

## Data Analytics III

1. Implement Simple Naïve Bayes classification algorithm using Python/R on iris.csv dataset.
2. Compute Confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset.

Step 1: Import libraries and create alias for Pandas, Numpy and Matplotlib

Step 2: Import the Iris dataset by calling URL.

Step 3: Initialize the data frame

Step 4: Perform Data Preprocessing

- Convert Categorical to Numerical Values if applicable
- Check for Null Value
- Divide the dataset into Independent(X) and Dependent(Y) variables.
- Split the dataset into training and testing datasets
- Scale the Features if necessary.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
df=pd.read_csv("Iris.csv")
df
```

Saving... X

			WidthCm	PetalLengthCm	PetalWidthCm	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
...	...	...	...	...	...	...
145	146	6.7	3.0	5.2	2.3	Iris-virginica
146	147	6.3	2.5	5.0	1.9	Iris-virginica
147	148	6.5	3.0	5.2	2.0	Iris-virginica
148	149	6.2	3.4	5.4	2.3	Iris-virginica
149	150	5.9	3.0	5.1	1.8	Iris-virginica

150 rows × 6 columns

```
df.isnull().sum()
```

```
Id          0
SepalLengthCm  0
SepalWidthCm  0
PetalLengthCm  0
PetalWidthCm  0
Species      0
dtype: int64
```

```
#Removing null values
columns=['SepalLengthCm','SepalWidthCm','PetalLengthCm','PetalWidthCm']
for col in columns:
    df[col]=df[col].fillna(df[col].mean())
df
```

	<b>Id</b>	<b>SepalLengthCm</b>	<b>SepalWidthCm</b>	<b>PetalLengthCm</b>	<b>PetalWidthCm</b>	<b>Species</b>
<b>0</b>	1	5.1	3.5	1.4	0.2	Iris-setosa
<b>1</b>	2	4.9	3.0	1.4	0.2	Iris-setosa
<b>2</b>	3	4.7	3.2	1.3	0.2	Iris-setosa
<b>3</b>	4	4.6	3.1	1.5	0.2	Iris-setosa
<b>4</b>	5	5.0	3.6	1.4	0.2	Iris-setosa
...	...	...	...	...	...	...
<b>145</b>	146	6.7	3.0	5.2	2.3	Iris-virginica
<b>146</b>	147	6.3	2.5	5.0	1.9	Iris-virginica
<b>147</b>	148	6.5	3.0	5.2	2.0	Iris-virginica
<b>148</b>	149	6.2	3.4	5.4	2.3	Iris-virginica
<b>149</b>	150	5.9	3.0	5.1	1.8	Iris-virginica

150 rows × 6 columns

```
#Converting categorical values ot numeric values
from sklearn import preprocessing
label_encoder = preprocessing.LabelEncoder()
df['Species']= label_encoder.fit_transform(df['Species'])
df
```

	<b>Id</b>	<b>SepalLengthCm</b>	<b>SepalWidthCm</b>	<b>PetalLengthCm</b>	<b>PetalWidthCm</b>	<b>Species</b>
<b>0</b>	1	5.1	3.5	1.4	0.2	0
	2	4.9	3.0	1.4	0.2	0
	3	4.7	3.2	1.3	0.2	0
<b>3</b>	4	4.6	3.1	1.5	0.2	0
<b>4</b>	5	5.0	3.6	1.4	0.2	0
...	...	...	...	...	...	...
<b>145</b>	146	6.7	3.0	5.2	2.3	2
<b>146</b>	147	6.3	2.5	5.0	1.9	2
<b>147</b>	148	6.5	3.0	5.2	2.0	2
<b>148</b>	149	6.2	3.4	5.4	2.3	2
<b>149</b>	150	5.9	3.0	5.1	1.8	2

150 rows × 6 columns

```
y=df['Species']
x=df.drop('Species',axis=1)
```

```
#Spliting data for training and testing
#Here, 20% data used for testing and 80% data used for training
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size = 0.2,random_state=2)
```

```
# import the class
from sklearn.naive_bayes import GaussianNB
gaussian = GaussianNB()
gaussian.fit(x_train, y_train)
```

```
▼ GaussianNB
GaussianNB()
```

```
y_pred = gaussian.predict(x_test)
```

```
from sklearn.metrics import accuracy_score, precision_score, recall_score
accuracy = accuracy_score(y_test, y_pred)
accuracy
```

```
1.0
```

```
precision = precision_score(y_test, y_pred, average='micro')
precision
```

```
1.0
```

```
recall = recall_score(y_test, y_pred, average='micro')
recall
```

```
1.0
```

```
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```

```
array([[14,  0,  0],
       [ 0,  8,  0],
       [ 0,  0,  8]])
```

Saving...

