris['Petal.WidthCm'] == 0 ] # no missing values

Out[15]:
```

Id	Sepal.LengthCm	Sepal.WidthCm	Petal.LengthCm	Petal.WidthCm	Species
----	----------------	---------------	----------------	---------------	---------

```


In [16]: #Visualization for understanding and analysing the distribution of data for different variables
iris.hist(figsize = (15,10),grid=True,color = 'blue')

Out[16]:
```



```


In [17]: #box plots are used to show distributions of numeric data values,
#especially when you want to compare them between multiple groups.
iris.plot(kind='box', subplots=True, layout=(3,3), sharex=False, sharey=False, figsize=(15,15))
sns.set(font_scale = 1.5)

Out[17]:
```



```


In [18]: from scipy.stats import norm
fig, ax = plt.subplots(2,2,figsize=(10,10))
sns.set(font_scale = 1)
sns.distplot(iris.Sepal.LengthCm, ax = ax[0,0], fit = norm, color = 'royalblue')
sns.distplot(iris.Sepal.WidthCm, ax = ax[0,1], color = 'black')
sns.distplot(iris.Petal.LengthCm, ax = ax[1,0], color = 'pink')
sns.distplot(iris.Petal.WidthCm, ax = ax[1,1], color = 'purple')

Out[18]:
```



```


In [19]: #bi variate analysis of data
#ascertain the correlation between the length and the width of the sepals and petals of three species of Iris flower.
corr = iris.corr()

Out[19]:
```

	Id	Sepal.LengthCm	Sepal.WidthCm	Petal.LengthCm	Petal.WidthCm
Id	1.000000	0.716676	-0.397729	0.882747	0.899759
Sepal.LengthCm	0.716676	1.000000	-0.109369	0.871754	0.817954
Sepal.WidthCm	-0.397729	-0.109369	1.000000	-0.420516	-0.356544
Petal.LengthCm	0.882747	0.871754	-0.420516	1.000000	0.982757
Petal.WidthCm	0.899759	0.817954	-0.356544	0.982757	1.000000

```


In [20]: # heat map
# heatmap is a graphical representation of data where values are depicted by color.
plt.figure(figsize=(15,10))
sns.set(font_scale = 1)
sns.heatmap(corr, annot = True, cmap = 'viridis', vmin = -1, vmax = 1, linecolor='white', linewidths=1)

Out[20]:
```



```


In [21]: #box plot of petal length
plt.figure(figsize=(15,10))
sns.set(font_scale = 1.5)
sns.boxplot(iris['Species'], iris['Petal.LengthCm'], palette='Set3')
sns.set(font_scale = 1.5)

Out[21]:
```



```


In [22]: plt.figure(figsize=(15,10))
sns.set(font_scale = 1.5)
sns.boxplot(iris['Species'], iris['Sepal.WidthCm'], palette='Set3')
sns.set(font_scale = 1.5)

Out[22]:
```



```


In [23]: plt.figure(figsize=(15,6))
sns.set(font_scale = 1.5)
sns.countplot(x = 'Petal.LengthCm', hue = 'Species', data = iris, palette = 'Set1')

Out[23]:
```



```


In [24]: sns.set(font_scale = 2)
sns.pairplot(data = iris, hue = 'Species', diag_kind = 'kde', palette = 'Set1')

Out[24]:
```



```


In [25]: # a relationship between predictor variables and a categorical Gaussian variable.
from sklearn.linear_model import LogisticRegression
#Gaussian Naive Bayes is a variant of Naive Bayes that follows Gaussian normal distribution and supports continuous data.
from sklearn.naive_bayes import GaussianNB
# for evaluating the quality of a model's predictions
# Classification report is used to measure the quality of predictions from a classification algorithm.
# as many predictions are True and the many are False.
from sklearn.metrics import classification_report
# It is mainly used for numerical and predictive analysis
from sklearn.metrics import roc_curve, auc

In [26]: #to split the dataset for training and testing
from sklearn.model_selection import train_test_split
X = iris.drop('Species', axis = 1)
y = iris['Species']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 1)

In [27]: print(X_train.shape)
print(X_test.shape)
print(y_train.size)
print(y_test.size)

Out[27]:
```

	(112, 5)	(38, 5)
0	112	38

```


In [28]: model = LogisticRegression()

In [29]: model.fit(X_train, y_train)

Out[29]: LogisticRegression()

In [30]: y_pred = model.predict(X_test)

In [31]: # as x y matrix used for evaluating the performance of a classification model
confusion = metrics.confusion_matrix(y_test, y_pred)
confusion

Out[31]:
```

	[13, 0, 0]	[ 5, 34, 2]	[ 0, 0, 9]
0	13	0	0
1	0	16	2
2	0	0	9

```
dtype: int64

In [32]: #heat map of the above confusion matrix
plt.figure(figsize=(15,6))
sns.heatmap(confusion, annot=True, linecolor='white', linewidths=1)

Out[32]:
```



```


In [33]: print('Accuracy of Logistic Regression is: ', model.score(X_test,y_test) * 100, '%')

Accuracy of Logistic Regression is: 94.73684210526315 %

In [34]: print(classification_report(y_test,y_pred))

precision    recall    f1-score   support

Iris-setosa      1.00      1.00      1.00       13
Iris-versicolor  0.88      0.88      0.88       16
Iris-virginica   0.82      0.90      0.86        9

accuracy         0.94
macro avg        0.94      0.96      0.94       38
weighted avg     0.96      0.95      0.95       38

In [35]: #f1 score
metrics.f1_score(y_test, y_pred, average = 'weighted')

Out[35]: 0.9482456140358877

In [36]: # Naive Bayes Classifier
nbModel = GaussianNB()

In [37]: nbModel = GaussianNB()

In [38]: #fitting the model
nbModel.fit(X_train, y_train)

Out[38]: GaussianNB()

In [39]: #predicting the x-test
nb_y_pred = nbModel.predict(X_test)

In [40]: #matrix of y-test and y-pred
nbConfusion = metrics.confusion_matrix(y_test, nb_y_pred)
nbConfusion

Out[40]:
```

	[13, 0, 0]	[ 5, 34, 2]	[ 0, 0, 9]
0	13	0	0
1	0	16	2
2	0	0	9

```
dtype: int64

In [41]: #heat map of the above confusion matrix
plt.figure(figsize=(15,6))
sns.heatmap(nbConfusion, annot=True, linecolor='white', linewidths=1)

Out[41]:
```



```


In [42]: #accuracy
print('Accuracy of Naive Bayes Classifier is: ', nbModel.score(X_test,y_test)*100, '%')

Accuracy of Naive Bayes Classifier is: 100.0 %

In [ ]:

In [ ]:
```