

Breast Cancer Prediction

Importing the necessary libraries

```
In [1]: 1 # Data Analysis Libraries
2
3 import pandas as pd
4
5 # Visualisation Libraries
6 import matplotlib.pyplot as plt
7 import plotly.express as px
8 import plotly.graph_objs as go
9 import plotly.io as pio
10
11 # Library for Splitting Dataset
12 from sklearn.model_selection import train_test_split
13
14 #Library for Implementing models
15 from sklearn.linear_model import LogisticRegression
16 from sklearn.naive_bayes import GaussianNB
17 from sklearn.ensemble import RandomForestClassifier
18 from sklearn.tree import DecisionTreeClassifier
19
20 # Library for metrics generation
21 from sklearn.metrics import accuracy_score
```

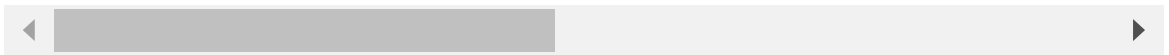
Importing the Dataset and Visualising it

```
In [2]: 1 cancer = pd.read_csv(r'C:\Users\tanmayee\OneDrive\Documents\Personal\0  
2 cancer
```

Out[2]:

	patient_id	clump_thickness	cell_size_uniformity	cell_shape_uniformity	marginal_adhes
0	1000025	5.0	1.0	1	
1	1002945	5.0	4.0	4	
2	1015425	3.0	1.0	1	
3	1016277	6.0	8.0	8	
4	1017023	4.0	1.0	1	
...
694	776715	3.0	1.0	1	
695	841769	2.0	1.0	1	
696	888820	5.0	10.0	10	
697	897471	4.0	8.0	6	
698	897471	4.0	8.0	8	

699 rows × 12 columns



Deleting the columns from the dataset that are no more required

```
In [3]: 1 cancer = cancer.drop(columns=['patient_id', 'doctor_name'], axis=1)
        2 cancer
```

```
Out[3]:
```

	clump_thickness	cell_size_uniformity	cell_shape_uniformity	marginal_adhesion	single_
0	5.0	1.0	1	1	
1	5.0	4.0	4	5	
2	3.0	1.0	1	1	
3	6.0	8.0	8	1	
4	4.0	1.0	1	3	
...
694	3.0	1.0	1	1	
695	2.0	1.0	1	1	
696	5.0	10.0	10	3	
697	4.0	8.0	6	4	
698	4.0	8.0	8	5	

699 rows × 10 columns



The value counts() method is used to first calculate the counts of each distinct value in the 'cell size uniformity' column. Using the idxmax() function, it then returns the index (i.e., the value) corresponding to the highest count. This index corresponds to the first value in the array returned by the mode() function, which is the mode of the 'cell size uniformity' column.

```
In [4]: 1 r1 = cancer['cell_size_uniformity'].value_counts().idxmax()
        2 r1
```

```
Out[4]: 1.0
```

The 'cell size uniformity' column's mode is first calculated using the mode() function, and then the iloc method is used to choose the first value in the resultant array. The fillna() function, with the inplace argument set to True, utilises this mode value to replace missing values in the 'cell size uniformity' column. This is comparable to the supplied code snippet, except the mode value is computed inline in the fillna() function call.

```
In [5]: 1 cancer['cell_size_uniformity'].fillna(cancer['cell_size_uniformity'].m
```

```
In [6]: 1 r2 = cancer.isna().sum()
        2 r2
```

```
Out[6]: clump_thickness      1
        cell_size_uniformity  0
        cell_shape_uniformity 0
        marginal_adhesion     0
        single_ep_cell_size    0
        bare_nuclei            2
        bland_chromatin        4
        normal_nucleoli        1
        mitoses                0
        class                  0
        dtype: int64
```

Fills missing values in the 'clump_thickness', 'bare_nuclei', 'bland_chromatin', and 'normal_nucleoli' columns of the 'cancer' dataframe with the mode of each respective column.

```
In [7]: 1 cols = ['clump_thickness', 'bare_nuclei', 'bland_chromatin', 'normal_n
        2 cancer[cols] = cancer[cols].apply(lambda x: x.fillna(x.mode()[0]))
```

```
In [8]: 1 r3 = cancer.isna().sum()
        2 r3
```

```
Out[8]: clump_thickness      0
        cell_size_uniformity  0
        cell_shape_uniformity 0
        marginal_adhesion     0
        single_ep_cell_size    0
        bare_nuclei            0
        bland_chromatin        0
        normal_nucleoli        0
        mitoses                0
        class                  0
        dtype: int64
```

Counting total number of patients registered in a hospital

```
In [9]: 1 total_patients = len(cancer)
        2 total_patients
```

```
Out[9]: 699
```

Seperating the normal patients and Patients suffering from Cancer

```
In [10]: 1 normal_patients = cancer["class"].value_counts()[0]
        2 normal_patients
```

```
Out[10]: 458
```

```
In [11]: 1 breast_cancer_patients = cancer["class"].value_counts()[1]
         2 breast_cancer_patients
```

Out[11]: 241

Using the formula to find total number of people suffering from cancer

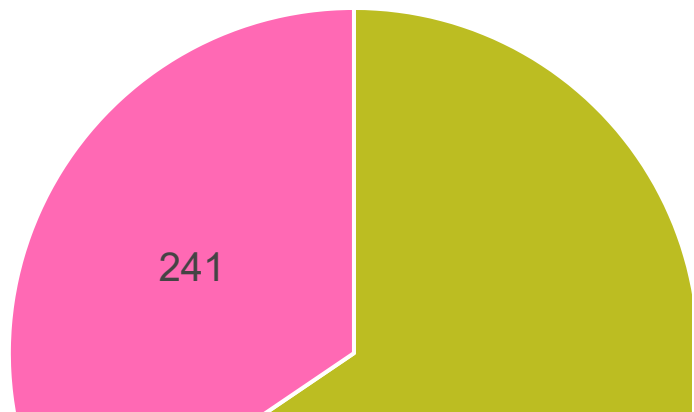
```
In [12]: 1 percentage_breast_cancer = round((breast_cancer_patients / total_patients) * 100)
         2 print("Percentage of people with breast cancer: ", percentage_breast_cancer)
```

Percentage of people with breast cancer: 34.48 %

Plotting the Pie Chart to show the number of people suffering from breast cancer and the normal people

```
In [13]: 1 labels = ['Without breast cancer', 'With breast cancer']
2 sizes = [normal_patients, breast_cancer_patients]
3
4 fig1 = go.Figure(data=[go.Pie(labels=labels, values=sizes)])
5
6 fig1.update_traces(hoverinfo='label+percent', textinfo='value', textfont=
7                     marker=dict(colors=['#bcbd22', '#FF69B4'], line=dict
8
9 fig1.update_layout(title='Breast Cancer Cases', title_font_size=30, ti
10                    legend=dict(title='Cases', font=dict(size=20)),
11                    font=dict(family='Arial', size=20))
12
13 fig1.show()
```

Breast Cancer Cases



The values "benign" and "malignant" are mapped to 0 and 1, respectively, in a dictionary class mapping that we construct. The values in the "class" column of the "cancer" dataframe are then replaced with the equivalent values from the class mapping dictionary using the map() function.

```
In [14]: 1 class_mapping = {'benign': 0, 'malignant': 1}
2 cancer['class'] = cancer['class'].map(class_mapping)
```

Code replaces insufficient data in the 'bare nuclei' column of the 'cancer' dataframe with the mode of the column using the fillna() method.

```
In [15]: 1 mode = cancer['bare_nuclei'].mode()[0]
2 cancer['bare_nuclei'] = cancer['bare_nuclei'].replace("?", mode)
3 print(cancer)
```

	clump_thickness	cell_size_uniformity	cell_shape_uniformity	\
0	5.0	1.0	1	
1	5.0	4.0	4	
2	3.0	1.0	1	
3	6.0	8.0	8	
4	4.0	1.0	1	
..	
694	3.0	1.0	1	
695	2.0	1.0	1	
696	5.0	10.0	10	
697	4.0	8.0	6	
698	4.0	8.0	8	

	marginal_adhesion	single_ep_cell_size	bare_nuclei	bland_chromatin
\				
0	1	2	1	3.0
1	5	7	10	3.0
2	1	2	2	3.0
3	1	3	4	3.0
4	3	2	1	3.0
..
694	1	3	2	1.0
695	1	2	1	1.0
696	3	7	3	8.0
697	4	3	4	10.0
698	5	4	5	10.0

	normal_nucleoli	mitoses	class
0	1.0	1	0
1	2.0	1	0
2	1.0	1	0
3	7.0	1	0
4	1.0	1	0
..
694	1.0	1	0
695	1.0	1	0
696	10.0	2	1
697	6.0	1	1
698	4.0	1	1

[699 rows x 10 columns]

This line of code converts the 'bare_nuclei' column of the 'cancer' dataframe from a string data type to an integer data type using the astype() method.

Creating a new dataframe that contains first 9 columns to divide it into train and test

```
In [16]: 1 cr_x = cancer.loc[:, cancer.columns[:9]]
          2 cr_x
```

```
Out[16]:
```

	clump_thickness	cell_size_uniformity	cell_shape_uniformity	marginal_adhesion	single_
0	5.0	1.0	1	1	
1	5.0	4.0	4	5	
2	3.0	1.0	1	1	
3	6.0	8.0	8	1	
4	4.0	1.0	1	3	
...
694	3.0	1.0	1	1	
695	2.0	1.0	1	1	
696	5.0	10.0	10	3	
697	4.0	8.0	6	4	
698	4.0	8.0	8	5	

699 rows × 9 columns



Creating a new dataframe that contains last column of dataset to divide it into train and test

```
In [17]: 1 cr_y = cancer['class']
          2 cr_y
```

```
Out[17]: 0      0
          1      0
          2      0
          3      0
          4      0
          ..
          694    0
          695    0
          696    1
          697    1
          698    1
          Name: class, Length: 699, dtype: int64
```

Dividing the dataset into 80:20 ratio for Train and Test respectively

```
In [18]: 1 split = train_test_split(cr_x, cr_y, test_size=0.2, random_state=1)
          2 crx_train, crx_test, cry_train, cry_test = split
```



```
In [19]: 1 crx_train
```

Out[19]:

	clump_thickness	cell_size_uniformity	cell_shape_uniformity	marginal_adhesion	single_
617	1.0	1.0	1	1	
107	1.0	6.0	8	10	
17	4.0	1.0	1	1	
441	5.0	2.0	2	4	
365	2.0	1.0	1	1	
...
144	2.0	1.0	1	1	
645	3.0	1.0	1	1	
72	1.0	3.0	3	2	
235	3.0	1.0	4	1	
37	6.0	2.0	1	1	

559 rows × 9 columns



Logistic Regression

```
In [20]: 1 logi_reg = LogisticRegression()  
2 logi_reg.fit(crx_train, cry_train)  
3 LR_p = logi_reg.predict(crx_test)
```

```
In [21]: 1 LRaccuracy = accuracy_score(cry_test, LR_p)  
2 print("Accuracy on Logistic Regression:", LRaccuracy)
```

Accuracy on Logistic Regression: 0.9642857142857143

Naive Bayes

```
In [22]: 1 nbx = GaussianNB()  
2 nbx.fit(crx_train, cry_train)  
3 NB_predx = nbx.predict(crx_test)
```

```
In [23]: 1 NBaccuracy = accuracy_score(cry_test, NB_predx)  
2 print("Accuracy of Naive Bayes:", NBaccuracy)
```

Accuracy of Naive Bayes: 0.9642857142857143

Random Forest

```
In [24]: 1 rfx = RandomForestClassifier()
2 rfx.fit(crx_train, cry_train)
3 RF_predx = rfx.predict(crx_test)
```

```
In [25]: 1 RFaccuracy = accuracy_score(cry_test, RF_predx)
2 print("Accuracy of Random Forest", RFaccuracy)
```

Accuracy of Random Forest 0.9642857142857143

Decision Tree

```
In [26]: 1 dtx = DecisionTreeClassifier()
2 dtx.fit(crx_train, cry_train)
3 DT_predx = dtx.predict(crx_test)
```

```
In [27]: 1 DTaccuracy = accuracy_score(cry_test, DT_predx)
2 print("Accuracy of Decision Tree", DTaccuracy)
```

Accuracy of Decision Tree 0.8928571428571429

Comparison of Accuracy of Implemented Algorithms

```
In [28]: 1 print("Logistic Regression Accuracy:", LRaccuracy)
2 print("Naive Bayes Accuracy:", NBaccuracy)
3 print("Random Forest Accuracy:", RFaccuracy)
4 print("Decision Tree Accuracy:", DTaccuracy)
```

Logistic Regression Accuracy: 0.9642857142857143

Naive Bayes Accuracy: 0.9642857142857143

Random Forest Accuracy: 0.9642857142857143

Decision Tree Accuracy: 0.8928571428571429

Code to convert the accuracy into dataframe

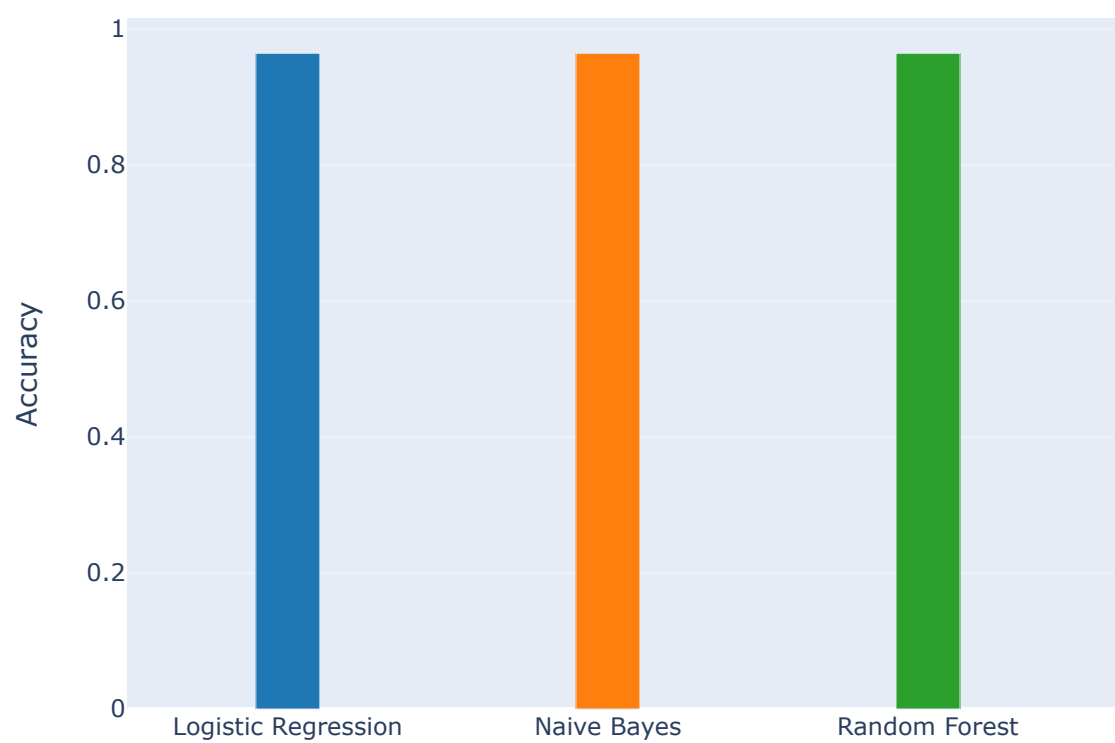
```
In [29]: 1 data = {
2     'Model': ['Logistic Regression', 'Naive Bayes', 'Random Forest', '
3     'Accuracy': [LRaccuracy, NBaccuracy, RFaccuracy, DTaccuracy]
4 }
5
6 # create a pandas DataFrame from the dictionary
7 df = pd.DataFrame(data)
8
9 # print the DataFrame
10 print(df)
```

	Model	Accuracy
0	Logistic Regression	0.964286
1	Naive Bayes	0.964286
2	Random Forest	0.964286
3	Decision Tree	0.892857

Plot of Accuracy of Implemented Algorithms

```
In [30]: 1 # set the width and color of the bars
2 bar_width = 0.2
3 # define a list of colors for the bars
4 colors = ['#1f77b4', '#ff7f0e', '#2ca02c', '#d62728']
5
6 # create the trace
7 trace = go.Bar(
8     x=df['Model'],
9     y=df['Accuracy'],
10    width=bar_width,
11    marker=dict(
12        color=colors
13    )
14 )
15
16 # create the layout
17 layout = go.Layout(
18     title='Model Comparison',
19     yaxis=dict(title='Accuracy')
20 )
21
22 # create the figure
23 fig = go.Figure(data=[trace], layout=layout)
24
25 # adjust the plot width
26 fig.update_layout(width=800)
27
28 fig.show()
```

Model Comparison



◀

▶

In []:

1